

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6260384号
(P6260384)

(45) 発行日 平成30年1月17日 (2018. 1. 17)

(24) 登録日 平成29年12月22日 (2017. 12. 22)

(51) Int. Cl.

F I

G 0 6 F 3 / 0 6 (2006.01)

G 0 6 F 3 / 0 6 3 0 1 N

G 0 6 F 3 / 0 6 3 0 4 N

請求項の数 8 (全 33 頁)

(21) 出願番号 特願2014-56799 (P2014-56799)
 (22) 出願日 平成26年3月19日 (2014. 3. 19)
 (65) 公開番号 特開2015-179425 (P2015-179425A)
 (43) 公開日 平成27年10月8日 (2015. 10. 8)
 審査請求日 平成28年12月6日 (2016. 12. 6)

(73) 特許権者 000005223
 富士通株式会社
 神奈川県川崎市中原区上小田中4丁目1番
 1号
 (74) 代理人 100092978
 弁理士 真田 有
 (74) 代理人 100112678
 弁理士 山本 雅久
 (72) 発明者 大江 和一
 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
 (72) 発明者 岩田 聡
 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

最終頁に続く

(54) 【発明の名称】 ストレージ制御装置、制御プログラム、及び制御方法

(57) 【特許請求の範囲】

【請求項 1】

第1の記憶装置の記憶領域を所定の大きさで分割した複数の単位領域について、入力された要求に対する応答性能を監視する監視部と、

前記第1の記憶装置の移動対象の単位領域に記憶されたデータを前記第1の記憶装置とは異なる性能の第2の記憶装置に移動する移動処理において、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第2の記憶装置に移動する分割部と、

前記監視部が監視した前記移動処理実行中の第1の応答性能に基づいて、前記所定の分割数を変更する変更部と、
 をそなえることを特徴とする、ストレージ制御装置。

【請求項 2】

前記変更部は、前記第1の応答性能と、前記移動処理実行前の第2の応答性能とに基づいて、前記所定の分割数を変更することを特徴とする、請求項1記載のストレージ制御装置。

【請求項 3】

前記所定の分割数の初期値は、前記移動処理実行前の第2の応答性能を基に求められることを特徴とする、請求項2記載のストレージ制御装置。

【請求項 4】

前記変更部は、前記第1の応答性能と前記移動処理実行前の第2の応答性能とを比較し

、前記第 1 の応答性能が前記第 2 の応答性能よりも劣化していると判断した場合に前記所定の分割数を増加させる一方、前記第 1 の応答性能が前記第 2 の応答性能よりも優れていると判断した場合に前記所定の分割数を減少させる、
ことを特徴とする、請求項 1 ～ 3 のいずれか 1 項記載のストレージ制御装置。

【請求項 5】

前記変更部は、前記監視部が前記移動処理実行中の複数の時点で監視した複数の応答性能を取得し、前記第 2 の応答性能と前記複数の応答性能の平均値とを比較して、前記平均値が前記第 2 の応答性能よりも劣化しているか否かに応じて前記所定の分割数を増減させる、
ことを特徴とする、請求項 4 記載のストレージ制御装置。

10

【請求項 6】

前記監視部は、前記複数の単位領域について単位領域ごとに入出力数を集計し、
前記ストレージ制御装置は、
前記監視部により集計された入出力数が第 1 の閾値より大きな単位領域と所定の距離内にある単位領域を繋ぎ合わせた拡張領域と、該拡張領域と繋がる他の拡張領域と、を合わせた移動領域を特定する特定部、をさらにそなえ、
前記分割部は、前記移動領域に記憶されたデータを前記第 2 の記憶装置に移動する移動処理において、前記移動領域に含まれる複数の単位領域の各々を所定の分割数により複数の分割領域に分割し、前記移動領域に記憶されたデータを前記分割領域の単位で前記第 2 の記憶装置に移動する、
ことを特徴とする、請求項 1 ～ 5 のいずれか 1 項記載のストレージ制御装置。

20

【請求項 7】

第 1 の記憶装置及び第 2 の記憶装置の制御を行なうコンピュータに、
前記第 1 の記憶装置の記憶領域を所定の大きさを分割した複数の単位領域について、入力された要求に対する応答性能を監視し、
前記第 1 の記憶装置の移動対象の単位領域に記憶されたデータを前記第 1 の記憶装置とは異なる性能の第 2 の記憶装置に移動する移動処理を行ない、
前記移動処理において、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第 2 の記憶装置に移動し、
前記監視により監視した前記移動処理実行中の第 1 の応答性能に基づいて、前記所定の分割数を変更する、
処理を実行させることを特徴とする、制御プログラム。

30

【請求項 8】

第 1 の記憶装置及び第 2 の記憶装置の制御を行なうストレージ制御装置における制御方法であって、
前記第 1 の記憶装置の記憶領域を所定の大きさを分割した複数の単位領域について、入力された要求に対する応答性能を監視し、
前記第 1 の記憶装置の移動対象の単位領域に記憶されたデータを前記第 1 の記憶装置とは異なる性能の第 2 の記憶装置に移動する移動処理を行ない、
前記移動処理において、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第 2 の記憶装置に移動し、
前記監視により監視した前記移動処理実行中の第 1 の応答性能に基づいて、前記所定の分割数を変更する、
ことを特徴とする、制御方法。

40

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ストレージ制御装置、制御プログラム、及び制御方法に関する。

【背景技術】

【0002】

50

データを格納するストレージシステムとして、複数の記憶媒体（記憶装置）を組み合わせた階層ストレージシステムが用いられることがある。階層ストレージシステムは、例えば、高速アクセスが可能であるが比較的低容量、高価格なＳＳＤ（Solid State Drive）と、大容量で低価格であるが比較的低速なＨＤＤ（Hard Disk Drive）とを含む。

階層ストレージシステムでは、アクセス頻度が低い領域をＨＤＤに配置する一方、アクセス頻度が高い領域をＳＳＤに配置することで、ＳＳＤの使用効率を高め、システム全体の性能を高めることができる。つまり、階層ストレージシステムの性能を向上させるには、アクセス頻度が高い領域を効率的にＳＳＤに配置することが望ましい。

【０００３】

アクセス頻度が高い領域をＳＳＤに配置する手法としては、例えば、前日のアクセス頻度に応じて、１日単位でアクセス頻度が高い領域をＳＳＤに配置する手法がある。具体的には、階層ストレージシステムは、ユーザのアクセス頻度が低い深夜時間帯に２４時間分のアクセス頻度を集計し、アクセス頻度が高い領域から順にＳＳＤへ配置する。毎日ほぼ同じ領域へアクセス集中が発生するワークロードにおいてはこの手法で十分である。

【０００４】

しかし、数分から数十分程度の比較的短時間にアクセス集中（負荷）が移動するワークロードにおいては、１日単位のアクセス頻度の集計では追従できない場合が多い。なお、ワークロードとは、記憶装置へのアクセス分布をいい、時間の経過と記憶装置のオフセット位置（領域）とに応じて変化する。短時間に負荷が移動するワークロードに対応するためには、アクセス頻度が高くなった領域をリアルタイムで把握して、当該領域をＳＳＤに移動することが好ましい。

【０００５】

また、ＨＤＤからＳＳＤへデータを移動する際にも、当該データに対するユーザからのＩＯ（Input Output）（以下、ユーザＩＯという）が発生し得る。このユーザＩＯへの対処としては、例えば、ストレージシステムが、第１のボリュームから第２のボリュームへのデータの移行中に、共有メモリ等に対象データを転送してアクセス応答可能な状態にする技術が知られている（例えば、特許文献１参照）。これにより、移行中における対象データに対するアクセス性能が確保される。

【０００６】

また、記憶制御装置が、データ処理装置から論理的記憶装置にアクセスがあったとき、アクセス位置が再配置完了領域か再配置未完了領域かに応じて、再配置先又は当該論理的記憶装置にアクセスさせる技術も知られている（例えば、特許文献２参照）。

さらに、ストレージ管理装置が、アクセス対象範囲の論理セグメント（segment）と物理セグメントとを、それぞれサブ論理セグメントとサブ物理セグメントとに分割する技術も知られている（例えば、特許文献３参照）。この技術では、ストレージ管理装置は、サブセグメント（sub-segment）単位で対象データを再配置することで、ストレージ装置にかかる負荷を分散し、アクセス性能を向上させることができる。

【先行技術文献】

【特許文献】

【０００７】

【特許文献１】特開２００８－２９９５５９号公報

【特許文献２】特開２００３－２７１４２５号公報

【特許文献３】国際公開第２００８／１２６２０２号パンフレット

【発明の概要】

【発明が解決しようとする課題】

【０００８】

データの移行中に、対象データを一時バッファに転送してユーザＩＯを処理する手法では、移動完了後に一時バッファと移動先のＳＳＤとの間でsync（同期）等の処理が発生するため、移動時間が伸びてしまう。

最も短時間にデータを移動する手法は、データの移動中は当該データへのユーザＩＯを

10

20

30

40

50

ブロックすることが考えられるが、ユーザＩＯがブロックされることで、ユーザＩＯのレスポンスは悪化する。一方、階層ストレージシステムが一度に移動する領域を小さくすれば、ユーザＩＯのレスポンス悪化を小さくできるが、代わりに移動時間が大きくなる。

【０００９】

１つの側面では、本発明は、第１の記憶装置から第２の記憶装置へのデータ移動にかかる処理時間を低減させつつ、入力される要求への応答性能の劣化を抑制することを目的とする。

なお、前記目的に限らず、後述する発明を実施するための形態に示す各構成により導かれる作用効果であって、従来の技術によっては得られない作用効果を奏することも本発明の他の目的の１つとして位置付けることができる。

【課題を解決するための手段】

【００１０】

本件のストレージ制御装置は、第１の記憶装置の記憶領域を所定の大きさを分割した複数の単位領域について、入力された要求に対する応答性能を監視する監視部をそなえる。また、このストレージ制御装置は、前記第１の記憶装置の移動対象の単位領域に記憶されたデータを前記第１の記憶装置とは異なる性能の第２の記憶装置に移動する移動処理において、以下の処理を行なう分割部をそなえる。ここで、分割部は、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第２の記憶装置に移動する。さらに、このストレージ制御装置は、前記監視部が監視した前記移動処理実行中の第１の応答性能に基づいて、前記所定の分割数を変更する変更部をそなえる。

【発明の効果】

【００１１】

一実施形態によれば、第１の記憶装置から第２の記憶装置へのデータ移動にかかる処理時間を低減させつつ、入力される要求への応答性能の劣化を抑制することができる。

【図面の簡単な説明】

【００１２】

【図１】階層移動中に移動領域に発生したユーザＩＯをブロックする手法の一例を示す図である。

【図２】移動領域をサブセグメント単位で移動する例を示す図である。

【図３】サブセグメントの分割数と領域移動時間との関係の一例を示す図である。

【図４】サブセグメントの分割数とＩＯ応答時間との関係の一例を示す図である。

【図５】サブセグメントの分割数を２５６にした場合のＩＯ応答時間の一例を示す図である。

【図６】サブセグメントの分割数を２０４８にした場合のＩＯ応答時間の一例を示す図である。

【図７】一実施形態に係る階層ストレージシステムの構成例を示す図である。

【図８】図７に示すデータベースの一例を示す図である。

【図９】図７に示す階層テーブルの一例を示す図である。

【図１０】データ収集部によるデータ収集処理の動作例を示すフローチャートである。

【図１１】ワークロード分析部による移動判定処理の動作例を示すフローチャートである。

【図１２】移動指示部による移動指示通知処理の動作例を示すフローチャートである。

【図１３】分割数判定部による分割数判定処理の動作例を示すフローチャートである。

【図１４】階層ドライバによる転送指示通知処理の動作例を示すフローチャートである。

【図１５】階層ドライバによる転送完了受信処理の動作例を示すフローチャートである。

【図１６】分割部による転送指示受信処理の動作例を示すフローチャートである。

【図１７】分割部による分割数更新処理の動作例を示すフローチャートである。

【図１８】ＩＯマップ部によるＩＯ受信処理の動作例を示すフローチャートである。

【図１９】図７に示す階層ストレージ制御装置のハードウェア構成例を示す図である。

10

20

30

40

50

【図 20】適用例に係る階層ストレージ制御装置による動的階層制御を説明するための図である。

【図 21】適用例に係る階層ストレージ制御装置による動的階層制御を説明するための図である。

【図 22】適用例に係る階層ストレージシステムの構成例を示す図である。

【図 23】適用例に係るデータ収集部によるデータ収集処理の動作例を示すフローチャートである。

【図 24】図 22 に示すデータベースの一例を示す図である。

【図 25】適用例に係るワークロード分析部による移動判定処理の動作例を示すフローチャートである。

【図 26】図 22 に示す候補テーブルの一例を示す図である。

【図 27】図 22 に示す管理テーブルの一例を示す図である。

【図 28】適用例に係る移動指示部による移動指示通知処理の動作例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】

【0013】

以下、図面を参照して実施の形態を説明する。

〔1〕一実施形態

〔1-1〕対比例

はじめに、図 1 及び図 2 に示す対比例を説明する。図 1 は、階層移動中に移動領域に発生したユーザ I/O をブロックする手法の一例を示す図であり、図 2 は、移動領域をサブセグメント単位で移動する例を示す図である。なお、図 1 及び図 2 では、階層ストレージ制御装置 100 が Linux（登録商標） device-mapper の機能を用いるものとする。この例では、device-mapper がストレージボリュームをセグメント単位で監視し、高負荷となったセグメントのデータを HDD 300 から SSD 200 へ移動することで高負荷セグメントへの I/O を処理する。

【0014】

まず、図 1 において、階層ストレージ制御装置 100 のユーザ空間で実行されるアプリケーションがデータ格納先の変更要求としてコピー指示を発行する（図 1 の（1）参照）。OS（Operating System）空間で実行される階層ドライバ 110 は、コピー指示を受け取り、格納先を変更するため、デバイス間のデータコピーを非同期で実行する kcopyd に、SSD 200 - HDD 300 間のコピー（移動）を指示する。kcopyd による移動中、ユーザから I/O 要求が発行されると（図 1 の（2）参照）、階層ドライバ 110 は、メモリ等のペンディングキューに I/O 要求を格納し、移動完了まで待ち合わせを行なう（図 1 の（3）参照）。なお、device-mapper および kcopyd はコンピュータプログラムとして実装されている。

【0015】

移動が完了すると（図 1 の（4）参照）、階層ドライバ 110 は、移動先の SSD 200 又は HDD 300 を選択し、SSD ドライバ 120 又は HDD ドライバ 130 を介してペンディングキューで保留していた I/O 要求を発行する（図 1 の（5）参照）。そして、I/O 要求を受けた移動先の SSD 200 又は HDD 300 は、ユーザへ I/O 応答を返す（図 1 の（6）参照）。

【0016】

図 1 に示す例では、I/O 要求がペンディングキューに保留されていた時間は、そのままユーザにレスポンスの悪化として見えてしまう。例えば、セグメント = 1 GB（Byte），HDD 300 のスループット性能 = 100 MB / sec，SSD 200 のスループット性能 = 1000 MB / sec であると仮定すると、移動時間は、1 [GB] / 100 [MB / sec] = 10 秒となる。すなわち、ユーザ I/O が最大 10 秒待たされる可能性があることが分かる。一時的であっても、階層ストレージシステムにおいて、ユーザ I/O が最大 10 秒待たされることは許容されない場合が多い。

10

20

30

40

50

【 0 0 1 7 】

一方、図 2 では、セグメントをさらにサブセグメントというより小さい単位に分割し、階層ストレージ制御装置 1 0 0 がサブセグメント単位で領域移動を行なう。これにより、ユーザ I O の待ち時間を、セグメント全体の移動時間よりも小さいサブセグメントの移動時間に抑えることができる。

しかし、図 2 に示すように、サブセグメント単位で領域移動を行なう場合、セグメント単位での領域移動よりも全体での移動コストがかかってしまう。これは、移動時間以外のオーバーヘッドがサブセグメントの分割数分だけ増えるためである。

【 0 0 1 8 】

そこで、セグメント移動時のユーザ I O 保留時間をできるだけ短くしつつ、可能な限り短時間にセグメント移動を完了させることが望ましい。上述のように、セグメントをサブセグメントに分割し、さらにその分割数を増やすことでユーザ I O の保留時間を小さくすることは可能である。しかし、セグメント全体の移動時間は大きく増加してしまう。

図 3 ~ 図 6 は、サブセグメントの分割数を増やした場合に移動時間がどの程度増加するかを評価した評価結果である。図 3 はサブセグメントの分割数と領域移動時間との関係の一例を示す図であり、図 4 はサブセグメントの分割数と I O 応答時間との関係の一例を示す図である。図 5 及び図 6 は、それぞれ、サブセグメントの分割数を 2 5 6 及び 2 0 4 8 にした場合の I O 応答時間の一例を示す図である。なお、図 3 ~ 図 6 は、セグメント = 1 G B , H D D 3 0 0 のスループット性能 = 1 0 0 M B / s e c , S S D 2 0 0 のスループット性能 = 1 0 0 0 M B / s e c の環境での評価結果である。

【 0 0 1 9 】

例えば、図 3 に示すように、S S D 2 0 0 から H D D 3 0 0 への領域移動は、セグメントが未分割の場合、1 0 秒程度で完了する。しかし、分割数を増やしていくことで、領域移動時間は最大 1 2 倍 (1 2 0 秒) 程度にまで増加することが分かる。

一方、ユーザ I O への応答時間に関しては、図 4 に示すように、セグメントが未分割の場合にはコピー中の I O 応答時間は 1 0 秒以上であるが、2 5 6 分割以上の場合には平均 0 . 4 秒未満となる。これは、コピーなしの場合の平均レスポンスと同等の応答時間である。但し、図 5 に示すように、分割数が 2 5 6 の場合でも、応答時間が 1 秒を超えるケースが一部に発生し得るため、分割数をある程度増加させてもセグメントの移動の影響を完全に隠ぺいできるとは限らない。

【 0 0 2 0 】

なお、応答時間が 1 秒を超えてしまうリクエストに関しても、分割数をより増やすことで応答時間を減らすことは可能であるが、その分コピー時間が増加してしまう。例えば、図 6 に示すように、2 0 4 8 分割にすると応答時間が 1 秒以上となるリクエストは発生しなくなる。しかし、S S D 2 0 0 から H D D 3 0 0 へ移動する場合の領域移動時間は、1 4 秒 (2 5 6 分割) から 4 5 秒 (2 0 4 8 分割) へと大幅に増加してしまう (図 3 参照) 。

【 0 0 2 1 】

ユーザ I O の応答時間を 1 秒未満にする場合、図 5 及び図 6 の例では、2 0 4 8 分割以上とすることが望ましいが、階層ストレージシステムにおいてそこまでの応答時間保証が要求されていない場合には、2 5 6 分割で十分といえる。このように、保証したいユーザ I O の応答時間を予め定めておき、これを上回らない程度の移動単位を決定することで、移動時間を抑えることができる。

【 0 0 2 2 】

なお、図 3 ~ 図 6 で示した評価 (実験) 結果は、階層ストレージシステムで使用される機器 (例えば S S D 2 0 0 や H D D 3 0 0 , バス等) やワークロードによって変動することが容易に想像できる。この評価結果では、図 4 より 2 5 6 分割が最適値であるが、使用する機器等の条件が変化すればこの分割数も増減する可能性がある。

ところで、移動単位とユーザ I O に加わるオーバーヘッド、つまりレスポンス悪化との関係は、階層ストレージシステムで使用される機器やワークロードによって変化する。す

10

20

30

40

50

なわち、閾値制御のような単純な手法では、機器の性能やワークロードの変化等に応じて、階層ストレージシステムに最適な移動単位で領域移動を行なうことは困難である。

【 0 0 2 3 】

〔 1 - 2 〕 階層ストレージシステムの説明

上述した点に鑑み、本実施形態に係る階層ストレージシステム 1（図 7 参照）は、以下に詳述するように、ユーザ I/O の平均レスポンスを監視することにより、ユーザ I/O のレスポンスの変化に応じて移動単位を動的に変更することができる。これにより、階層ストレージシステム 1 は、直前の平均レスポンスと同等のレスポンスを保ちつつ、分割数を自律的に最小にすることが可能となる。すなわち、第 1 の記憶装置から第 2 の記憶装置へのデータ移動にかかる処理時間を低減させつつ、入力される要求への応答性能の劣化を抑制することができる。

10

【 0 0 2 4 】

図 7 は、一実施形態に係る階層ストレージシステム 1 の構成例を示す図である。図 7 に示すように、階層ストレージシステム（ストレージ装置）1 は、階層ストレージ制御装置 10、SSD 20、及び HDD 30 をそなえる。

階層ストレージ制御装置 10 は、図示しない入力装置や、ネットワークを介したホスト装置からのユーザ I/O に応じて、SSD 20 及び HDD 30 への種々のアクセスを行なうことができる。例えば、階層ストレージ制御装置 10 は、SSD 20 及び HDD 30 へのリード又はライト等のアクセスを行なうことができる。階層ストレージ制御装置 10 としては、PC（Personal Computer）やサーバ、又はコントローラモジュール（CM；Controller Module）等の情報処理装置が挙げられる。

20

【 0 0 2 5 】

また、本実施形態に係る階層ストレージ制御装置 10 は、ユーザ I/O のアクセス頻度に応じて、アクセス頻度が低い領域を HDD 30 に配置する一方、アクセス頻度が高い領域を SSD 20 に配置する、動的階層制御を行なうことができる。

HDD 30 は、種々のデータやプログラム等を格納する記憶装置の一例であり、SSD 20 は、HDD 30 とは異なる性能の（例えばより高速な）記憶装置の一例である。本実施形態において、互いに異なる記憶装置（以下、便宜上、第 1 及び第 2 の記憶装置と表記する場合がある）として、HDD 30 等の磁気ディスク装置、SSD 20 等の半導体ドライブ装置をそれぞれ例に挙げているが、これに限定されるものではない。第 1 及び第 2 の記憶装置として、互いに性能差（例えばリード/ライトの速度差）のある種々の記憶装置が用いられればよい。

30

【 0 0 2 6 】

SSD 20 及び HDD 30 は、階層ストレージシステム 1 におけるストレージボリュームを構成する。SSD 20 及び HDD 30 の各々は、ストレージボリューム上のセグメント（単位領域）のデータを格納可能な記憶領域を含む。セグメントは、階層ストレージ制御装置 10 による階層移動の最小単位であり、図 7 では、1 セグメントが 1 GB であるものとする。階層ストレージ制御装置 10 は、セグメント単位で、SSD 20 - HDD 30 間の領域移動を制御する。

【 0 0 2 7 】

なお、図 7 では、階層ストレージシステム 1 がそれぞれ 1 つの SSD 20 及び HDD 30 をそなえるものとしているが、これに限定されるものではなく、それぞれ複数の SSD 20 及び HDD 30 をそなえてもよい。

40

〔 1 - 3 〕 階層ストレージ制御装置の説明

次に、階層ストレージ制御装置 10 の詳細について説明する。

【 0 0 2 8 】

階層ストレージ制御装置 10 は、一例として、図 7 に示すように、階層管理部 11、階層ドライバ 12、SSD ドライバ 13、及び HDD ドライバ 14 をそなえる。なお、階層管理部 11 は、ユーザ空間で実行されるプログラムとして実現され、階層ドライバ 12、SSD ドライバ 13、及び HDD ドライバ 14 は、OS 空間で実行されるプログラムとし

50

て実現される。

【0029】

階層管理部11は、blktraceを用いて、SSD20又は/及びHDD30についてトレースされたI/Oの情報に基づいて、領域移動を行なうセグメントを判定し、判定したセグメントのデータの移動を階層ドライバ12に指示する。ここで、blktraceは、ブロックI/OレベルでのI/Oをトレースするコマンドである。階層管理部11は、blktraceに代えて、ディスクI/Oの利用状況を確認するコマンドであるiostatを用いてもよい。なお、blktrace及びiostatはOS空間で実行される。

【0030】

階層管理部11は、データ収集部11a, データベース11b, ワークロード分析部11c, 移動指示部11d, 及び分割数判定部11eをそなえる。

なお、階層管理部11の動作を実現するために、階層ストレージシステム1の管理者等は、予め以下の情報を決定しておくことが好ましい。

- ・全セグメントがHDD30にある場合の平均レスポンス(コピーなし), 並びに平均レスポンスがコピーなしの場合と同等になる分割数。図4に示す例では、分割数は256である。なお、平均レスポンス及び分割数は、階層ストレージシステム1で使用予定の機器を用いて、図4に示すような実験を行なうことで求められる。

【0031】

- ・平均レスポンスを求める期間。例えば60s程度。
- ・平均レスポンスの誤差範囲。例えば50ms程度。

データ収集部11aは、blktraceを用いてSSD20又は/HDD30についてトレースされたI/Oの情報を所定間隔(例えば1分間隔)で収集する。また、データ収集部11aは、収集した情報に基づいて、セグメントごとに、例えば、セグメントを特定する情報、合計I/O数(iopm; I/O per minute), 及び平均レスポンス(応答性能)を集計する。そして、データ収集部11aは、集計結果をタイムスタンプとともにデータベース11bに書き込む。なお、セグメントを特定する情報としては、ボリューム上のオフセットに関する情報を用いることができる。

【0032】

また、データ収集部11aは、全セグメントを対象にした(全セグメントの)合計I/O数及び平均レスポンスも集計し、タイムスタンプとともにデータベース11bに書き込むことができる。このとき、データ収集部11aは、全セグメントを対象にした情報をデータベース11bに追加したことを分割数判定部11eに通知してもよい。

なお、データ収集部11aは、各セグメント又は/及び全セグメントへのI/Oのリードライト比(rw比)を集計し、上述した情報に含めてもよい。

【0033】

このように、データ収集部11aは、SSD20又はHDD30で使用される領域を所定の大きさで分割した複数の単位領域について、入力された要求に対する応答性能を監視する監視部の一例である。

データベース11bは、データ収集部11aにより集計されたセグメントに関する情報を記憶するものであり、例えば図示しないメモリ等により実現される。

【0034】

図8は、図7に示すデータベース11bの一例を示す図である。図8に示すように、データベース11bは、セグメントごとに、セグメントを特定する情報、I/O数、平均レスポンス、及びタイムスタンプを対応させて記憶するテーブルである。例えば、セグメント“1”であるセグメントは、合計I/O数が“1000”, 平均レスポンスが“0.6”(秒), タイムスタンプが“1”である。

【0035】

なお、セグメントを特定する情報として、セグメントの番号を用いているが、これに代えてストレージボリュームの先頭オフセットが用いられてもよい。ここで、I/O数は、セ

10

20

30

40

50

グメントに対して1分間に行なわれたI/Oの合計数であり、平均レスポンスは、階層ストレージ制御装置10がセグメントへのI/Oを受け取ってからレスポンスを送信するまでにかかった時間の平均である。タイムスタンプは、時刻を識別する識別子であり、例えば時刻そのものが設定されてもよい。

【0036】

また、図8において、セグメントが“all”のエントリは、全セグメントを対象にした集計結果である。セグメントが“all”のエントリについては、後述する分割数判定部11eにより、過去n個のデータが参照されるため、複数の“all”のエントリを追加できるようになっている。なお、“all”のエントリの新旧は、タイムスタンプにより識別可能である。一方、個々のセグメントのエントリについては、同一セグメントのデータを上書きできるようにしてもよいし、“all”と同様に複数のデータを登録できるようにしてもよい。

10

【0037】

ワークロード分析部11cは、データベース11bが記憶するセグメントから、SSD20又はHDD30にデータを移動するセグメントを選択し、選択したセグメントに関する情報を移動指示部11dに渡す。

一例として、ワークロード分析部11cは、セグメント数が同時に階層移動を行なう最大のセグメント数(所定数)に達するまで、I/O数が多い順にセグメントを抽出することができる。又は、ワークロード分析部11cは、SSD20にデータを移動するセグメントとして、I/O数又はアクセスの集中度(全体に対するI/O数の割合)が所定の閾値よりも高いセグメントを抽出してもよい。

20

【0038】

また、ワークロード分析部11cは、HDD30にデータを移動するセグメントとして、例えばI/O数が上記所定数に入らなかった、又はI/O数若しくはアクセスの集中度が所定の閾値以下となった、SSD20上のセグメントを抽出することができる。

なお、ワークロード分析部11cは、上記のSSD20又はHDD30にデータを移動するセグメントの抽出条件に所定回数以上連続して該当したときに、当該セグメントを、SSD20又はHDD30にデータを移動するセグメントとして抽出してもよい。また、ワークロード分析部11cは、上記I/O数等の他に、リードライト比(rw比)に基づいてセグメントを選択してもよい。

30

【0039】

ここで、ワークロード分析部11cは、移動指示部11dに対して、HDD30内のセグメントについてSSD20への階層移動を指示した後、SSD20内の他のセグメントについてHDD30への階層移動を指示することができる。一方、ワークロード分析部11cは、あるセグメントについてSSD20への階層移動を行なっている間に当該セグメントへの負荷が下がると予測される場合には、他のセグメントについてのみHDD30への階層移動を指示してもよい。

【0040】

例えば、ワークロード分析部11cは、スパイクの平均余命時間と階層移動にかかる時間とに基づいて、階層移動中のセグメントへの負荷が下がるか否かを判断できる。なお、スパイクとは、一部のセグメントに負荷が集中することであり、平均余命時間は、負荷が継続する継続時間から既に実行済みの実行時間を減じた時間であり、ワークロードに応じて定まる値である。管理者等は予め平均余命時間を求め、階層ストレージ制御装置10に設定しておくことができる。

40

【0041】

具体的には、ワークロード分析部11cは、SSD20にデータを移動するセグメントを抽出し、抽出したセグメントについてSSD20にデータを移動するコスト(時間)を計算する。そして、ワークロード分析部11cは、平均余命時間が移動時間以下になる場合には、SSD20からHDD30への階層移動のみを行なうと判断することができる。

移動指示部11dは、ワークロード分析部11cからの指示に基づいて、階層ドライバ

50

12に、選択されたセグメントのデータの、HDD30からSSD20への移動、又は、SSD20からHDD30への移動を指示する。このとき、移動指示部11dは、選択されたセグメントのストレージボリューム上のオフセットをHDD30上のオフセットに変換してセグメントごとにデータの移動を指示する。例えば、HDD30のセクターサイズが512Bである場合に、ボリューム上のオフセットが1GBであれば、HDD30上でのオフセットは $1 \times 1024 \times 1024 \times 1024 / 512 = 2097152$ となる。

【0042】

また、移動指示部11dは、移動指示を出したセグメント数と、ワークロード分析部11cで移動の判定が行なわれたデータのタイムスタンプ（直近のタイムスタンプ）とを含む移動開始通知を、分割数判定部11eに通知する。

10

分割数判定部11eは、セグメントの分割数を決定するとともに、セグメントのデータの移動開始前後のI/Oレスポンスの変化に基づいて、動的に分割数を変更する。

【0043】

具体的には、分割数判定部11eは、移動指示部11dから移動開始通知を受け取ると、直前の分割数判定結果（立ち上げ直後なら初期値、例えば256）を階層ドライバ12（分割部12d）に通知する。階層ドライバ12は、通知された分割数に従ってセグメントの分割を行ない、データ移動を進める。

また、分割数判定部11eは、データ収集部11aにより求められたセグメント移動直前の平均レスポンスを取得し、予め決められた平均レスポンス誤差範囲値より、セグメント移動中の平均レスポンスの期待値を求めておく。なお、セグメント移動直前の平均レスポンスは、移動開始通知に含まれるタイムスタンプに対応する、データベース11b内の「全セグメントを対象にした平均レスポンス」を取り出すことで取得できる。例えば、平均レスポンスが400msであり、平均レスポンス誤差範囲値が50msである場合、350～450msの範囲が期待値となる。

20

【0044】

さらに、分割数判定部11eは、セグメント移動中に、新たにレスポンスが入力されると、期待値の範囲に平均レスポンスが収まっているか否かを評価する。

具体的には、分割数判定部11eは、セグメント移動中に、データ収集部11aにより全セグメントを対象にした平均レスポンスが所定数（例えばn個）求められると、これらn個のデータの平均値を計算し、セグメント移動中の平均レスポンスとする。そして、分割数判定部11eは、セグメント移動中の平均レスポンスが期待値よりも大きい（応答性能が劣化している）場合には、セグメントの分割数を現在の設定値よりも大きくして、階層ドライバ12に通知する。一方、分割数判定部11eは、セグメント移動中の平均レスポンスが期待値よりも小さい（応答性能が優れている）場合には、セグメントの分割数を現在の設定値よりも小さくして、階層ドライバ12に通知する。

30

【0045】

例えば、分割数判定部11eは、分割数を現在の設定値よりも大きくする場合、2倍（例えば256から512）、3倍、・・・のように変化させてもよいし、所定値を加算させてもよい。また、分割数判定部11eは、分割数を現在の設定値よりも小さくする場合、1/2倍（例えば256から128にして）、1/3倍、・・・のように変化させてもよいし、所定値を減算させてもよい。

40

【0046】

このように、分割数判定部11eは、セグメント移動中の平均レスポンスが期待値から外れる場合には、領域移動時の平均レスポンスが移動直前のレスポンスに基づく期待値内に収まるように、セグメント移動中に分割数を動的に変更するのである。これにより、セグメント移動によって、レスポンスを基準としてのセグメント移動前のレスポンスから大きく劣化させずに済み、システムの安定性を保証することができる。

【0047】

なお、分割数判定部11eは、データ収集部11aが階層移動中の複数の時点で監視した複数（n個）の平均レスポンスの平均値を用いることで、ごく短時間にアクセスが集中

50

した場合のようなワークロードの突発的な変化の影響を緩和することができる。

なお、ここまで、分割数判定部 11 e は、データ収集部 11 a が監視した移動処理実行中の第 1 の応答性能と移動処理実行前の第 2 の応答性能とに基づいて、分割数を変更するものとして説明したが、これに限定されるものではない。

【0048】

例えば、分割数判定部 11 e は、移動処理実行中の応答性能に基づいて、分割数を変化させてもよい。一例として、分割数判定部 11 e は、セグメント移動中に、階層ドライバ 12 に通知した分割数でのレスポンスを取得するとともに、分割数を大きく又は小さくさせ、そのときのレスポンスを取得する。そして、分割数判定部 11 e は、セグメント移動中に取得した複数のレスポンスを比較し、直近のレスポンスが 1 つ前のレスポンスよりも大きいのか（応答性能が劣化しているのか）に応じて、上述の如く分割数を変化させることができる。なお、分割数判定部 11 e は、比較したレスポンスの差が平均レスポンスの誤差範囲値を超える場合に、分割数を変化させてもよい。

10

【0049】

以上のように、分割数判定部 11 e は、データ収集部 11 a が監視した移動処理実行中の第 1 の応答性能に基づいて、分割数を変更する変更部の一例であるといえる。

なお、分割数判定部 11 e は、セグメントの分割数を決定し、決定した分割数を階層ドライバ 12（分割部 12 d）に通知するものとして説明したが、移動対象の単位領域を移動する際の移動単位を決定し、通知してもよい。つまり、分割数判定部 11 e は、移動対象のセグメントについて一度に転送するデータサイズ（移動単位の大きさ）を、上述した判定により決定して（変化させて）、階層ドライバ 12 に通知してもよい。

20

【0050】

階層ドライバ 12 は、I/O マップ部 12 a、ペンディングキュー 12 b、階層テーブル 12 c、及び分割部 12 d をそなえる。

I/O マップ部 12 a は、ユーザからのストレージボリュームに対する I/O 要求を階層テーブル 12 c を用いて SSD ドライバ 13 又は HDD ドライバ 14 に振り分け、SSD ドライバ 13 又は HDD ドライバ 14 からの I/O レスポンスをユーザに返す。

【0051】

ペンディングキュー 12 b は、I/O 要求を一時的に格納する保持部であり、図示しないメモリ等により実現される。I/O マップ部 12 a は、階層移動中のセグメントに対して I/O 要求が発行されると、当該セグメントのデータの移動が完了するまで、当該 I/O 要求をペンディングキュー 12 b に格納し、I/O 要求を保留する。データの移動が完了すると、I/O マップ部 12 a は、ペンディングキュー 12 b から当該 I/O 要求を読み出して、SSD ドライバ 13 又は HDD ドライバ 14 への振り分けを再開する。

30

【0052】

階層テーブル 12 c は、I/O マップ部 12 a による I/O 要求の振り分け及び分割部 12 d による階層制御に用いられるテーブルであり、例えば図示しないメモリ等により実現される。

図 9 は、図 7 に示す階層テーブル 12 c の一例を示す図である。図 9 に示すように、階層テーブル 12 c は、SSD 20 にデータが移動されたセグメントごとに、SSD オフセットと、HDD オフセットと、状態とを対応させて記憶するテーブルである。

40

【0053】

SSD オフセットは、SSD 20 にデータが移動されたセグメントの SSD 20 におけるオフセットを示す。SSD オフセットは、ボリューム上のサイズ 1 GB に対応するオフセット“2097152”を単位とする固定値であり、例えば、“0”、“2097152”、“4194304”、“6291456”、...となる。

HDD オフセットは、SSD 20 にデータが移動されたセグメントの HDD 30 におけるオフセットを示す。HDD オフセットの値“NULL”は、SSD オフセットで指定される SSD 20 の領域が未使用であることを示す。

【0054】

50

状態は、セグメントの状態を示し、“allocated”、“Moving (HDD SSD)”、“Moving (SSD HDD)”、又は“free”である。“allocated”はセグメントがSSD 20に割り当てられていることを示し、“Moving (HDD SSD)”はセグメントのデータがHDD 30からSSD 20に転送中であることを示す。“Moving (SSD HDD)”はセグメントのデータがSSD 20からHDD 30に転送中であることを示し、“free”はSSDオフセットで指定されるSSD 20の領域が未使用であることを示す。

【0055】

IOマップ部12aは、上述した階層テーブル12cを参照することで、IO要求をSSDドライバ13又はHDDドライバ14のいずれに振り分けるかを判定することができる。10

図7の説明に戻り、階層ドライバ12は、移動指示部11dからセグメント移動指示を受け取ると、HDD 30又はSSD 20の移動対象の単位領域に記憶されたデータをSSD 20又はHDD 30に移動する移動処理を実行する。具体的には、階層ドライバ12は、階層テーブル12c及び分割部12dにより、セグメント移動指示で指定されたセグメントのデータをSSD 20 - HDD 30間で移動する。

【0056】

より具体的に、階層ドライバ12は、セグメント移動指示を受け取ると、階層テーブル12c内のHDDオフセットより“NULL”となっているエントリを探し、セグメント移動指示で指定されたHDDオフセット情報と、状態とを登録する。なお、このとき登録される状態は、“Moving (HDD SSD)”又は“Moving (SSD HDD)”である。そして、階層ドライバ12は、SSD 20 - HDD 30間のデータの転送指示を分割部12dに送出する。20

【0057】

また、階層ドライバ12は、データの転送完了を通知されると、階層テーブル12cから転送が完了したエントリを探し、状態が“Moving (HDD SSD)”の場合、当該状態を“allocated”に変更する。一方、階層テーブル12cは、状態が“Moving (SSD HDD)”である場合、当該状態を“free”に変更するとともに、対応するHDDオフセットを“NULL”に設定する。

【0058】

分割部12dは、階層ドライバ12からのSSD 20 - HDD 30間のデータの転送指示に応じて、分割数判定部11eから指示された分割数でセグメントを分割し、セグメントのデータの階層移動を行なう。30

具体的には、分割部12dは、階層ドライバ12から転送指示を受け取ると、転送指示に係る各セグメントを、それぞれ分割数判定部11eから指示された分割数mmで分割し、分割した単位でkcopydに転送指示を発行する。そして、分割部12dは、kcopydにより、分割した全ての領域内のデータの転送が完了すると、階層ドライバ12にデータの転送完了を通知する。

【0059】

また、分割部12dは、分割数判定部11eから分割数mmの更新要求を受け取ると、当該要求に応じて、分割数mmの更新を行なう。例えば、分割部12dは、分割数を2倍にする指示を受け取ると、 $mm = mm * 2$ を計算して、分割数mmを更新する。また、分割部12dは、分割数を1/2倍にする指示を受け取ると、 $mm = mm / 2$ を計算して、分割数mmを更新する。40

【0060】

なお、分割部12dは、階層ドライバ12からデータの転送指示を受けたときの分割数（移動単位）として、分割数判定部11eによりセグメント移動前（例えば移動直前）のレスポンスを基に求められた所定の分割数（移動単位）を用いることができる。これにより、分割数（移動単位）はセグメント移動前のレスポンスを考慮して設定されるので、kcopydによるデータの転送が開始したときの急激なレスポンスの低下を抑制することができる。50

【 0 0 6 1 】

このように、分割部 1 2 d は、移動対象の単位領域を所定の分割数により複数の分割領域に分割し、移動対象の単位領域に記憶されたデータを分割領域の単位で S S D 2 0 又は H D D 3 0 に移動するのである。換言すれば、分割部 1 2 d は、分割数判定部 1 1 e からの指示に応じて移動対象の単位領域の移動単位を変更し、変更した移動単位で k c o p y d にデータ転送を指示するのである。

【 0 0 6 2 】

S S D ドライバ 1 3 は、階層ドライバ 1 2 の指示に基づいて S S D 2 0 へのアクセスを制御する。H D D ドライバ 1 4 は、階層ドライバ 1 2 の指示に基づいて H D D 3 0 へのアクセスを制御する。

10

以上のように、本実施形態に係る階層ストレージシステム 1 によれば、ユーザ I O の平均レスポンスを監視し、ユーザ I O のレスポンス変化に応じてレスポンス悪化が収束する領域サイズに移動単位を動的に設定（変更）することができる。従って、ユーザ I O へのレスポンス悪化と階層移動時間のバランスを適切にとることができ、セグメントの階層移動時の平均レスポンスを可能な限り小さくしつつ、できるだけ少ない分割数で（短時間で）セグメントの階層移動を実現することができる。

【 0 0 6 3 】

すなわち、本実施形態に係る階層ストレージシステム 1 によれば、分割数判定部 1 1 e 及び分割部 1 2 d により、使用する機器の性能やワークロードに応じた最適な移動単位で、S S D 2 0 - H D D 3 0 間のデータの階層移動を行なうことができるのである。

20

〔 1 - 4 〕 階層ストレージシステムの動作例

次に、上述の如く構成された階層ストレージシステム 1 の動作例を、図 1 0 ~ 図 1 8 を参照して説明する。

【 0 0 6 4 】

はじめに、図 1 0 を参照してデータ収集部 1 1 a の動作を説明する。図 1 0 はデータ収集部 1 1 a によるデータ収集処理の動作例を示すフローチャートである。なお、データ収集部 1 1 a は、b l k t r a c e コマンドを 6 0 秒間実行して、終了することを条件として起動される。

図 1 0 に示すように、データ収集部 1 1 a により、b l k t r a c e コマンドの実行により得られたトレース結果が取り出される（ステップ S 1）。次いで、データ収集部 1 1 a により、1 G B オフセット単位すなわちセグメント単位で各セグメントの I O 数及び平均レスポンスが集計され、タイムスタンプとともにデータベース 1 1 b に書き込まれる（ステップ S 2）。

30

【 0 0 6 5 】

そして、データ収集部 1 1 a により、全セグメントを対象にした合計 I O 数、平均レスポンスが集計され、タイムスタンプとともにデータベース 1 1 b に格納される（ステップ S 3）。なお、データ収集部 1 1 a は、ステップ S 3 の処理を実行した旨を分割数判定部 1 1 e に通知してもよい。

このように、データ収集部 1 1 a は、定期的に全てのセグメントの平均レスポンスを監視することで、流動的に変化するワークロードがユーザ I O に与える影響を分割数判定部 1 1 e にフィードバックすることができる。

40

【 0 0 6 6 】

次に、図 1 1 を参照してワークロード分析部 1 1 c の動作を説明する。図 1 1 はワークロード分析部 1 1 c による移動判定処理の動作例を示すフローチャートである。

図 1 1 に示すように、ワークロード分析部 1 1 c により、データベース 1 1 b から直近のタイムスタンプのセグメントについて I O 数が取り出される（ステップ S 1 1）。そして、ワークロード分析部 1 1 c により、セグメント数が所定数に達するまで、I O 数が多い順に候補セグメントが抽出される（ステップ S 1 2）。

【 0 0 6 7 】

次いで、ワークロード分析部 1 1 c により、予め求められた平均余命時間が全候補セグ

50

メントにかかる移動時間よりも大きいか否かが判定される（ステップS 1 3）。平均余命時間が当該移動時間以下の場合（ステップS 1 3のNoルート）、処理がステップS 1 5に移行する。一方、平均余命時間が当該移動時間よりも大きい場合（ステップS 1 3のYesルート）、ワークロード分析部1 1 cにより、候補セグメントの情報が移動指示部1 1 dに通知され、データの移動（HDD 3 0からSSD 2 0）が指示される（ステップS 1 4）。

【0068】

ステップS 1 5では、ワークロード分析部1 1 cにより、SSD 2 0上のセグメントから候補セグメントに含まれないセグメント、つまりI/O数が比較的少ないセグメントが抽出される。そして、ワークロード分析部1 1 cにより、抽出したセグメントの情報が移動指示部1 1 dに通知され、データの移動（SSD 2 0からHDD 3 0）が指示される（ステップS 1 6）。

【0069】

そして、ワークロード分析部1 1 cは、所定時間、例えば60秒スリープし（ステップS 1 7）、処理がステップS 1 1に移行する。

なお、ワークロード分析部1 1 cは、ステップS 1 2において、I/O数又はアクセスの集中率（全体に対するI/O数の割合）が所定の閾値よりも高いセグメントを抽出してもよい。また、ワークロード分析部1 1 cは、ステップS 1 5において、HDD 3 0にデータを移動するセグメントとして、例えばI/O数又はアクセスの集中率が所定の閾値以下となったSSD 2 0上のセグメントを抽出してもよい。さらに、ワークロード分析部1 1 cは、ステップS 1 2及びS 1 5で抽出するセグメントとして、当該抽出条件に所定回数以上連続して該当したセグメントを選択してもよい。

【0070】

このように、ワークロード分析部1 1 cは、I/Oの集中度が高いセグメントのデータをHDD 3 0からSSD 2 0に移動するように移動指示部1 1 dに指示することによって、ユーザはHDD 3 0のデータに高速にアクセスすることができる。また、ワークロード分析部1 1 cは、I/Oの集中度が低くなったセグメントのデータをSSD 2 0からHDD 3 0に移動するように移動指示部1 1 dに指示することによって、比較的高価格、低容量のSSD 2 0を有効利用することができる。

【0071】

次に、図12を参照して移動指示部1 1 dの動作を説明する。図12は移動指示部1 1 dによる移動指示通知処理の動作例を示すフローチャートである。

図12に示すように、移動指示部1 1 dにより、ワークロード分析部1 1 cからの移動指示が待ち受けられる（ステップS 2 1）。移動指示を受け取ると、移動指示部1 1 dにより、各セグメントのボリューム上のオフセットがHDD 3 0上のオフセットに変換される（ステップS 2 2）。

【0072】

そして、移動指示部1 1 dにより、セグメントごとに、HDD 3 0上のオフセットと、データの移動方向とが階層ドライバ1 2に通知される（ステップS 2 3）。ここで、データの移動方向は、HDD 3 0からSSD 2 0か、又は、SSD 2 0からHDD 3 0である。そして、移動指示部1 1 dにより、移動指示を出したセグメント数と移動判断を行なったデータのタイムスタンプ（直近のタイムスタンプ）とが分割数判定部1 1 eに通知され（ステップS 2 4）、処理がステップS 2 1に移行する。

【0073】

このように、移動指示部1 1 dが各セグメントのボリューム上のオフセットをHDD 3 0上のオフセットに変換することによって、階層ドライバ1 2はSSD 2 0とHDD 3 0との間でデータを移動することができる。

次に、図13を参照して分割数判定部1 1 eの動作を説明する。図13は分割数判定部1 1 eによる分割数判定処理の動作例を示すフローチャートである。

【0074】

図 1 3 に示すように、分割数判定部 1 1 e により、移動指示部 1 1 d からのセグメントの移動情報（セグメント数、タイムスタンプ）が待ち受けられる（ステップ S 3 1）。移動情報を受け取ると、分割数判定部 1 1 e により、受け取ったタイムスタンプ $timestamp_org$ に対応するデータベース 1 1 b のデータへのアクセスが行なわれ、全セグメントの平均レスポンス $resp_org$ が取り出される（ステップ S 3 2）。

【 0 0 7 5 】

次いで、分割数判定部 1 1 e は、 $timestamp_org$ よりも新しいデータが n 個データベース 1 1 b に登録されるまで（例えば $60 \times n + 10$ 秒）スリープする（ステップ S 3 3）。

新しいデータが n 個データベース 1 1 b に登録されると、分割数判定部 1 1 e により、新しい n 個のデータにアクセスされ、全データの全セグメントの平均レスポンスが取り出される。そして、分割数判定部 1 1 e により、取り出した平均レスポンスの平均値 $resp_new$ が求められる（ステップ S 3 4）。

【 0 0 7 6 】

次に、分割数判定部 1 1 e により、 $resp_new > resp_org + m$ が成立するか否かが判定される（ステップ S 3 5）。 $resp_new > resp_org + m$ が成立する場合（ステップ S 3 5 の Yes ルート）、分割数判定部 1 1 e により、分割部 1 2 d に分割数を増加させる（例えば現在の 2 倍にする）指示が発行され（ステップ S 3 6）、処理がステップ S 3 1 に移行する。なお、 m は、平均レスポンスの誤差範囲値であり、例えば 50 ms に設定することができる。

【 0 0 7 7 】

一方、 $resp_new > resp_org + m$ が成立しない場合（ステップ S 3 5 の No ルート）、分割数判定部 1 1 e により、 $resp_new < resp_org + m$ が成立するか否かが判定される（ステップ S 3 7）。 $resp_new < resp_org + m$ が成立する場合（ステップ S 3 7 の Yes ルート）、分割数判定部 1 1 e により、分割部 1 2 d に分割数を減少させる（例えば現在の $1/2$ 倍にする）指示が発行され（ステップ S 3 8）、処理がステップ S 3 1 に移行する。なお、 $resp_new < resp_org + m$ が成立しない場合（ステップ S 3 7 の No ルート）、 $resp_new$ は平均レスポンスの期待値の範囲内にあるため、分割数の更新は行なわれず、処理がステップ S 3 1 に移行する。

【 0 0 7 8 】

このように、分割数判定部 1 1 e が、セグメント移動前及び移動中の平均レスポンスに基づいて、セグメント移動による平均レスポンスの悪化を抑制するように分割数を決定することができる。従って、階層ドライバ 1 2 は、移動中のセグメントを動的に最適な分割数に分割することができるため、移動時間を低減させつつ、対象データに対するユーザ I/O へのレスポンス悪化を抑制することができる。

【 0 0 7 9 】

次に、図 1 4 及び図 1 5 を参照して階層ドライバ 1 2 の動作を説明する。

はじめに、移動指示を受け取ったときの階層ドライバ 1 2 の動作を説明する。図 1 4 は、階層ドライバ 1 2 による転送指示通知処理の動作例を示すフローチャートである。

図 1 4 に示すように、階層ドライバ 1 2 により、移動指示部 1 1 d からの移動指示が待ち受けられ（ステップ S 4 1）、移動指示を受け取ると、HDD 3 0 から SSD 2 0 へのデータの移動であるか否かが判定される（ステップ S 4 2）。

【 0 0 8 0 】

HDD 3 0 から SSD 2 0 へのデータの移動である場合（ステップ S 4 2 の Yes ルート）、階層ドライバ 1 2 により、移動を指示されたセグメントが SSD 2 0 へ移動済みであるか否かが判定される（ステップ S 4 3）。移動済みである場合（ステップ S 4 3 の Yes ルート）、処理がステップ S 4 1 に移行する。

一方、移動済みでない場合（ステップ S 4 3 の No ルート）、階層ドライバ 1 2 により、階層テーブル 1 2 c 内の HDD オフセットより “NULL” となっているエントリが探索さ

10

20

30

40

50

れ、HDDオフセット情報と状態とが登録される。このとき階層ドライバ12が登録する状態は、“Moving(HDD SSD)”である。そして、階層ドライバ12により、HDD30からSSD20へのデータの転送指示が分割部12dに発行され(ステップS44)、処理がステップS41に移行する。

【0081】

また、HDD30からSSD20へのデータの移動でない場合(ステップS42のNoルート)、階層ドライバ12により、階層テーブル12c内のHDDオフセットよりセグメントが探索され、HDDオフセット情報と状態とが登録される。このとき階層ドライバ12が登録する状態は、“Moving(SSD HDD)”である。そして、階層ドライバ12により、SSD20からHDD30へのデータの転送指示が分割部12dに発行され(ステップS45)、処理がステップS41に移行する。

10

【0082】

次に、転送指示後に転送完了通知を受け取ったときの階層ドライバ12の動作を説明する。図15は、階層ドライバ12による転送完了受信処理の動作例を示すフローチャートである。

図15に示すように、階層ドライバ12により、分割部12dからの転送完了通知が待ち受けられる(ステップS51)。転送完了通知を受け取ると、階層ドライバ12により、転送が完了した階層テーブル12cのエントリがHDDオフセットを用いて探索され、状態が“Moving(HDD SSD)”の場合は状態が“allocated”に変更される。一方、階層テーブル12cにより、状態が“Moving(SSD HDD)”の場合は状態が“free”に変更され、且つ、対応するHDDオフセットが“NULL”に設定され(ステップS52)、処理がステップS51に移行する。

20

【0083】

このように、階層ドライバ12が階層テーブル12cを用いてSSD20とHDD30との間でデータを転送することにより、I/Oが集中するセグメントのデータをSSD20に置くことができる。

次に、図16及び図17を参照して分割部12dの動作を説明する。

はじめに、転送指示を受け取ったときの分割部12dの動作を説明する。図16は、分割部12dによる転送指示受信処理の動作例を示すフローチャートである。

【0084】

30

図16に示すように、分割部12dにより、階層ドライバ12からのSSD20-HDD30間の転送指示が待ち受けられる(ステップS61)。転送指示を受け取ると、分割部12dにより、転送指示で指定された移動する各セグメントが分割数mmで分割され、分割単位でkcopydに転送指示が発行される(ステップS62)。

全データの転送が終わると、分割部12dにより、階層ドライバ12にデータの転送完了が通知され(ステップS63)、処理がステップS61に移行する。

【0085】

次に、分割数更新指示を受け取ったときの分割部12dの動作を説明する。図17は、分割部12dによる分割数更新処理の動作例を示すフローチャートである。

図17に示すように、分割部12dにより、分割数判定部11eからの分割数更新指示が待ち受けられる(ステップS71)。分割数更新指示を受け取ると、分割部12dにより、当該指示に応じて分割数mmが更新され(ステップS72)、処理がステップS71に移行する。

40

【0086】

このように、分割部12dが階層ドライバ12からのセグメント単位の転送指示を、さらに小さい移動単位に分割することによって、ユーザI/Oのレスポンス悪化を抑制することができる。また、分割部12dが分割数判定部11eからの分割数更新指示に応じて分割数mmを適宜更新することができるため、ワークロードの変化に柔軟に対応することができる。

【0087】

50

次に、図 18 を参照して I O マップ部 12 a の動作を説明する。図 18 は I O マップ部 12 a による I O 受信処理の動作例を示すフローチャートである。

図 18 に示すように、I O マップ部 12 a により、ユーザ I O の受信が待ち受けられる (ステップ S 81)。ユーザ I O を受け取ると、I O マップ部 12 a により、ユーザ I O で指定されるオフセットと階層テーブル 12 c に登録されている各オフセット + セグメントサイズとが比較され (ステップ S 82)。

【0088】

そして、I O マップ部 12 a により、比較の結果、階層テーブル 12 c に一致するオフセットが存在し、且つ状態が “allocated” であるか否かが判定される (ステップ S 83)。一致するオフセットが存在し、且つ状態が “allocated” である場合 (ステップ S 83 の Yes ルート)、I O マップ部 12 a により、SSD ドライバ 13 へ I O 要求が送付され (ステップ S 84)、処理がステップ S 81 に移行する。

10

【0089】

一方、一致するオフセットが存在せず、又は状態が “allocated” ではない場合 (ステップ S 83 の No ルート)、I O マップ部 12 a により、状態が “Moving(HDD SSD)” 又は “Moving(SSD HDD)” であるか否かが判定される (ステップ S 85)。状態が “Moving(HDD SSD)” 又は “Moving(SSD HDD)” ではない場合 (ステップ S 85 の No ルート)、I O マップ部 12 a により、HDD ドライバ 14 へ I O 要求が送付され (ステップ S 86)、処理がステップ S 81 に移行する。

【0090】

20

また、状態が “Moving(HDD SSD)” 又は “Moving(SSD HDD)” である場合 (ステップ S 85 の Yes ルート)、I O マップ部 12 a により、当該状態が “free” 又は “allocated” に変化するまで I O 要求がペンディングキュー 12 b に格納される。すなわち、I O マップ部 12 a により、I O 要求に係るセグメントの階層移動が完了するまで、I O 要求が保留にされる (ステップ S 87)。階層移動が完了すると、I O マップ部 12 a によりペンディングキュー 12 b に格納された I O 要求が取り出され、処理がステップ S 83 に移行する。

【0091】

〔1-5〕ハードウェア構成例

次に、図 19 を参照して、図 7 に示す階層ストレージ制御装置 10 のハードウェア構成について説明する。図 19 は、図 7 に示す階層ストレージ制御装置 10 のハードウェア構成例を示す図である。

30

階層ストレージ制御装置 10 は、図 19 に示すように、CPU (Central Processing Unit) 10 a、メモリ 10 b、記憶部 10 c、インタフェース部 10 d、入出力部 10 e、記録媒体 10 f、及び読取部 10 g をそなえることができる。

【0092】

CPU 10 a は、対応する各ブロック 10 b ~ 10 g と接続され、種々の制御や演算を行なう演算処理装置 (プロセッサ) である。CPU 10 a は、メモリ 10 b、記憶部 10 c、記録媒体 10 f や 10 h、又は図示しない ROM (Read Only Memory) 等に格納されたプログラムを実行することにより、階層ストレージ制御装置 10 における種々の機能を実現することができる。

40

【0093】

メモリ 10 b は、種々のデータやプログラムを格納する記憶装置である。CPU 10 a は、プログラムを実行する際に、メモリ 10 b にデータやプログラムを格納し展開する。なお、メモリ 10 b としては、例えば RAM (Random Access Memory) 等の揮発性メモリが挙げられる。

記憶部 10 c は、種々のデータやプログラム等を格納するハードウェアである。記憶部 10 c としては、例えば HDD 等の磁気ディスク装置、SSD 等の半導体ドライブ装置、フラッシュメモリ等の不揮発性メモリ等の各種デバイスが挙げられる。なお、記憶部 10 c として複数のデバイスが用いられてもよく、これらのデバイスで RAID (Redundant

50

Arrays of Inexpensive Disks) が構成されてもよい。また、記憶部 10c は、図 7 に示す SSD 20 及び HDD 30 を含んでもよい。

【0094】

インタフェース部 10d は、有線又は無線による、ネットワーク(図示省略)や他の情報処理装置との間の接続及び通信の制御等を行なうものである。インタフェース部 10d としては、例えば、LAN(Local Area Network)、ファイバチャネル(Fibre Channel; FC)、インフィニバンド(InfiniBand)等に準拠したアダプタが挙げられる。

入出力部 10e は、マウスやキーボード等の入力装置及びディスプレイやプリンタ等の出力装置の少なくとも一方を含むことができる。例えば、入出力部 10e は、階層ストレージ制御装置 10 の使用者又は管理者等による種々の作業に用いられる。

10

【0095】

記録媒体 10f は、例えばフラッシュメモリや ROM 等の記憶装置であり、種々のデータやプログラムを記録することができる。読取部 10g は、コンピュータ読取可能な記録媒体 10h に記録されたデータやプログラムを読み出す装置である。記録媒体 10f 及び 10h の少なくとも一方には、本実施形態に係る階層ストレージ制御装置 10 の各種機能の全部もしくは一部を実現する制御プログラムが格納されてもよい。例えば、CPU 10a は、記録媒体 10f から読み出したプログラム、又は、読取部 10g を介して記録媒体 10h から読み出したプログラムを、メモリ 10b 等の記憶装置に展開して実行することができる。これにより、コンピュータ(CPU 10a、情報処理装置、各種端末を含む)は、上述した階層ストレージ制御装置 10 の機能を実現することができる。

20

【0096】

なお、記録媒体 10h としては、例えばフレキシブルディスク、CD(Compact Disc)、DVD(Digital Versatile Disc)、ブルーレイディスク等の光ディスクや、USB(Universal Serial Bus)メモリやSDカード等のフラッシュメモリが挙げられる。なお、CDとしては、CD-ROM、CD-R(CD-Recordable)、CD-RW(CD-Rewritable)等が挙げられる。また、DVDとしては、DVD-ROM、DVD-RAM、DVD-R、DVD-RW、DVD+R、DVD+RW等が挙げられる。

【0097】

なお、上述した各ブロック 10a ~ 10g 間はバスで相互に通信可能に接続される。例えば CPU 10a と記憶部 10c との間は、ディスクインタフェースを介して接続される。また、階層ストレージ制御装置 10 の上述したハードウェア構成は例示である。従って、階層ストレージ制御装置 10 内でのハードウェアの増減(例えば任意のブロックの追加や省略)、分割、任意の組み合わせでの統合、バスの追加又は省略等は適宜行なわれてもよい。

30

【0098】

〔1-6〕適用例

上述のように、階層ストレージ制御装置 10 は、リアルタイムに測定した負荷に基づいて高負荷領域のデータを SSD 20 に移動する動的階層制御に用いて好適である。

ここで、階層ストレージ制御装置 10 は、高負荷領域の近傍のデータを SSD 20 に移動するために、近傍として適切な領域を選択する機能をさらにそなえてもよい。すなわち、階層ストレージ制御装置 10 を、以下に詳述する階層ストレージ制御装置 10A(図 22 参照)に適用してもよい。

40

【0099】

まず、本適用例に係る階層ストレージ制御装置 10A による動的階層制御について説明する。図 20 及び図 21 は、本適用例に係る階層ストレージ制御装置 10A による動的階層制御を説明するための図である。

図 20 は、本適用例に係る階層ストレージシステム 1A(図 22 参照)のワークロードの分析例を示す図であり、縦軸は下に向かってオフセットを示し、横軸は経過時間を示す。図 20 において、網掛けの領域 1 が高負荷領域を示す。階層ストレージ制御装置 10A は、図の矢印 2 が示すように、高負荷領域からある決めた範囲を拡張領域とする。

50

【 0 1 0 0 】

そして、階層ストレージ制御装置 1 0 A は、拡張領域と、拡張領域とつながる別の拡張領域とを合わせて 1 つの拡張領域とみなす。そして、階層ストレージ制御装置 1 0 A は、拡張領域を S S D 2 0 にデータを移動する移動領域とする。図 2 0 では、上下の破線の間の領域が移動領域である。

また、階層ストレージ制御装置 1 0 A は、ある時点で移動領域を決定すると、一定時間高負荷が発生しなくなるまで、移動領域を維持する。すなわち、階層ストレージ制御装置 1 0 A は、高負荷領域が消滅し、一定時間高負荷が発生しないと高負荷は消滅したとみなす。図 2 0 では、タイムアウトの矢印が高負荷が発生しない一定時間を示す。

【 0 1 0 1 】

10

図 2 1 は、階層ストレージシステム 1 A のワークロードの他の分析例を示す図であり、縦軸は上に向かってオフセットを示し、横軸は経過時間を示す。また、ボリュームは 1 G B 単位のセグメントに分割され、経過時間は 1 分間を単位としている。すなわち、図 2 1 において、網掛けの正方形領域 3 は、1 つのセグメントが 1 分間高負荷であったことを示す。また、s は高負荷領域から拡張領域として拡張するセグメント数を示し、図 2 1 では $s = 1$ である。

【 0 1 0 2 】

そして、階層ストレージ制御装置 1 0 A は、高負荷領域のセグメント間で距離が s 以内のものを繋ぎ合わせて n __セグメントを作成する。n __セグメントは、S S D 2 0 にデータが移動される移動領域であり、データの移動について一体制机制される。n __セグメントのセグメント数は $2 s + 1$ 以上となる。図 2 1 では、セグメント数が 5 の 2 つの n __セグメントが特定されている。

20

【 0 1 0 3 】

このように、階層ストレージ制御装置 1 0 A は、S S D 2 0 にデータを移動する領域として n __セグメントを特定することにより、高負荷領域の近傍として適切な領域を選択することができる。

次に、適用例に係る階層ストレージ制御装置 1 0 A の機能構成について説明する。図 2 2 は、適用例に係る階層ストレージシステム 1 A の構成例を示す図である。図 2 2 に示すように、階層ストレージ制御装置 1 0 A は、階層管理部 1 1 A , 階層ドライバ 1 2 , S S D ドライバ 1 3 , 及び H D D ドライバ 1 4 をそなえることができる。なお、以下の説明において、階層ストレージ制御装置 1 0 と同様の機能については、重複した説明を省略する。例えば、階層ドライバ 1 2 , S S D ドライバ 1 3 , 及び H D D ドライバ 1 4 は、図 7 に示す階層ストレージ制御装置 1 0 の構成と略同一である。また、図の簡略化のため、図 2 2 において階層ドライバ 1 2 がそなえる機能ブロックの図示を省略している。

30

【 0 1 0 4 】

以下、図 2 2 に示す階層ストレージ制御装置 1 0 A のうち、主に階層管理部 1 1 A の機能及び動作について、図 2 3 ~ 図 2 8 を参照してフローチャートに沿って説明する。

階層管理部 1 1 A は、H D D 3 0 についてトレースされた I O の情報に基づいて、S S D 2 0 にデータを移動する n __セグメントを決定し、決定した n __セグメントのデータの移動を階層ドライバ 1 2 に指示する。図 2 2 に示すように、階層管理部 1 1 A は、データ収集部 1 5 a , データベース 1 5 b , ワークロード分析部 1 5 c , 移動指示部 1 5 d , 及び分割数判定部 1 1 e をそなえる。なお、分割数判定部 1 1 e は、図 7 に示す階層ストレージ制御装置 1 0 の構成と略同一である。

40

【 0 1 0 5 】

はじめに、データ収集部 1 5 a の処理手順について説明する。図 2 3 は、適用例に係るデータ収集部 1 5 a によるデータ収集処理の動作例を示すフローチャートであり、図 2 4 は、図 2 2 に示すデータベース 1 5 b の一例を示す図である。なお、データ収集部 1 5 a は、b l k t r a c e コマンドを 6 0 秒間実行して、終了したことを条件として起動される。

【 0 1 0 6 】

50

図 2 3 に示すように、データ収集部 1 5 a は、b l k t r a c e コマンドの実行により得られたトレース結果を取り出し、1 G B オフセット単位すなわちセグメント単位で各セグメントの I O 数を抽出する (ステップ S 1 0 1)。

そして、データ収集部 1 5 a は、セグメントごとに I O 数が閾値 p を上回るかどうかを判定し、p を上回ったセグメントの抽出を行なう (ステップ S 1 0 2)。I O 数が閾値 p を上回ったセグメントは高負荷領域である。

【 0 1 0 7 】

そして、データ収集部 1 5 a は、抽出したセグメントに関して、隣接間距離が s 以内となったセグメントを繋ぎ合わせていく (ステップ S 1 0 3)。そして、データ収集部 1 5 a は、繋ぎ合わせたセグメントとその外側の s までの範囲のセグメントを n __セグメントと定義し、抽出順に n __セグメント番号を採番する (ステップ S 1 0 4)。

そして、データ収集部 1 5 a は、n __セグメントごとに n __セグメント番号、セグメント範囲、I O 数、平均レスポンスをタイムスタンプと共にデータベース 1 5 b に書き込む (ステップ S 1 0 5)。

【 0 1 0 8 】

ここで、データベース 1 5 b は、データ収集部 1 1 1 により特定された n __セグメントに関する情報を記憶する。図 2 4 に示すように、データベース 1 5 b は、n __セグメントごとに、n __セグメント番号、セグメント範囲、I O 数、平均レスポンス、及びタイムスタンプを対応させて記憶する。例えば、n __セグメント番号が“1”である n __セグメントは、先頭セグメントのオフセットが“3”であり、最終セグメントのオフセットが“5”であり、平均レスポンスが“0.6”(秒)であり、I O 数が“1000”であり、タイムスタンプが“1”である。

【 0 1 0 9 】

図 2 3 の説明に戻り、データ収集部 1 5 a は、図 7 に示すデータ収集部 1 1 a と同様に、全セグメントを対象にした合計 I O 数、平均レスポンスを集計し、タイムスタンプとともにデータベース 1 5 b に格納する (ステップ S 1 0 6)。なお、ステップ S 1 0 6 でデータベース 1 5 b に格納される情報は、図 2 4 における n __セグメント番号“all”のエントリに対応する。

【 0 1 1 0 】

以上により、データ収集部 1 5 a の処理が終了する。

このように、データ収集部 1 5 a は、高負荷のセグメントに関して隣接間距離が s 以内となったセグメントを繋ぎ合わせて n __セグメントを抽出することにより、高負荷のセグメントの近傍を適切に選択することができる。

次に、ワークロード分析部 1 5 c の処理手順について説明する。図 2 5 は、適用例に係るワークロード分析部 1 5 c による移動判定処理の動作例を示すフローチャートであり、図 2 6 及び図 2 7 は、それぞれ図 2 2 に示す候補テーブル 1 5 1 及び管理テーブル 1 5 2 の一例を示す図である。

【 0 1 1 1 】

図 2 5 に示すように、ワークロード分析部 1 5 c は、データベース 1 5 b から直近のタイムスタンプの n __セグメントについて I O 数を取り出し (ステップ S 1 1 1)、I O 数が多い順に n __セグメントを並べ替える (ステップ S 1 1 2)。

そして、ワークロード分析部 1 5 c は、各 n __セグメントの I O 数を合計することで i o __a l l を求める (ステップ S 1 1 3)。そして、ワークロード分析部 1 5 c は、以下の式 (1) の計算を m が m a x __s e g __n u m に到達するか、i o __r a t e が i o __r a t e __v a l u e を超えるまで行なう (ステップ S 1 1 4)。

【 0 1 1 2 】

【数 1】

$$\begin{aligned} \text{io_concentration} &= \sum_{k=1}^m \text{seg_sort}(k) \quad \dots (1) \\ \text{io_rate} &= (\text{io_concentration} * 100) / \text{io_all} \end{aligned}$$

【0113】

ここで、 max_seg_num は、同時にSSD20ヘデータの移動を行なう n __セグメント数である。また、 $\text{seg_sort}(k)$ は、 k 番目にアクセス数が多い n __セグメントのIO数である。 io_concentration は、トップ k 個の n __セグメントのIO数の合計であり、この数が多いほどトップ k 個の n __セグメントにアクセスが集中していることを示す。また、 io_all は、IO数を全 n __セグメントについて合計した総数であり、 io_rate は、トップ k 個の n __セグメントのIO数の合計の総数に対する割合を%表示したものである。従って、 io_rate の値が大きいほどトップ k 個の n __セグメントへのアクセスの集中度が高いことを示す。

10

【0114】

io_rate_value は、SSD20ヘデータを移動する候補としてトップ k 個の n __セグメントを選択するか否かの閾値である。

そして、 io_rate が io_rate_value を超えた場合には、ワークロード分析部15cは、以下のステップS115～ステップS122を行ない、 m が max_seg_num に到達した場合には、ステップS123に移動する。すなわち、 io_rate が io_rate_value を超えた場合には、ワークロード分析部15cは、対応する n __セグメント番号が連続して何回このトップ k に入ったのかを候補テーブル151に記録する(ステップS115)。

20

【0115】

ここで、候補テーブル151は、ワークロード分析部15cがそなえるテーブルであり、SSD20ヘデータを移動する候補を記憶する。図26に示すように、候補テーブル151は、 n __セグメントごとに、 n __セグメント番号、先頭セグメント番号、セグメント数、及び連続数を対応させて記憶する。ここで、先頭セグメント番号は、 n __セグメントの先頭セグメントのオフセットである。セグメント数は、 n __セグメントに含まれるセグメントの個数である。連続数は、候補として連続して候補テーブル151に登録された回数を示す。

30

【0116】

図25の説明に戻り、ワークロード分析部15cは、前タイムスライスでトップ k に入った n __セグメントで今回トップ k から外れた n __セグメントは、連続数をリセットする(ステップS116)。

そして、ワークロード分析部15cは、連続数が所定の閾値 t_1 を超えた n __セグメントを移動候補として抽出し、抽出した n __セグメントに含まれるセグメント数を n とし、 n __セグメントのデータの移動時間 Tiering_time を計算する(ステップS117)。

【0117】

40

ここで、 $\text{Tiering_time} = \text{seg_move_time} \times n + \text{検出遅延}$ であり、 seg_move_time は、1セグメントのデータをHDD30からSSD20へ移動するのにかかる時間である。また、検出遅延は、移動候補の検出にかかる時間であり、ここではデータの収集間隔の60秒とする。

そして、ワークロード分析部15cは、 Tiering_time とIOの集中度が高い状態が続くと期待される時間 Life_ex_time (平均余命時間)とを比較する(ステップS118)。 Tiering_time が Life_ex_time 以上の場合(ステップS118のNoルート)、処理がステップS121に移行する。一方、 Tiering_time が Life_ex_time より小さい場合(ステップS118のYesルート)、ワークロード分析部15cは、移動候補 n __セグメントの情報を移動指

50

示部 15d へ通知し、移動候補 n __セグメントのデータの HDD30 から SSD20 への移動を指示する (ステップ S119)。また、ワークロード分析部 15c は、SSD20 へのデータの移動を指示した n __セグメントの情報を管理テーブル 152 に記録する (ステップ S120)。

【0118】

ここで、管理テーブル 152 は、ワークロード分析部 15c がそなえるテーブルであり、SSD20 へデータを移動する対象として選択した n __セグメントを記憶する。図 27 に示すように、管理テーブル 152 は、 n __セグメントごとに、 n __セグメント番号、先頭セグメント番号、セグメント数、及び連続数を対応させて記憶する。ここで、連続数は、トップ k 個の候補が選択された場合に、連続して候補として選択されなかった回数を示す。

10

【0119】

図 25 の説明に戻り、ワークロード分析部 15c は、トップ k に入った n __セグメント番号と管理テーブル 152 に登録されている n __セグメント番号の突合せを行なう。また、ワークロード分析部 15c は、管理テーブル 152 に登録されている n __セグメントごとにトップ k に入らなかった n __セグメント番号の連続数を “+1” し、トップ k に入っていたら連続数を “0” にリセットする (ステップ S121)。

【0120】

次いで、ワークロード分析部 15c は、管理テーブル 152 に登録されている n __セグメント番号ごとに連続数が所定の閾値 t_2 を超えているか否かの判断を行なう。連続数が所定の閾値 t_2 を超えている場合、ワークロード分析部 15c は、 n __セグメント番号を移動指示部 15d に通知して SSD20 から HDD30 へのデータの移動を指示する。また、ワークロード分析部 15c は、管理テーブル 152 に登録されている n __セグメントの情報を削除する (ステップ S122)。そして、ワークロード分析部 15c は、60 秒スリープし (ステップ S123)、処理がステップ S111 に移行する。

20

【0121】

このように、ワークロード分析部 15c が I/O の集中度が高い n __セグメントのデータを HDD30 から SSD20 に移動するように移動指示部 15d に指示することによって、ユーザは HDD30 のデータに高速にアクセスすることができる。

以上のように、ワークロード分析部 15c は、データ収集部 15a により集計された入出力数が第 1 の閾値より大きな単位領域と所定の距離内にある単位領域を繋ぎ合わせた拡張領域と、当該拡張領域と繋がる他の拡張領域と、を合わせた移動領域を特定する特定部の一例であるといえる。

30

【0122】

次に、移動指示部 15d の処理手順について説明する。図 28 は、適用例に係る移動指示部 15d による移動指示通知処理の動作例を示すフローチャートである。

図 28 に示すように、移動指示部 15d は、ワークロード分析部 15c からの移動指示を待つ (ステップ S131)。移動指示があると、移動指示部 15d は、 n __セグメント番号に属する各セグメントのボリューム上のオフセットを HDD30 上のオフセットに変換する (ステップ S132)。

40

【0123】

そして、移動指示部 15d は、セグメントごとに、セグメント番号に対応する HDD30 上のオフセットと、データの移動方向を階層ドライバ 12 に通知する (ステップ S133)。ここで、データの移動方向は、HDD30 から SSD20 か、又は、SSD20 から HDD30 である。

また、移動指示部 15d は、移動指示を出したセグメント数と移動判断を行なったデータのタイムスタンプ (直近のタイムスタンプ) とを分割数判定部 11e に通知し (ステップ S134)、処理がステップ S131 に移行する。

【0124】

このように、移動指示部 15d が各セグメントのボリューム上のオフセットを HDD30

50

0上のオフセットに変換することによって、階層ドライバ12はSSD20とHDD30との間でデータを移動することができる。

以上のように、本適用例に係る階層ストレージ制御装置10Aによれば、データベース15bが、I/O数が閾値pを上回るセグメントに関して、隣接間距離s以内のセグメントを繋ぎ合わせる。そして、データ収集部15aは、繋ぎ合わせたセグメントとその外側のsまでの範囲をn__セグメントとして抽出する。また、ワークロード分析部15cが、n__セグメントを単位として、HDD30からSSD20にデータを移動する対象を決定する。

【0125】

このとき、本適用例に係る分割数判定部11e及び分割部12dは、移動単位のn__セグメントに属する複数のセグメントの各々について、動的に決定した分割数に分割して階層移動を行なうことができる。例えば、分割部12dは、n__セグメントの移動領域に記憶されたデータをSSD20に移動する移動処理において、移動領域に含まれる複数の単位領域の各々を所定の分割数により複数の分割領域に分割する。そして、分割部12dは、移動領域に記憶されたデータを分割領域の単位でSSD20に移動する。

【0126】

従って、階層ストレージシステム1Aは、高負荷領域の近傍を適切に選択した上で、使用する機器の性能やワークロードに応じた最適な移動単位で、HDD30からSSD20にデータを移動することができ、HDD30へのアクセスを高速化することができる。

〔2〕その他

以上、本発明の好ましい実施形態について詳述したが、本発明は、係る特定の実施形態に限定されるものではなく、本発明の趣旨を逸脱しない範囲内において、種々の変形、変更して実施することができる。

【0127】

例えば、一実施形態において、SSD20及びHDD30を用いた階層ストレージシステム1及び1Aについて説明したが、これに限定されるものではなく、例えばキャッシュメモリと主記憶装置とを用いた階層記憶システムにも同様に適用することができる。すなわち、本発明は、不揮発性記憶装置の階層記憶システムだけでなく、揮発性記憶装置を含む階層記憶システムにも同様に適用することができる。

【0128】

また、一実施形態に係る階層ストレージシステム1及び1Aは、SSD20及びHDD30の他に、速度差のある記憶装置にも適用することが可能である。例えばHDDと、HDDよりも大容量だが低速なテープドライブ等の磁気記録装置とを用いた階層ストレージシステム等にも適用することが可能である。

さらに、一実施形態において、階層ストレージ制御装置10及び10Aの動作を1つのSSD20及び1つのHDD30に着目して説明したが、複数のSSD20及び複数のHDD30が階層ストレージシステム1及び1Aにそなえられる場合も同様である。

【0129】

〔3〕付記

以上の実施形態に関し、更に以下の付記を開示する。

（付記1）

第1の記憶装置の記憶領域を所定の大きさを分割した複数の単位領域について、入力された要求に対する応答性能を監視する監視部と、

前記第1の記憶装置の移動対象の単位領域に記憶されたデータを前記第1の記憶装置とは異なる性能の第2の記憶装置に移動する移動処理において、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第2の記憶装置に移動する分割部と、

前記監視部が監視した前記移動処理実行中の第1の応答性能に基づいて、前記所定の分割数を変更する変更部と、

をそなえることを特徴とする、ストレージ制御装置。

【 0 1 3 0 】

(付 記 2)

前記変更部は、前記第 1 の応答性能と、前記移動処理実行前の第 2 の応答性能とに基づいて、前記所定の分割数を変更することを特徴とする、付記 1 記載のストレージ制御装置。

(付 記 3)

前記所定の分割数は、前記移動処理実行前の第 2 の応答性能を基に求められることを特徴とする、付記 2 記載のストレージ制御装置。

【 0 1 3 1 】

(付 記 4)

前記変更部は、前記第 1 の応答性能と前記第 2 の応答性能とを比較し、前記第 1 の応答性能が前記第 2 の応答性能よりも劣化していると判断した場合に前記所定の分割数を増加させる一方、前記第 1 の応答性能が前記第 2 の応答性能よりも優れていると判断した場合に前記所定の分割数を減少させる、ことを特徴とする、付記 1 ～ 3 のいずれか 1 項記載のストレージ制御装置。

【 0 1 3 2 】

(付 記 5)

前記変更部は、前記監視部が前記移動処理実行中の複数の時点で監視した複数の第 1 の応答性能を取得し、前記第 2 の応答性能と前記複数の応答性能の平均値とを比較して、前記平均値が前記第 2 の応答性能よりも劣化しているか否かに応じて前記所定の分割数を増減させる、ことを特徴とする、付記 4 記載のストレージ制御装置。

【 0 1 3 3 】

(付 記 6)

前記監視部は、前記複数の単位領域について単位領域ごとに入出力数を集計し、前記ストレージ制御装置は、

前記監視部により集計された入出力数が第 1 の閾値より大きな単位領域と所定の距離内にある単位領域を繋ぎ合わせた拡張領域と、該拡張領域と繋がる他の拡張領域と、を合わせた移動領域を特定する特定部、をさらにそなえ、

前記分割部は、前記移動領域に記憶されたデータを前記第 2 の記憶装置に移動する移動処理において、前記移動領域に含まれる複数の単位領域の各々を所定の分割数により複数の分割領域に分割し、前記移動領域に記憶されたデータを前記分割領域の単位で前記第 2 の記憶装置に移動する、ことを特徴とする、付記 1 ～ 5 のいずれか 1 項記載のストレージ制御装置。

【 0 1 3 4 】

(付 記 7)

第 1 の記憶装置及び第 2 の記憶装置の制御を行なうコンピュータに、

前記第 1 の記憶装置の記憶領域を所定の大きさを分割した複数の単位領域について、入力された要求に対する応答性能を監視し、

前記第 1 の記憶装置の移動対象の単位領域に記憶されたデータを前記第 1 の記憶装置とは異なる性能の第 2 の記憶装置に移動する移動処理を行ない、

前記移動処理において、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第 2 の記憶装置に移動し、

前記監視により監視した前記移動処理実行中の第 1 の応答性能に基づいて、前記所定の分割数を変更する、処理を実行させることを特徴とする、制御プログラム。

【 0 1 3 5 】

(付 記 8)

前記変更する処理において、前記第 1 の応答性能と、前記移動処理実行前の第 2 の応答性能とに基づいて、前記所定の分割数を変更する、

ことを特徴とする、付記 7 記載の制御プログラム。

(付記 9)

前記所定の分割数は、前記移動処理実行前の第 2 の応答性能を基に求められる、ことを特徴とする、付記 8 記載の制御プログラム。

【0136】

(付記 10)

前記変更する処理において、前記第 1 の応答性能と前記第 2 の応答性能とを比較し、前記第 1 の応答性能が前記第 2 の応答性能よりも劣化していると判断した場合に前記所定の分割数を増加させる一方、前記第 1 の応答性能が前記第 2 の応答性能よりも優れていると判断した場合に前記所定の分割数を減少させる、

10

ことを特徴とする、付記 7 ~ 9 のいずれか 1 項記載の制御プログラム。

【0137】

(付記 11)

前記変更する処理において、前記監視により前記移動処理実行中の複数の時点で監視した複数の第 1 の応答性能を取得し、前記第 2 の応答性能と前記複数の応答性能の平均値とを比較して、前記平均値が前記第 2 の応答性能よりも劣化しているか否かに応じて前記所定の分割数を増減させる、

ことを特徴とする、付記 10 記載の制御プログラム。

【0138】

(付記 12)

20

前記コンピュータに、

前記複数の単位領域について単位領域ごとに入出力数を集計し、

前記集計された入出力数が第 1 の閾値より大きな単位領域と所定の距離内にある単位領域を繋ぎ合わせた拡張領域と、該拡張領域と繋がる他の拡張領域と、を合わせた移動領域を特定し、

前記移動領域に記憶されたデータを前記第 2 の記憶装置に移動する移動処理において、前記移動領域に含まれる複数の単位領域の各々を所定の分割数により複数の分割領域に分割し、前記移動領域に記憶されたデータを前記分割領域の単位で前記第 2 の記憶装置に移動する、

処理を実行させることを特徴とする、付記 7 ~ 11 のいずれか 1 項記載の制御プログラム

30

【0139】

(付記 13)

第 1 の記憶装置及び第 2 の記憶装置の制御を行なうストレージ制御装置における制御方法であって、

前記第 1 の記憶装置の記憶領域を所定の大きさを分割した複数の単位領域について、入力された要求に対する応答性能を監視し、

前記第 1 の記憶装置の移動対象の単位領域に記憶されたデータを前記第 1 の記憶装置とは異なる性能の第 2 の記憶装置に移動する移動処理を行ない、

前記移動処理において、前記移動対象の単位領域を所定の分割数により複数の分割領域に分割し、前記データを前記分割領域の単位で前記第 2 の記憶装置に移動し、

40

前記監視により監視した前記移動処理実行中の第 1 の応答性能に基づいて、前記所定の分割数を変更する、

ことを特徴とする、制御方法。

【0140】

(付記 14)

前記変更する処理において、前記第 1 の応答性能と、前記移動処理実行前の第 2 の応答性能とに基づいて、前記所定の分割数を変更する、

ことを特徴とする、付記 13 記載の制御方法。

(付記 15)

50

前記所定の分割数は、前記移動処理実行前の第 2 の応答性能を基に求められる、
ことを特徴とする、付記 1 4 記載の制御方法。

【 0 1 4 1 】

(付記 1 6)

前記変更する処理において、前記第 1 の応答性能と前記第 2 の応答性能とを比較し、前記第 1 の応答性能が前記第 2 の応答性能よりも劣化していると判断した場合に前記所定の分割数を増加させる一方、前記第 1 の応答性能が前記第 2 の応答性能よりも優れていると判断した場合に前記所定の分割数を減少させる、
ことを特徴とする、付記 1 3 ~ 1 5 のいずれか 1 項記載の制御方法。

【 0 1 4 2 】

10

(付記 1 7)

前記変更する処理において、前記監視により前記移動処理実行中の複数の時点で監視した複数の第 1 の応答性能を取得し、前記第 2 の応答性能と前記複数の応答性能の平均値とを比較して、前記平均値が前記第 2 の応答性能よりも劣化しているか否かに応じて前記所定の分割数を増減させる、
ことを特徴とする、付記 1 6 記載の制御方法。

【 0 1 4 3 】

(付記 1 8)

前記複数の単位領域について単位領域ごとに入出力数を集計し、
前記集計された入出力数が第 1 の閾値より大きな単位領域と所定の距離内にある単位領域を繋ぎ合わせた拡張領域と、該拡張領域と繋がる他の拡張領域と、を合わせた移動領域を特定し、

20

前記移動領域に記憶されたデータを前記第 2 の記憶装置に移動する移動処理において、前記移動領域に含まれる複数の単位領域の各々を所定の分割数により複数の分割領域に分割し、前記移動領域に記憶されたデータを前記分割領域の単位で前記第 2 の記憶装置に移動する、
ことを特徴とする、付記 1 3 ~ 1 7 のいずれか 1 項記載の制御方法。

【 符号の説明 】

【 0 1 4 4 】

- 1 , 1 A 階層ストレージシステム (ストレージ装置)
- 1 0 , 1 0 A 階層ストレージ制御装置 (ストレージ制御装置)
- 1 0 a C P U
- 1 0 b メモリ
- 1 0 c 記憶部
- 1 0 d インタフェース部
- 1 0 e 入出力部
- 1 0 f , 1 0 h 記録媒体
- 1 0 g 読取部
- 1 1 , 1 1 A 階層管理部
- 1 1 a , 1 5 a データ収集部
- 1 1 b , 1 5 b データベース
- 1 1 c , 1 5 c ワークロード分析部
- 1 1 d , 1 5 d 移動指示部
- 1 1 e 分割数判定部
- 1 2 , 1 1 0 階層ドライバ
- 1 2 a I O マップ部
- 1 2 b ペンディングキュー
- 1 2 c 階層テーブル
- 1 2 d 分割部
- 1 3 , 1 2 0 S S D ドライバ

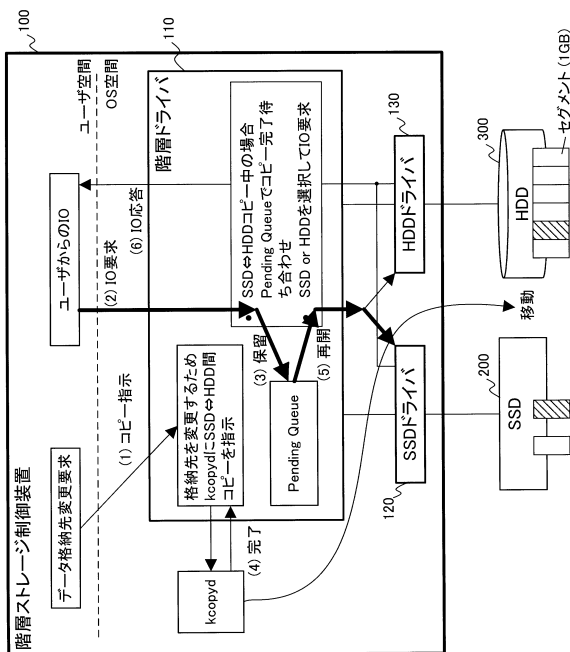
30

40

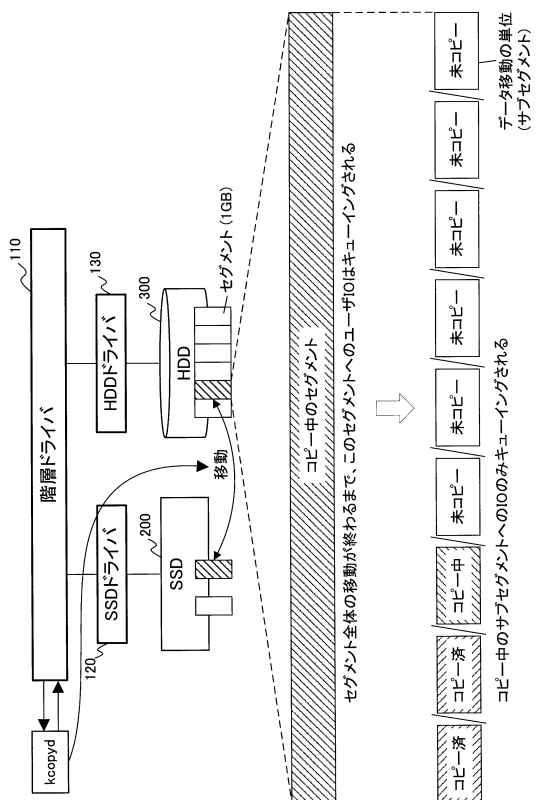
50

- 1 4 , 1 3 0 H D D ドライバ
- 1 5 1 候補テーブル
- 1 5 2 管理テーブル
- 2 0 , 2 0 0 S S D
- 3 0 , 3 0 0 H D D
- 1 0 0 階層ストレージ制御装置

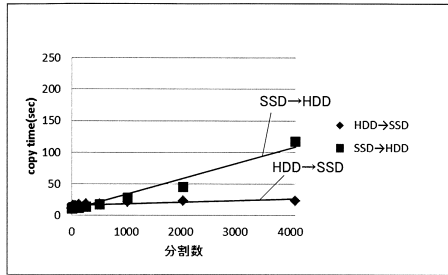
【 図 1 】



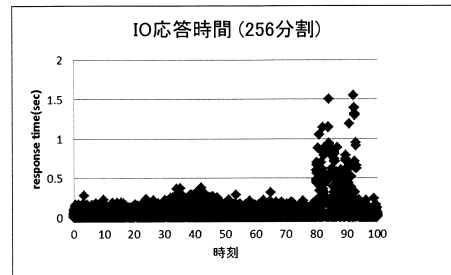
【 図 2 】



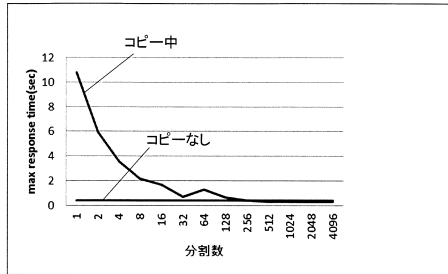
【図 3】



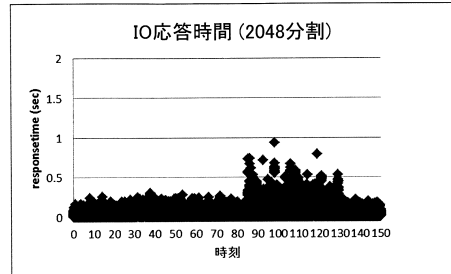
【図 5】



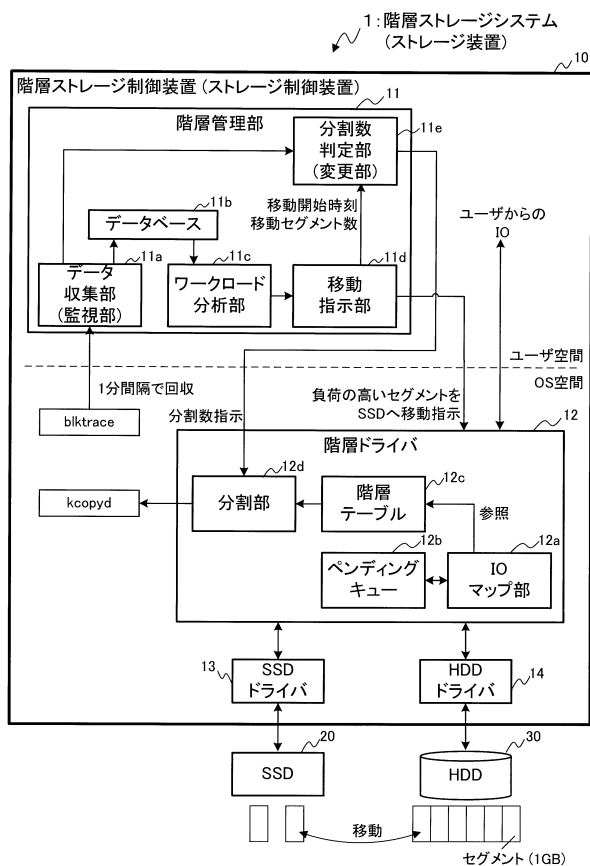
【図 4】



【図 6】



【図 7】



【図 8】

データベース ^{11b}

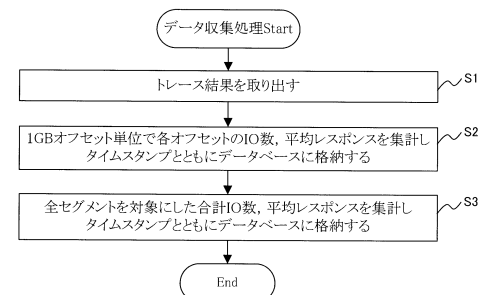
セグメント	IO数	平均レスポンス(s)	タイムスタンプ
0	1000	0.6	1
1	1200	0.4	3
⋮	⋮	⋮	⋮
all	50000	0.8	10
all	55000	0.9	11

【図 9】

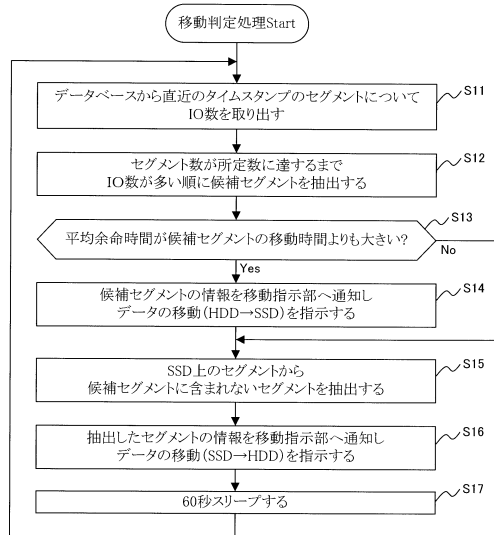
階層テーブル ^{12c}

SSDオフセット	HDDオフセット	状態
0	268435456	allocated
2097152	306184192	Moving(HDD→SSD)
4194304	505413632	Moving(SSD→HDD)
6291456	NULL	free
⋮	⋮	⋮

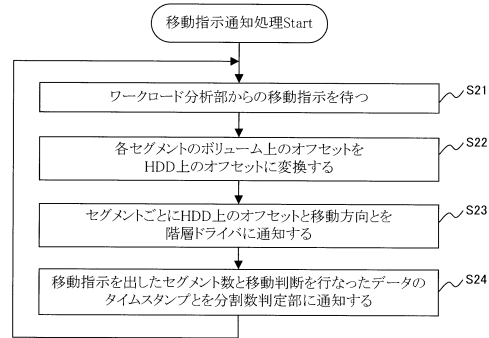
【図 10】



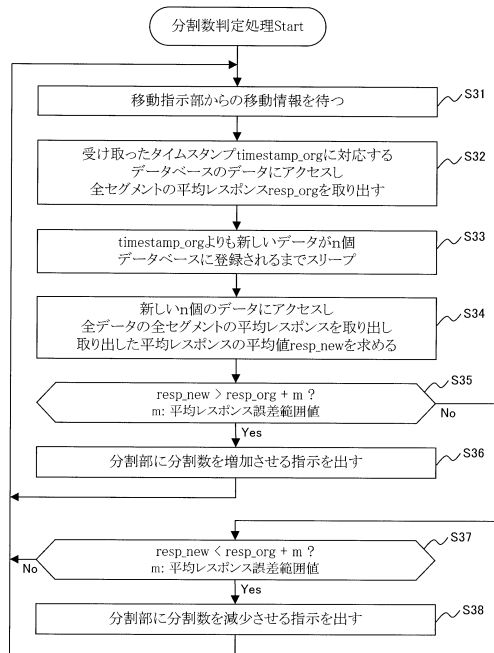
【図 1 1】



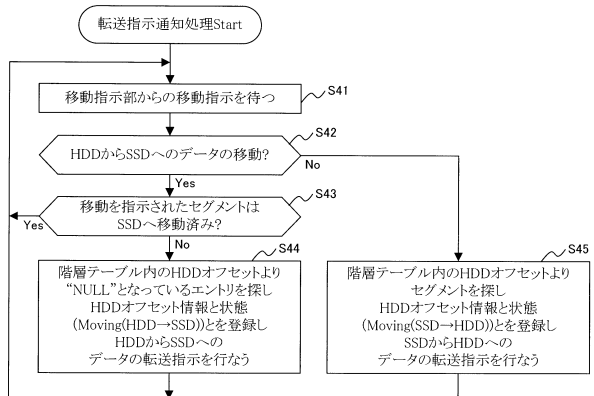
【図 1 2】



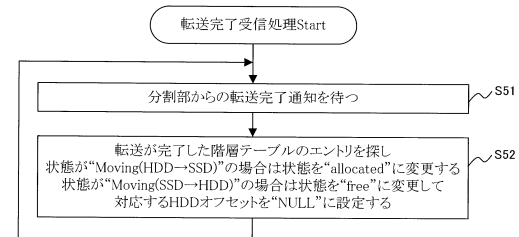
【図 1 3】



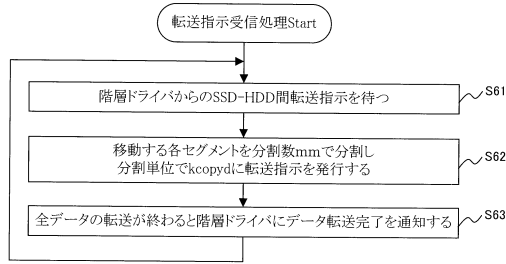
【図 1 4】



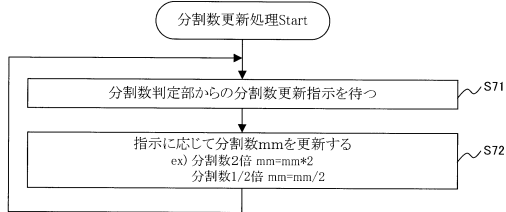
【図 1 5】



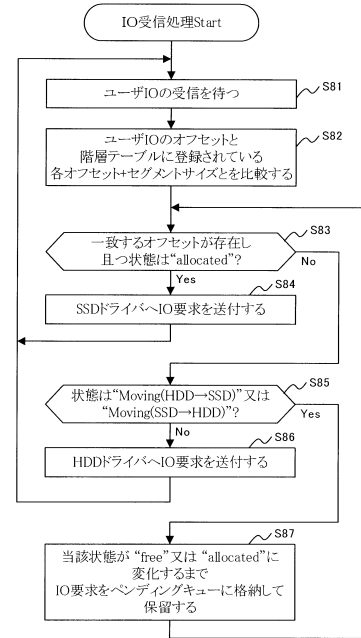
【図 16】



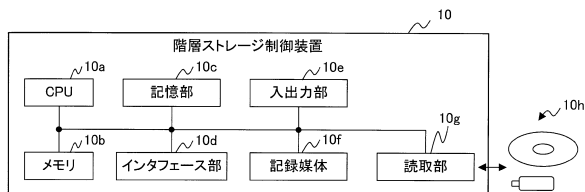
【図 17】



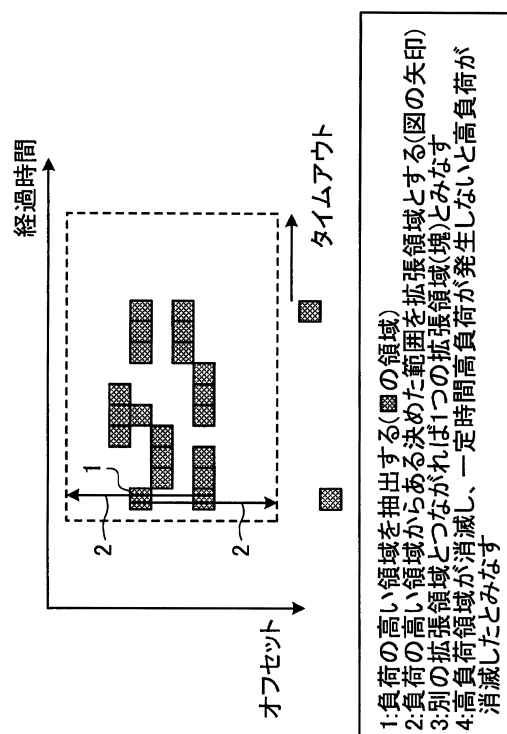
【図 18】



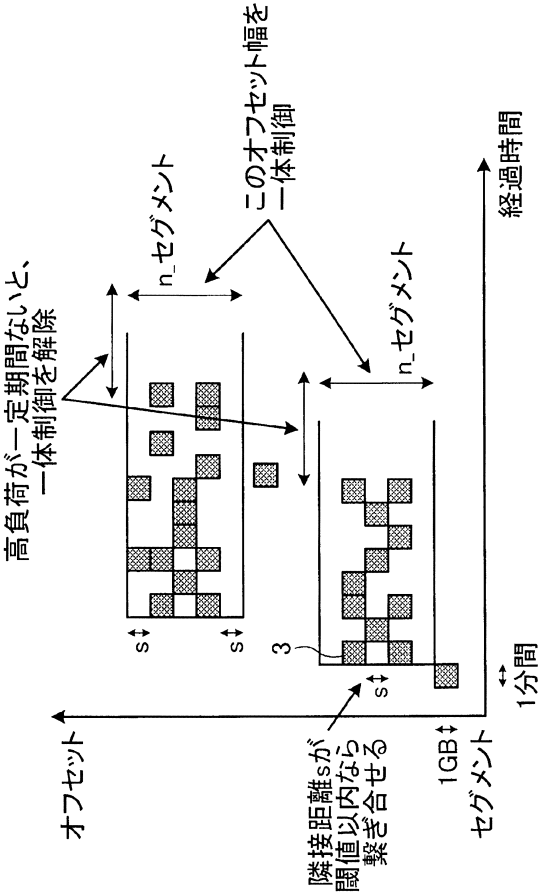
【図 19】



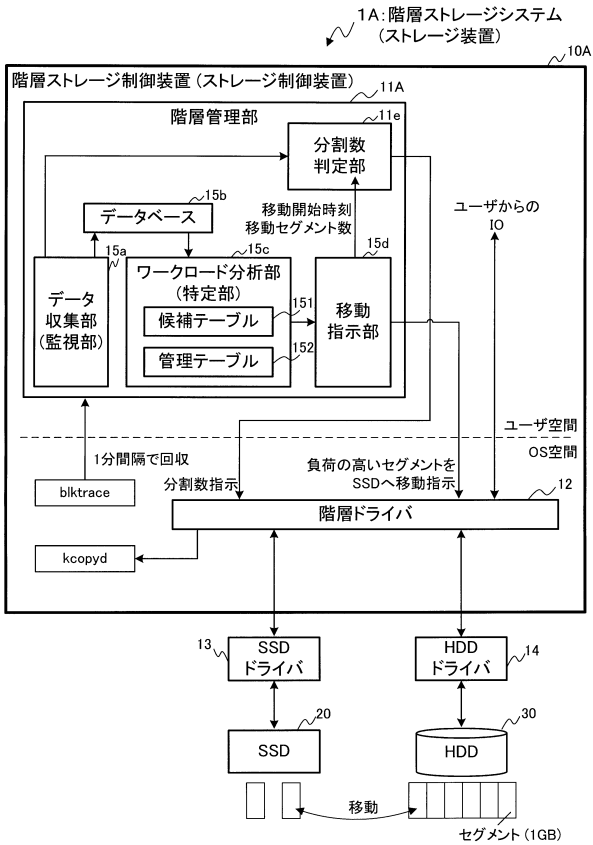
【図 20】



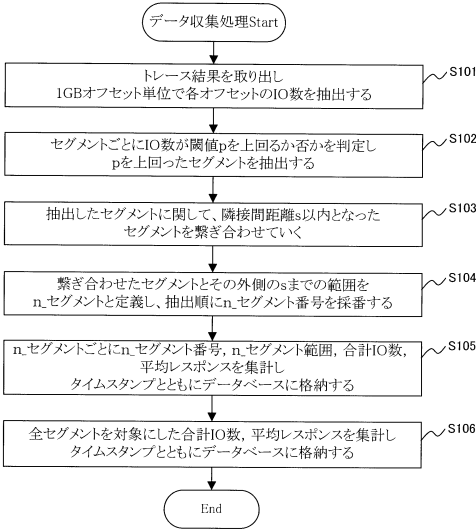
【図 2 1】



【図 2 2】



【図 2 3】

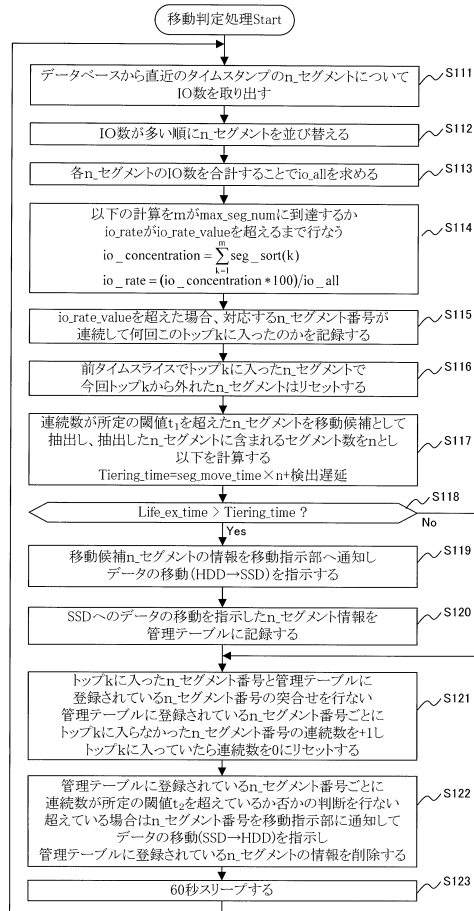


【図 2 4】

データベース

n_セグメント番号	セグメント範囲	IO数	平均レスポンス(s)	タイムスタンプ
0	3	5	1000	1
...
all	1	xx	50000	10
all	1	xx	55000	11

【図 25】



【図 26】

候補テーブル¹⁵¹

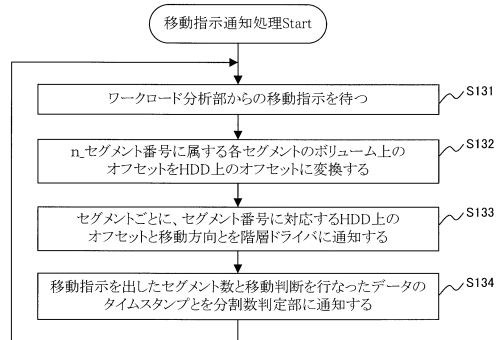
n_セグメント番号	先頭セグメント番号	セグメント数	連続数
0	10	5	1
:	:	:	:

【図 27】

管理テーブル¹⁵²

n_セグメント番号	先頭セグメント番号	セグメント数	連続数
0	10	5	2
:	:	:	:

【図 28】



フロントページの続き

(72)発明者 河場 基行

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 田名網 忠雄

(56)参考文献 特開2006-113882(JP,A)

特開2010-146450(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