

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2017-204068

(P2017-204068A)

(43) 公開日 平成29年11月16日(2017.11.16)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 12/0804 (2016.01)	G06F 12/08 501C	5B005
G06F 12/12 (2016.01)	G06F 12/12 551	5B205
G06F 12/08 (2016.01)	G06F 12/08 543B	

審査請求 未請求 請求項の数 4 O L (全 18 頁)

(21) 出願番号 特願2016-94628 (P2016-94628)
 (22) 出願日 平成28年5月10日 (2016.5.10)

(71) 出願人 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 100104190
 弁理士 酒井 昭徳
 (72) 発明者 加藤 純
 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
 Fターム(参考) 5B005 JJ13 MM01 PP03 QQ02 VV03
 5B205 JJ13 MM01 NN89 PP03 QQ02
 QQ11 VV02 VV03

(54) 【発明の名称】 情報処理装置、キャッシュメモリ制御方法、およびキャッシュメモリ制御プログラム

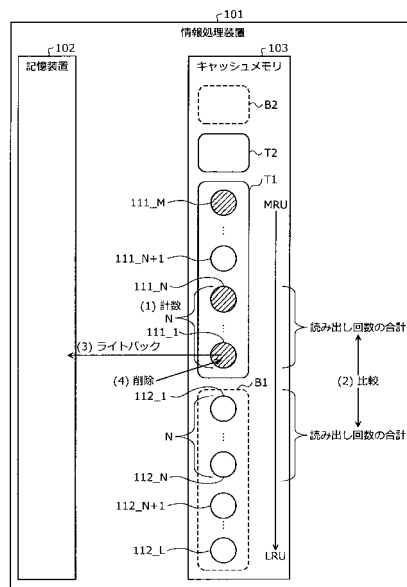
(57) 【要約】

【課題】 キャッシュメモリ上の更新されたデータをライトバックする適切なタイミングを決定すること。

【解決手段】 情報処理装置101は、図1の(1)で示すように、リストT1の削除対象となるエントリ情報111_1がDirtyページである場合、LRUに従った順序におけるエントリ情報111_1からDirtyページが連続する数を計数する。次に、情報処理装置101は、図1の(2)で示すように、第1の合計として、エントリ情報111_1からN分のエントリ情報への読み出し回数の合計と、第2の合計として、エントリ情報112_1からN分のエントリ情報への読み出し回数の合計とを比較する。図1の例では、第1の合計よりも第2の合計が多いものとする。この場合、情報処理装置101は、図1の(3)、(4)で示すように、エントリ情報111_1のDirtyページをライトバックし、エントリ情報111_1のページを削除する。

【選択図】 図1

本実施の形態にかかる情報処理装置101の動作例を示す説明図



【特許請求の範囲】

【請求項 1】

複数のデータを記憶する記憶装置と、

前記複数のデータのうちのいずれかのデータと前記いずれかのデータの前の記憶装置上における位置を示す情報とを含むエントリ情報を有する第 1 のリストと、前記第 1 のリストから追い出されたエントリ情報に含まれた削除済みのデータの前の記憶装置上における位置を示す情報を含むエントリ情報を有する第 2 のリストとを記憶するキャッシュメモリと、

前記第 1 のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、前記優先度に従った順序における前記削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、前記削除対象のエントリ情報から前記数分のエントリ情報へのアクセス回数の第 1 の合計よりも前記優先度に従った順序における前記第 2 のリストに最も後に追加されたエントリ情報から前記数分のエントリ情報へのアクセス回数の第 2 の合計が多い場合、前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを前記キャッシュメモリから削除する制御部と、

を有することを特徴とする情報処理装置。

【請求項 2】

前記制御部は、

前記第 1 の合計よりも前記第 2 の合計が多い場合、前記第 1 のリストに含まれるデータが更新されているエントリ情報の割合と比較される閾値であって前記閾値の方が小さければ前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを削除する前記閾値を、現在の値より小さく設定し、

前記第 1 の合計が前記第 2 の合計よりも多い場合、前記閾値を現在の値より大きく設定する、

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】

コンピュータが、

記憶装置が記憶する複数のデータのうちのいずれかのデータと前記いずれかのデータの前の記憶装置上における位置を示す情報とを含むエントリ情報を有する第 1 のリストおよび前記第 1 のリストから追い出されたエントリ情報に含まれた削除済みのデータの前の記憶装置上における位置を示す情報を含むエントリ情報を有する第 2 のリストを記憶するキャッシュメモリの前記第 1 のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、前記優先度に従った順序における前記削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、

前記削除対象のエントリ情報から前記数分のエントリ情報へのアクセス回数の第 1 の合計よりも前記優先度に従った順序における前記第 2 のリストに最も後に追加されたエントリ情報から前記数分のエントリ情報へのアクセス回数の第 2 の合計が多い場合、前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを前記キャッシュメモリから削除する、

処理を実行することを特徴とするキャッシュメモリ制御方法。

【請求項 4】

コンピュータに、

記憶装置が記憶する複数のデータのうちのいずれかのデータと前記いずれかのデータの前の記憶装置上における位置を示す情報とを含むエントリ情報を有する第 1 のリストおよび前記第 1 のリストから追い出されたエントリ情報に含まれた削除済みのデータの前の記憶装置上における位置を示す情報を含むエントリ情報を有する第 2 のリストを記憶するキャッシュメモリの前記第 1 のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、前記優先度に

従った順序における前記削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、

前記削除対象のエントリ情報から前記数分のエントリ情報へのアクセス回数の第1の合計よりも前記優先度に従った順序における前記第2のリストに最も後に追加されたエントリ情報から前記数分のエントリ情報へのアクセス回数の第2の合計が多い場合、前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを前記キャッシュメモリから削除する、

処理を実行させることを特徴とするキャッシュメモリ制御プログラム。

【発明の詳細な説明】

【技術分野】

10

【0001】

本発明は、情報処理装置、キャッシュメモリ制御方法、およびキャッシュメモリ制御プログラムに関する。

【背景技術】

【0002】

従来、記憶装置より高速にアクセス可能なキャッシュメモリを用いて、記憶装置へのアクセス性能を向上させる技術がある。また、記憶装置への書き込み要求を受け付けた場合、キャッシュメモリにデータを書き込んでおき記憶装置には書き込まず、CPU (Central Processing Unit) の空き時間等に、キャッシュメモリ上の更新されたデータを記憶装置に書き込む、いわゆるライトバックと呼ばれる動作がある。また、キャッシュメモリに新たなデータを追加できない場合、今後最も必要とされないデータを予測し、予測したデータを削除して、代わりに新たなデータを追加する、キャッシュ置換方式と呼ばれる技術がある。

20

【0003】

関連する先行技術として、例えば、二次データストレージキャッシュの現在の状態のキャッシング効率に従ってヒートメトリック閾値を調整し、そのヒートメトリックが閾値を下回る、二次データストレージキャッシュに提供された候補データを拒否するものがある。また、ディスク上の連続する複数のブロックを1つのグループとし、キャッシュメモリ上でのみ更新されているブロックを含むグループがあるときは、そのグループ内の連続する複数のブロックを1回のアクセスによりディスク上に書き戻す技術がある。

30

【先行技術文献】

【特許文献】

【0004】

【特許文献1】特表2014-535106号公報

【特許文献2】特開平05-303528号公報

【発明の概要】

【発明が解決しようとする課題】

【0005】

しかしながら、従来技術によれば、キャッシュメモリ上の更新されたデータを、いつライトバックすればよいか決定することが困難である。例えば、ライトバックを行う間隔を長くすると、キャッシュ置換方式によって最も必要とされないと予測された更新されたデータの代わりに、更新されていないデータが削除されていき、キャッシュメモリのヒット率が低下する可能性がある。

40

【0006】

1つの側面では、本発明は、キャッシュメモリ上の更新されたデータをライトバックする適切なタイミングを決定することができる情報処理装置、キャッシュメモリ制御方法、およびキャッシュメモリ制御プログラムを提供することを目的とする。

【課題を解決するための手段】

【0007】

本発明の一側面によれば、記憶装置が記憶する複数のデータのうちのいずれかのデータ

50

といずれかのデータの記憶装置上における位置を示す情報とを含むエントリ情報を有する第1のリストと、第1のリストから追い出されたエントリ情報に含まれた削除済みのデータの記憶装置上における位置を示す情報を含むエントリ情報を有する第2のリストとを記憶するキャッシュメモリの第1のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、優先度に従った順序における削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、削除対象のエントリ情報から数分のエントリ情報へのアクセス回数の第1の合計よりも優先度に従った順序における第2のリストに最も後に追加されたエントリ情報から数分のエントリ情報へのアクセス回数の第2の合計が多い場合、削除対象のエントリ情報のデータを記憶装置に書き出して削除対象のエントリ情報のデータをキャッシュメモリから削除する情報処理装置、キャッシュメモリ制御方法、およびキャッシュメモリ制御プログラムが提案される。

10

【発明の効果】

【0008】

本発明の一態様によれば、キャッシュメモリ上の更新されたデータをライトバックする適切なタイミングを決定することができるという効果を奏する。

【図面の簡単な説明】

【0009】

【図1】図1は、本実施の形態にかかる情報処理装置101の動作例を示す説明図である。

20

【図2】図2は、ARCの動作例を示す説明図である。

【図3】図3は、ディスクアレイ装置300のハードウェア構成例を示す説明図である。

【図4】図4は、CM311のハードウェア構成例を示す説明図である。

【図5】図5は、CM311の機能構成例を示す説明図である。

【図6】図6は、読み出し時の動作例を示す説明図である。

【図7】図7は、書き込み時の動作例を示す説明図である。

【図8】図8は、ディスク313への書き込みの動作例を示す説明図である。

【図9】図9は、読み出し処理手順の一例を示すフローチャートである。

【図10】図10は、Watermark調整部505の処理手順の一例を示すフローチャートである。

30

【図11】図11は、キャッシュデータ管理部503の処理手順の一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】

【0010】

以下に図面を参照して、開示の情報処理装置、キャッシュメモリ制御方法、およびキャッシュメモリ制御プログラムの実施の形態を詳細に説明する。

【0011】

図1は、本実施の形態にかかる情報処理装置101の動作例を示す説明図である。情報処理装置101は、キャッシュメモリを制御するコンピュータである。情報処理装置101は、複数のデータを記憶する記憶装置102と、キャッシュメモリ103とにアクセス可能である。情報処理装置101は、例えば、ディスクアレイ装置、サーバ、PC(Personal Computer)、タブレット端末、携帯端末、携帯電話等である。キャッシュメモリ103は、記憶装置102よりアクセス性能が高い記憶装置である。キャッシュメモリ103は、記憶装置102のデータを一時的に記憶する。情報処理装置101内のCPUがキャッシュメモリを用いてデータにアクセスすることにより、記憶装置102とCPUとの性能差を埋めることができる。

40

【0012】

また、キャッシュメモリは、CPUとRAM(Random Access Memory)との間に設置され、RAMの一部のデータを記憶するが、このような使用形態に限られない。例えば、RAMをキャッシュメモリとみなし、RAMは、RAMより低速な記

50

憶装置のデータを一時的に記憶してもよい。RAMより低速な記憶装置としては、例えば、SSD (Solid State Drive)、HDD (Hard Disk Drive)、光ディスクドライブ、磁気テープを有するテープドライブ等である。

【0013】

また、キャッシュメモリは、記憶装置をページと呼ばれる一定のサイズのデータに分割して管理を行う。そして、キャッシュメモリには、複数のエントリ情報があり、複数のエントリ情報の各エントリ情報は、複数のページのいずれかのページと、いずれかのページのメタデータとを記憶する。メタデータには、いずれかのページの記憶装置上における位置を示す情報が含まれる。いずれかのページの記憶装置上における位置を示す情報は、例えば、いずれかのページが格納されるLUN (Logical Unit Number)、LBA (Logical Block Addressing)といった情報である。

10

【0014】

また、キャッシュメモリのエントリ情報が新たに追加できない場合、キャッシュ置換方式による優先度に基づいて、削除対象のエントリ情報を決定する。具体的には、キャッシュ置換方式は、今後最も必要とされないエントリ情報を予測し、予測したエントリ情報を、削除対象のエントリ情報として決定する方式である。キャッシュ置換方式としては、LRU (Least Recently Used)方式、LFU (Least Frequently Used)方式、ARC (Adaptive Replacement Cache)方式等がある。

20

【0015】

また、キャッシュメモリの動作方式の一つとして、ライトバックと呼ばれる動作がある。ライトバックは、記憶装置への書き込み要求を受け付けた場合、キャッシュメモリにデータを書き込んでおき記憶装置には書き込まず、CPUの空き時間等に、キャッシュメモリ上の更新されたページを記憶装置に書き込む方法である。以下、記憶装置にまだ書き込んでいない、キャッシュメモリ上の更新されたページを、「Dirtyページ」と呼称する。これに対し、キャッシュメモリ上のDirtyページ以外のページを、「Cleanページ」と呼称する。データの整合性を保つため、Dirtyページを削除する場合には、ライトバックを行ってからDirtyページを削除する。

【0016】

ライトバックを採用することにより、キャッシュメモリが記憶装置よりも高速に記憶することができるという特性を活かすことができ、書き込み性能を高速化することができる。

30

【0017】

しかしながら、ライトバックを採用すると、キャッシュメモリ上のDirtyページを、いつライトバックすればよいのか決定することが困難である。Dirtyページはライトバックしないと削除することができないため、キャッシュ置換方式によってDirtyページが最も必要とされないと予測された場合、そのDirtyページよりも必要とされると予測されたCleanページが削除されることになる。従って、例えば、ライトバックを行う間隔を長くすると、最も必要とされないと予測されたDirtyページの代わりにCleanページが削除されていき、キャッシュメモリのヒット率が低下する可能性がある。一方で、ライトバックを行う間隔を短くすると、ライトバックで得られる高速化の効果が低減することになる。具体的には、キャッシュメモリ上で上書きできれば、記憶装置への書き込み量を減らすことができるが、ライトバックを行う間隔を短くすると、キャッシュメモリ上で上書きできる機会を減らすことになる。また、記憶装置がSSDである場合、SSDには書き込み可能な最大回数がある。従って、頻りにライトバックが行われると、SSDの寿命が短くなるため、ライトバックの回数を低減することが好ましい。

40

【0018】

そこで、本実施の形態では、ARCで採用されているゴースト (Ghost) リストを活用する。具体的には、ARCでは、データとメタデータとを有するエントリ情報を有す

50

る通常のリストと、通常のリストにかつては入っていたが、データが削除され、削除済みデータのメタデータだけを含むエントリ情報を有するゴーストリストとを用いる。ARCの動作例については、図2で説明する。

【0019】

本実施の形態では、通常のリストの本来の削除対象からDirtyページが連続する数分のエントリ情報よりも、ゴーストリストの先頭から前述の数分のエントリ情報が多くアクセスされていれば、本来の削除対象のページをライトバックする方法について説明する。以下の説明では、ARCを採用しているものとする。そして、採用するキャッシュ置換方式としてはどのようなものでもよいが、所定のキャッシュ置換方式としてLRUを採用し、削除対象は、LRUの優先度によって決められるものとする。また、通常のリストと、ゴーストリストとのアクセスの比較は、読み出しの回数で比較するものとする。

10

【0020】

図1を用いて、情報処理装置101の動作例について説明する。キャッシュメモリ103は、通常のリストとなる第1のリストとしてリストT1と、ゴーストリストとなる第2のリストとしてリストB1とを記憶する。さらに、図1の例では、ARCを採用しているため、キャッシュメモリ103は、リストT2、B2を記憶する。図1、図2の例では、通常のリストを実線で示し、ゴーストリストを破線で示す。また、図1、図2の例では、リストT1、B1内のエントリ情報は、MRU(Most Recently Used)からLRUの順に並んでいるものとする。

【0021】

リストT1、B1は、1回もヒットしていないエントリ情報を有する。また、リストT2、B2は、1回以上ヒットしたエントリ情報を有する。リストT1、B1に対して本実施の形態を実施することもできるし、リストT2、B2に対して本実施の形態を実施することもできるし、リストT1、B1と、リストT2、B2との両方に本実施の形態を実施することもできる。図1の例では、リストT1、B1を用いて説明する。

20

【0022】

リストT1は、データとメタデータとを含むエントリ情報111__1~Mを有する。また、リストB1は、メタデータを含むエントリ情報112__1~Lを有する。また、エントリ情報111__1~N、Mは、Dirtyページを有する。図1の例では、Dirtyページを有するエントリ情報111には、網掛けを付与する。

30

【0023】

また、リストT1、B1は、MRUからLRUの順に並んでいるため、リストT1のMRUは、エントリ情報111__Mとなり、リストT1のLRUは、エントリ情報111__1となる。同様に、リストB1のMRUは、エントリ情報112__1となり、リストB1のLRUは、エントリ情報112__Lとなる。エントリ情報112__1は、エントリ情報112__1~Lの中で最も後にリストB1に追加されたエントリ情報である。

【0024】

情報処理装置101は、図1の(1)で示すように、リストT1の削除対象となるエントリ情報111__1がDirtyページである場合、LRUに従った順序におけるエントリ情報111__1からDirtyページが連続する数を計数する。図1の例では、Dirtyページが、エントリ情報111__1~Nまで連続するため、情報処理装置101は、Nと計数する。

40

【0025】

ここで、エントリ情報111__1~Nは、本来ならばリストB1にあるべきものであるが、Dirtyページを有するため、リストT1に留まっている。そして、エントリ情報111__1~Nと同数となるエントリ情報112__1~Nは、本来ならばリストT1にあるべきものであるが、エントリ情報111__1~NがリストT1に留まったため、代わりにリストT1を追い出されたものとなる。なお、エントリ情報111__N+1~Mと、エントリ情報112__N+1~Lとは、ARCの想定通りのリストに含まれている。

【0026】

50

次に、情報処理装置 101 は、図 1 の (2) で示すように、第 1 の合計として、エンタリ情報 111__1 から N 分のエンタリ情報への読み出し回数の合計と、第 2 の合計として、エンタリ情報 112__1 から N 分のエンタリ情報への読み出し回数の合計とを比較する。例えば、情報処理装置 101 は、エンタリ情報 111__1 ~ M、112__1 ~ L のそれぞれの読み出し回数をキャッシュメモリ 103 に記憶し、それぞれの読み出し回数から、第 1 の合計や第 2 の合計を算出してもよい。または、情報処理装置 101 は、第 1 の合計と第 2 の合計とをキャッシュメモリ 103 に記憶し、エンタリ情報 111__1 ~ N に読み出した際に第 1 の合計を 1 増加させ、エンタリ情報 112__1 ~ N に読み出した際に第 2 の合計を 1 増加させてもよい。

【0027】

第 1 の合計が第 2 の合計以上である場合には、実は Read ヒットが増えているケースであるため、ライトバックをすべきタイミングでなく、このままの状態とする。一方で、第 1 の合計よりも第 2 の合計が多い場合には、実は Read ヒットが減っているケースであるため、Read ヒット率を向上させるべく、Dirty ページをライトバックすべきタイミングとなる。

【0028】

図 1 の例では、第 1 の合計よりも第 2 の合計が多い場合であるとする。この場合、情報処理装置 101 は、図 1 の (3) で示すように、エンタリ情報 111__1 の Dirty ページをライトバックする。そして、情報処理装置 101 は、図 1 の (4) で示すように、エンタリ情報 111__1 のページをキャッシュメモリ 103 から削除する。エンタリ情報 111__1 のページを削除することにより、エンタリ情報 111__1 は、リスト B1 に追加されることになる。Dirty ページが減ることにより、Read ヒット率が改善することになる。

【0029】

また、Dirty ページの書き込みを、Watermark モデルで管理しているのであれば、第 1 の合計よりも第 2 の合計が多い場合に、情報処理装置 101 は、Watermark 値を変動し、Dirty ページの書き込みが行われるようにしてもよい。Watermark モデルについては、図 5 で説明する。

【0030】

以上により、情報処理装置 101 は、Read ヒット率を向上できるライトバックの適切なタイミングが判る。次に、ARC の動作例について、図 2 を用いて説明する。

【0031】

図 2 は、ARC の動作例を示す説明図である。ARC は、ゴーストリストを利用することにより、エンタリを含むリストのサイズを調整し、ヒット率を最大化する方法である。また、ARC は、ストレージ製品でよく使用される。

【0032】

ARC では、図 2 に示すように、4 つの LRU リストを管理する。1 つ目のリスト T1 には、データと、メタデータとを記憶しており、一度アクセスされたエンタリ情報が格納される。2 つ目のリスト B1 には、リスト T1 から追い出されたエンタリ情報であって、メタデータを記憶するエンタリ情報が格納される。3 つ目のリスト T2 には、リスト T1 にあるデータがヒットしたエンタリ情報が格納される。リスト T2 に格納されるエンタリ情報は、データと、メタデータとを記憶する。4 つ目のリスト B2 には、リスト T2 から追い出されたエンタリ情報であって、メタデータを記憶するエンタリ情報が格納される。

【0033】

以上により、リスト T1、B1 は、1 回もヒットしていないデータのエンタリ情報を有する。また、リスト T2、B2 は、1 回以上ヒットしたデータのエンタリ情報を有する。また、リスト T1 にあるエンタリ情報が 1 回ヒットした場合、ヒットしたエンタリ情報は、リスト T2 に移動する。また、リスト B1 のエンタリ情報がヒットした場合には、リスト T1 のサイズを大きくする。一方で、リスト B2 のエンタリ情報がヒットした場合には、リスト T2 のサイズを大きくする。

10

20

30

40

50

【0034】

図2の例では、リストT1、B1には、エントリ情報1～3が格納されており、リストT2、B2には、エントリ情報4～6が格納されている。そして、エントリ情報1～3のうち、エントリ情報3がMRUとなり、エントリ情報1がLRUとなる。また、エントリ情報1には、メタデータだけが格納されており、エントリ情報2、3は、データとメタデータとが格納される。また、エントリ情報4～6のうち、エントリ情報6がMRUとなり、エントリ情報4がLRUとなる。また、エントリ情報4、5には、メタデータだけが格納されており、エントリ情報6は、データとメタデータとが格納される。

【0035】

次に、情報処理装置101を、ディスクアレイ装置に適用した例を、図3を用いて説明する。

10

【0036】

(ディスクアレイ装置300のハードウェア構成例)

図3は、ディスクアレイ装置300のハードウェア構成例を示す説明図である。ディスクアレイ装置300は、CE(Controller Enclosure)301と、DE(Drive Enclosure)302とを含む。また、ディスクアレイ装置300は、ホスト装置331に接続する。ホスト装置331は、例えば、サーバである。

【0037】

そして、CE301は、CM(Controller Module)311と、CP SU(CE Power Supply Unit)312と、ディスク313とを有する。また、DE302は、IOM(I/O(Input/Output) Module)321と、DPSU(DE Power Supply Unit)322と、ディスク323とを含む。IOM321は、EXP(SAS(Serial Attached SCSI) Expander)324を含む。

20

【0038】

CE301は、CM311～ディスク313を含む筐体である。CM311は、ディスクアレイ装置300を制御する装置である。また、CM311は、CM間通信を行う。また、CM311は、ホスト装置331と接続する。CM311の内部のハードウェア構成は、図5で説明する。CP SU312は、CE301内部の装置に電源を供給するユニットである。ディスク313は、CM311が使用する記憶装置である。例えば、ディスク313は、SSDやHDDを採用することができる。

30

【0039】

DE302は、IOM321～ディスク323を含む筐体である。IOM321は、CM311とドライブ間とを制御するユニットである。DPSU322は、DE302内部の装置に電源を供給するユニットである。EXP324は、SAS接続用のexpanderチップである。図3に示すEXP324は、ディスク323のそれぞれと接続する。ディスク323は、CM311が使用する記憶装置である。例えば、ディスク323は、SSDやHDDを採用することができる。

【0040】

図4は、CM311のハードウェア構成例を示す説明図である。CM311は、CPU401と、メモリ402と、不揮発性メモリ403と、IOC(I/O Controller)404と、CA(Channel Adapter)405と、EXP406と、SCU(System Capacitor Unit)407とを含む。

40

【0041】

ここで、CM311が、図1に示した情報処理装置101に相当する。また、メモリ402が、図1で示したキャッシュメモリ103に相当する。また、ディスク313、323が、図1に示した記憶装置102に相当する。以下では、説明の簡略化のため、メモリ402が一時的に記憶するデータは、ディスク313のデータであるとする。

【0042】

CPU401は、CM311の全体の制御を司る演算処理装置である。また、CPU4

50

01は、他のCM311のCPU401と接続する。メモリ402は、CPU401のワークエリアとして使用される揮発性メモリである。例えば、メモリ402は、DRAM(Dynamic Random Access Memory)等を採用することができる。不揮発性メモリ403は、本実施の形態におけるキャッシュメモリ制御プログラムを記憶する不揮発性メモリである。不揮発性メモリ403の記憶媒体としては、例えば、NORフラッシュメモリ、NANDフラッシュメモリを採用することができる。

【0043】

IOC404は、CPU401からのI/Oを制御する。図4の例では、IOC404は、EXP406や、他のCM311のEXP406と接続し、CPU401からのディスク313、323へのI/Oを制御する。CA405は、ホスト装置331と通信する通信インターフェースである。EXP406は、SAS接続用のexpanderチップである。図4に示すEXP406は、ディスク313のそれぞれと、EXP324と接続する。SCU407は、停電時に、メモリ402のデータを、不揮発性メモリ403にバックアップするための電源を供給するユニットである。

【0044】

(CM311の機能構成例)

図5は、CM311の機能構成例を示す説明図である。情報処理装置101は、制御部500を有する。制御部500は、I/O受け付け部501と、読み出し/書き込み判定部502と、キャッシュデータ管理部503と、読み出し監視部504と、Watermark調整部505とを含む。制御部500は、記憶装置に記憶されたプログラムをCPU401が実行することにより、各部の機能を実現する。記憶装置とは、具体的には、例えば、図4に示したメモリ402、不揮発性メモリ403や、図3に示したディスク313、323などである。また、各部の処理結果は、メモリ402、CPU401のレジスタ等に格納される。

【0045】

I/O受け付け部501は、ホスト装置331からのアクセス要求を受け付ける。読み出し/書き込み判定部502は、I/O受け付け部501が受け付けたアクセス要求が、読み出し要求なのか書き込み要求なのかを判定する。

【0046】

キャッシュデータ管理部503は、リストT1、B1、T2、B2を管理する。読み出し監視部504は、本来ならばリストT1、T2にあるべきページへのReadミスと、本来ならばリストB1、B2にあるべきDirtyページへのReadヒットとを監視する。Watermark調整部505は、読み出し監視部504からの情報をもとに、Watermark値を変更する。

【0047】

より具体的なキャッシュデータ管理部503~Watermark調整部505の機能について説明する。ここで、キャッシュデータ管理部503~Watermark調整部505が行うリストT1、B1に対する処理と、リストT2、B2に対する処理とは同一であるため、説明の簡略化のため、リストT1、B1に対する処理を用いて説明する。

【0048】

読み出し監視部504は、リストT1のLRUエントリ情報がDirtyページを有する場合、LRUエントリ情報からDirtyページを有するエントリ情報が連続する数を計数する。これにより、本来ならばリストB1、B2にあるべきDirtyページを特定することができる。

【0049】

そして、読み出し監視部504は、本来ならばリストB1、B2にあるべきDirtyページへのReadヒットの回数の第1の合計と、本来ならばリストT1、T2にあるべきページへのReadミスの回数の第2の合計とを比較する。ここで、ReadヒットやReadミスの回数を計数する期間は、どのような長さでもよく、例えば、ディスクアレ

10

20

30

40

50

イ装置300の管理者によって決められた期間である。読み出し監視部504は、第1の合計と第2の合計とを、メモリ402に記憶する。または、リストT1、B1の各エントリ情報が、各エントリ情報の読み出し回数を記憶していてもよい。

【0050】

そして、第1の合計よりも第2の合計が多い場合、キャッシュデータ管理部503は、リストT1のLRUEントリ情報のDirtyページをライトバックし、リストT1のLRUEントリ情報のDirtyページを削除する。また、読み出し監視部504は、第1の合計と第2の合計とを管理するためのカウンタをメモリ402に記憶していてもよい。例えば、読み出し監視部504は、本来ならばリストB1、B2にあるべきDirtyページへのReadヒットがあればカウンタをデクリメントし、本来ならばリストT1、T2にあるべきページへのReadミスがあればカウンタをインクリメントする。そして、カウンタの正負により、キャッシュデータ管理部503は、第1の合計と第2の合計との比較結果を判断する。

10

【0051】

また、Watermark調整部505が、閾値となるWatermark値を変更し、キャッシュデータ管理部503は、変更したWatermark値に基づいて、リストT1のLRUEントリ情報のDirtyページをライトバックしてもよい。ここで、Watermarkモデルについて説明する。

【0052】

Watermarkモデルでは、Low Watermark値と、High Watermark値という2つのパラメータを使用する。そして、Low Watermark値 High Watermark値という関係を有する。Watermarkモデルでのライトバックを行うタイミングは、リストT1内のDirtyページの割合と、Low Watermark値、High Watermark値との関係により、以下の3つのうちのいずれかに分類される。

20

【0053】

1つ目の場合として、リストT1内のDirtyページの割合 < Low Watermark値となる場合、Watermarkモデルでは、何も行わない。場合によっては、idle時にライトバックを行ってもよい。

【0054】

2つ目の場合として、Low Watermark値 リストT1内のDirtyページの割合 High Watermark値となる場合、Watermarkモデルでは、処理を行うCPU等の判断に応じて、ライトバックを行ってもよいし行わなくてもよい。

30

【0055】

3つ目の場合として、High Watermark値 < リストT1内のDirtyページの割合となる場合、Watermarkモデルでは、ライトバックを行う。

【0056】

本実施の形態では、説明の簡略化のため、Low Watermark値 = High Watermark値とし、単純に、「Watermark値」と呼称する。従って、Watermark値は、リストT1内のDirtyページの割合と比較される閾値であって、閾値の方が小さければリストT1内のLRUEントリのDirtyページをライトバックして削除する閾値となる。また、Watermark値は、ページを有するエントリ情報を有するリストと、そのリストのページが削除されたエントリ情報を有するゴーストリストとの組ごとに存在する。従って、本実施の形態では、リストT1、B1との組に対応するWatermark値と、リストT2、B2との組に対応するWatermark値とが存在する。また、読み出し監視部504内のカウンタも、リストT1、B1との組に対応するものと、リストT2、B2との組に対応するものとが存在する。

40

【0057】

Watermark調整部505は、第1の合計よりも第2の合計が多い場合、Wat

50

ermark 値を、現在の値より小さく設定する。また、Watermark 調整部 505 は、第 1 の合計が第 2 の合計よりも多い場合、Watermark 値を、現在の値より大きく設定する。Watermark 値を小さくすると、Dirty ページのライトバックがされやすくなり、Watermark 値を大きくすると、Dirty ページのライトバックがされにくくなる。

【0058】

図 6 は、読み出し時の動作例を示す説明図である。ホスト装置 331 から読み出し要求が発行された場合、I/O 受け付け部 501、読み出し/書き込み判定部 502、キャッシュデータ管理部 503 は、読み出し要求を処理して、ホスト装置 331 に、読み出したデータを送信する。また、読み出し監視部 504 は、読み出しの結果を監視する。

10

【0059】

図 7 は、書き込み時の動作例を示す説明図である。ホスト装置 331 から書き込み要求が発行された場合、I/O 受け付け部 501、読み出し/書き込み判定部 502、キャッシュデータ管理部 503 は、書き込み要求を処理する。このとき、キャッシュデータ管理部 503 は、データを、メモリ 402 といったキャッシュに書き込んで、ホスト装置 331 に書き込み完了通知を行う。キャッシュデータ管理部 503 は、ディスク 313 への書き込みを後で行う。

【0060】

図 8 は、ディスク 313 への書き込みの動作例を示す説明図である。ディスク 313 への書き込み時に、Watermark 調整部 505 は、Watermark 値に基づいて、デステージ (destage) を決定する。ここで、デステージとは、キャッシュメモリとして使用されるメモリ 402 の内容を、ディスク 313 に書き込むことである。

20

【0061】

次に、CM311 が実行する処理を示すフローチャートを、図 9 ~ 図 11 を用いて説明する。ここで、本実施の形態におけるキャッシュメモリ制御方法は、リスト T1、B1 の組に対して行うこともできるし、リスト T2、B2 の組に対して行うこともできるし、リスト T1、B1 の組、リスト T2、B2 の組の両方に対して行うこともできる。図 9 ~ 図 11 では、本実施の形態にキャッシュメモリ制御方法を、リスト T1、B1 の組に対して行う例を用いて説明する。

【0062】

図 9 は、読み出し処理手順の一例を示すフローチャートである。CM311 は、ホスト装置からの読み出し要求を受け付ける (ステップ S901)。次に、CM311 は、リスト T1 における LRU 側から数えた Dirty の連続数 N を計数する (ステップ S902)。そして、CM311 は、Read ヒットしたか否かを判断する (ステップ S903)。ここで、Read ヒットとなる場合とは、読み出し要求の読み出し先のアドレスが、リスト T1 のいずれかのエントリ情報のアドレスに一致する場合である。以下、リスト T1 内のアドレスが一致したエントリ情報を、「Read ヒットしたエントリ情報」と呼ぶ。

30

【0063】

Read ヒットした場合 (ステップ S903: Yes)、CM311 は、続けて、Read ヒットしたエントリ情報が、LRU からみて N 以内か否かを判断する (ステップ S904)。Read ヒットしたエントリ情報が LRU からみて N 以内である場合 (ステップ S904: Yes)、CM311 は、読み出し監視部 504 内のカウンタをデクリメントする (ステップ S905)。

40

【0064】

ステップ S905 の処理終了後、または、Read ヒットしたエントリ情報が LRU からみて N 以内でない場合 (ステップ S904: No)、CM311 は、通常の Read ヒット処理を実行する (ステップ S906)。通常の Read ヒット処理としては、例えば、CM311 は、リスト T1 内の Read ヒットしたエントリ情報のページを読み出して、ホスト装置 331 に送信する。また、ARC の処理に従って、CM311 は、リスト T1 内の Read ヒットしたエントリ情報を、リスト T2 に移動させる。

50

【0065】

一方、読み出し要求が Read ヒットしていない場合（ステップ S 9 0 3 : No）、CM 3 1 1 は、Ghost ヒットしたか否かを判断する（ステップ S 9 0 7）。ここで、Ghost ヒットとなる場合とは、読み出し要求の読み出し先のアドレスが、リスト B 1 のいずれかのエントリ情報のアドレスに一致する場合である。リスト B 1 内のアドレスが一致したエントリ情報を、「Ghost ヒットしたエントリ情報」と呼ぶ。

【0066】

Ghost ヒットした場合（ステップ S 9 0 7 : Yes）、CM 3 1 1 は、続けて、Ghost ヒットしたエントリ情報が、Ghost の MRU からみて N 以内か否かを判断する（ステップ S 9 0 8）。Ghost ヒットしたエントリ情報が Ghost の MRU からみて N 以内である場合（ステップ S 9 0 8 : Yes）、CM 3 1 1 は、読み出し監視部 5 0 4 内のカウンタをインクリメントする（ステップ S 9 0 9）。

10

【0067】

ステップ S 9 0 9 の処理終了後、または、Ghost ヒットしていない場合（ステップ S 9 0 7 : No）、または、Ghost ヒットしたエントリ情報が Ghost の MRU からみて N 以内でない場合（ステップ S 9 0 8 : No）、CM 3 1 1 は、通常の read ミス処理を実行する（ステップ S 9 1 0）。通常の read ミス処理として、CM 3 1 1 は、ディスク 3 1 3 からデータを読み出して、ホスト装置 3 3 1 に送信する。また、ARC の処理に従って、CM 3 1 1 は、削除対象となったエントリ情報を削除し、読み出したデータをリスト T 1 に追加する。

20

【0068】

ステップ S 9 0 6、または、ステップ S 9 1 0 の処理終了後、CM 3 1 1 は、読み出し処理を終了する。

【0069】

図 1 0 は、Watermark 調整部 5 0 5 の処理手順の一例を示すフローチャートである。Watermark 調整部 5 0 5 は、読み出し監視部のカウンタが 0 以上か否かを判断する（ステップ S 1 0 0 1）。読み出し監視部のカウンタが 0 以上である場合（ステップ S 1 0 0 1 : Yes）、Watermark 調整部 5 0 5 は、Watermark 値を、現在の値より小さく設定する（ステップ S 1 0 0 2）。

【0070】

一方、読み出し監視部のカウンタが 0 未満である場合（ステップ S 1 0 0 1 : No）、Watermark 調整部 5 0 5 は、Watermark 値を、現在の値より大きく設定する（ステップ S 1 0 0 3）。

30

【0071】

ステップ S 1 0 0 2、または、ステップ S 1 0 0 3 の処理終了後、Watermark 調整部 5 0 5 は、一定時間後に Watermark 調整部 5 0 5 を再起動するように設定する（ステップ S 1 0 0 4）。ステップ S 1 0 0 4 の処理終了後、Watermark 調整部 5 0 5 は、一連の処理を終了する。

【0072】

図 1 1 は、キャッシュデータ管理部 5 0 3 の処理手順の一例を示すフローチャートである。図 1 1 で示す一連の処理は、図 1 0 で示したステップ S 1 0 0 2、S 1 0 0 3 の処理終了後に行ってもよいし、定期的に行ってもよい。

40

【0073】

キャッシュデータ管理部 5 0 3 は、Watermark 値がリスト T 1 における Dirty ページの割合以下かを判断する（ステップ S 1 1 0 1）。Watermark 値が Dirty ページの割合以下である場合（ステップ S 1 1 0 1 : Yes）、キャッシュデータ管理部 5 0 3 は、Dirty ページをディスク 3 1 3 に書き込む（ステップ S 1 1 0 2）。具体的には、キャッシュデータ管理部 5 0 3 は、リスト T 1 の LRU となるエントリの Dirty ページを、ディスク 3 1 3 に書き込む。また、書き込む Dirty ページの数は、1 つでもよいし、Watermark 値が Dirty ページの割合以下となるまで

50

の数でもよい。書き込む Dirty ページの数が複数である場合、キャッシュデータ管理部 503 は、リスト T1 の LRU の方から順に、Dirty ページをディスク 313 に書き込む。

【0074】

次に、キャッシュデータ管理部 503 は、書き込んだ Dirty ページを削除する（ステップ S1103）。Dirty ページを削除したエントリ情報は、メタデータだけとなるため、リスト B1 に移動することになる。ステップ S1103 の処理終了後、または、Watermark 値が Dirty ページの割合より大きい場合（ステップ S1101 : No）、キャッシュデータ管理部 503 は、一連の処理を終了する。

【0075】

以上説明したように、CM311 は、リスト T1 の LRU から Dirty ページが連続する数分のエントリ情報よりも、リスト B1 の先頭から前述する数分のエントリ情報への読み出しが多ければ、リスト T1 の LRU の Dirty ページをライトバックし削除する。これにより、CM311 は、Read ヒット率を向上できるライトバックの適切なタイミングが判る。また、上述した実施の形態では、リスト T1、B1 における読み出しの回数を比較していたが、書き込みの回数を比較してもよいし、読み出しの回数と書き込みの回数とを合わせたアクセス回数を比較してもよい。例えば、リスト T1、B1 の各エントリ情報が、各エントリ情報自身のアクセス回数を記憶していてもよい。

【0076】

また、CM311 は、Watermark 値を調整することにより、Dirty ページのライトバックを調整してもよい。これにより、本実施の形態は、Watermark モデルを採用している装置に対しても、容易に採用することができる。

【0077】

また、本実施の形態では、通常のリストとなる第 1 のリストは、ディスク 313 のデータのうちの 1 回のアクセスがあったデータと、そのメタデータを有するエントリ情報としてもよい。また、第 1 のリストは、ディスク 313 のデータのうちの 2 回以上のアクセスがあったデータと、そのメタデータを有するエントリ情報としてもよい。これにより、本実施の形態は、ARC を採用する装置に対しても、容易に採用することができる。

【0078】

また、本実施の形態は、ARC を採用していない、通常のリストだけしか有さない装置に対しても適用することができる。ここで、通常のリストは、所定のキャッシュ置換方式として、どのようなキャッシュ置換方式によるリストでもよく、例えば、LRU によるリストでもよいし、LFU によるリストでもよい。そして、通常のリストに対応するゴーストリストを用意することにより、ライトバックの適切なタイミングが判るようにすることができる。

【0079】

なお、本実施の形態で説明したキャッシュメモリ制御方法は、予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等のコンピュータで実行することにより実現することができる。本キャッシュメモリ制御プログラムは、ハードディスク、フレキシブルディスク、CD-ROM (Compact Disc-Read Only Memory)、DVD (Digital Versatile Disk) 等のコンピュータで読み取り可能な記録媒体に記録され、コンピュータによって記録媒体から読み出されることによって実行される。また本キャッシュメモリ制御プログラムは、インターネット等のネットワークを介して配布してもよい。

【0080】

上述した実施の形態に関し、さらに以下の付記を開示する。

【0081】

(付記 1) 複数のデータを記憶する記憶装置と、

前記複数のデータのうちのいずれかのデータと前記いずれかのデータの前記記憶装置上における位置を示す情報とを含むエントリ情報を有する第 1 のリストと、前記第 1 のリス

10

20

30

40

50

トから追い出されたエントリ情報に含まれた削除済みのデータの前記記憶装置上における位置を示す情報を含むエントリ情報を有する第2のリストとを記憶するキャッシュメモリと、

前記第1のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、前記優先度に従った順序における前記削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、前記削除対象のエントリ情報から前記数分のエントリ情報へのアクセス回数の第1の合計よりも前記優先度に従った順序における前記第2のリストに最も後に追加されたエントリ情報から前記数分のエントリ情報へのアクセス回数の第2の合計が多い場合、前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを前記キャッシュメモリから削除する制御部と、
を有することを特徴とする情報処理装置。

10

【0082】

(付記2) 前記制御部は、

前記第1の合計よりも前記第2の合計が多い場合、前記第1のリストに含まれるデータが更新されているエントリ情報の割合と比較される閾値であって前記閾値の方が小さければ前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを削除する前記閾値を、現在の値より小さく設定し、

前記第1の合計が前記第2の合計よりも多い場合、前記閾値を現在の値より大きく設定する、

20

ことを特徴とする付記1に記載の情報処理装置。

【0083】

(付記3) 前記第1のリストは、前記複数のデータのうちの1回のアクセスがあったデータと前記データの記憶装置上における位置を示す情報とを有するエントリ情報を有するリストである、

ことを特徴とする付記1または2に記載の情報処理装置。

【0084】

(付記4) 前記第1のリストは、前記複数のデータのうちの2回以上のアクセスがあったデータと前記データの記憶装置上における位置を示す情報とを有するエントリ情報を有するリストである、

30

ことを特徴とする付記1～3のいずれか一つに記載の情報処理装置。

【0085】

(付記5) コンピュータが、

記憶装置が記憶する複数のデータのうちのいずれかのデータと前記いずれかのデータの記憶装置上における位置を示す情報とを含むエントリ情報を有する第1のリストおよび前記第1のリストから追い出されたエントリ情報に含まれた削除済みのデータの記憶装置上における位置を示す情報を含むエントリ情報を有する第2のリストを記憶するキャッシュメモリの前記第1のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、前記優先度に従った順序における前記削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、

40

前記削除対象のエントリ情報から前記数分のエントリ情報へのアクセス回数の第1の合計よりも前記優先度に従った順序における前記第2のリストに最も後に追加されたエントリ情報から前記数分のエントリ情報へのアクセス回数の第2の合計が多い場合、前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを前記キャッシュメモリから削除する、

処理を実行することを特徴とするキャッシュメモリ制御方法。

【0086】

(付記6) コンピュータに、

記憶装置が記憶する複数のデータのうちのいずれかのデータと前記いずれかのデータの

50

前記記憶装置上における位置を示す情報とを含むエントリ情報を有する第1のリストおよび前記第1のリストから追い出されたエントリ情報に含まれた削除済みのデータの前記記憶装置上における位置を示す情報を含むエントリ情報を有する第2のリストを記憶するキャッシュメモリの前記第1のリストにおける所定のキャッシュ置換方式による優先度に基づいて決定される削除対象のエントリ情報のデータが更新されている場合、前記優先度に従った順序における前記削除対象のエントリ情報からデータが更新されているエントリ情報が連続する数を計数し、

前記削除対象のエントリ情報から前記数分のエントリ情報へのアクセス回数の第1の合計よりも前記優先度に従った順序における前記第2のリストに最も後に追加されたエントリ情報から前記数分のエントリ情報へのアクセス回数の第2の合計が多い場合、前記削除対象のエントリ情報のデータを前記記憶装置に書き出して前記削除対象のエントリ情報のデータを前記キャッシュメモリから削除する、

10

処理を実行させることを特徴とするキャッシュメモリ制御プログラム。

【符号の説明】

【0087】

T 1、T 2、B 1、B 2 リスト

1 0 1 情報処理装置

1 0 2 記憶装置

1 0 3 キャッシュメモリ

1 1 1、1 1 2 エントリ情報

20

3 1 1 C M

5 0 0 制御部

5 0 1 I / O 受け付け部

5 0 2 読み出し / 書き込み判定部

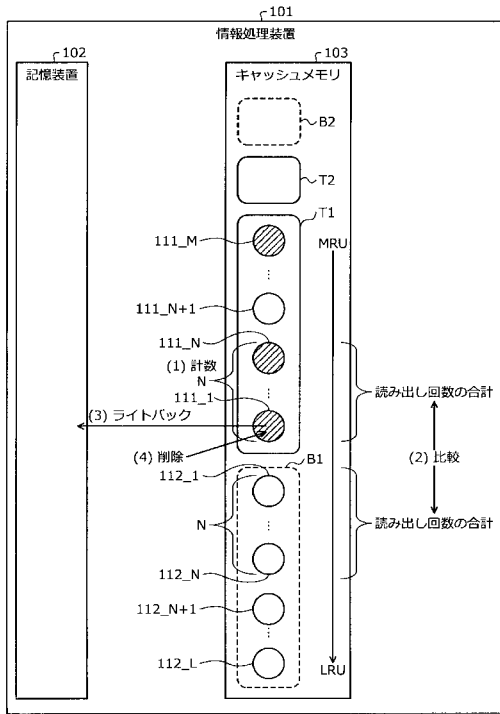
5 0 3 キャッシュデータ管理部

5 0 4 読み出し監視部

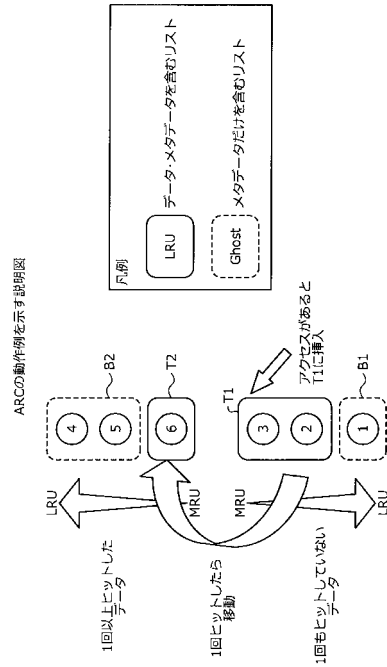
5 0 5 W a t e r m a r k 調整部

【 図 1 】

本実施の形態にかかる情報処理装置101の動作例を示す説明図

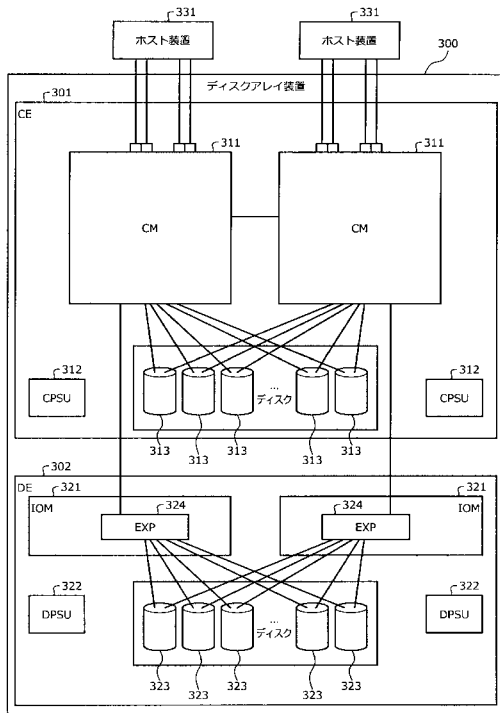


【 図 2 】



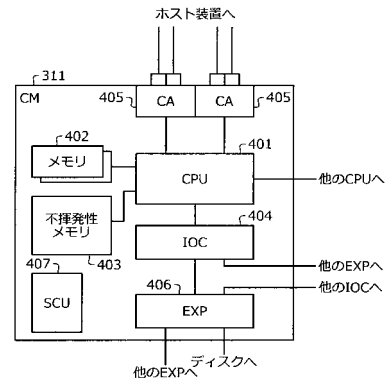
【 図 3 】

ディスクアレイ装置300のハードウェア構成例を示す説明図

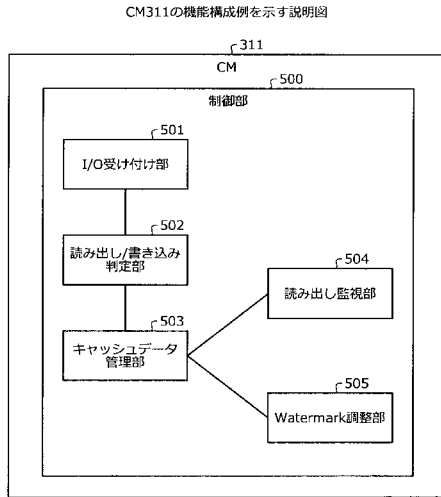


【 図 4 】

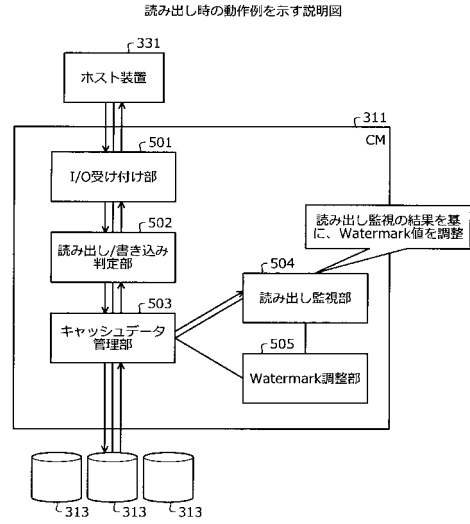
CM311のハードウェア構成例を示す説明図



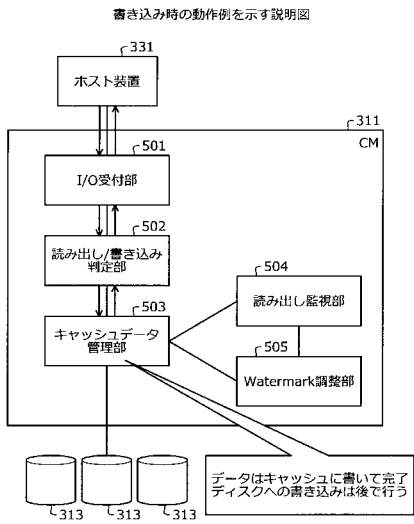
【 図 5 】



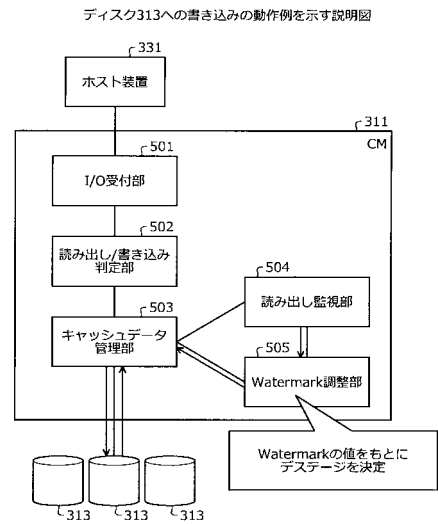
【 図 6 】



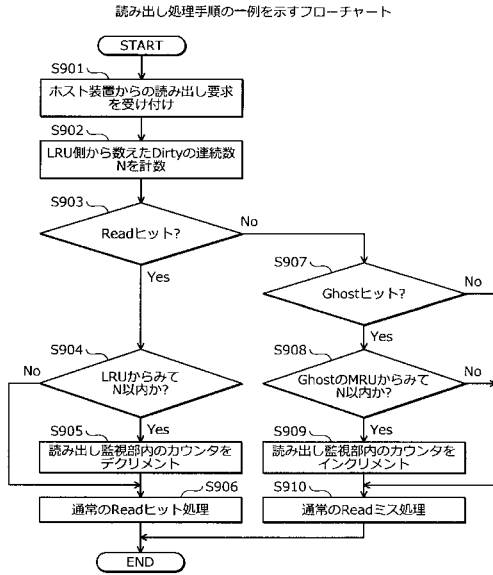
【 図 7 】



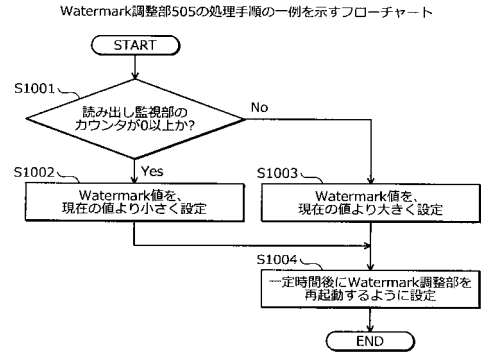
【 図 8 】



【 図 9 】



【 図 1 0 】



【 図 1 1 】

キャッシュデータ管理部503の処理手順の一例を示すフローチャート

