

(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2016-66259

(P2016-66259A)

(43) 公開日 平成28年4月28日 (2016.4.28)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 3/06 (2006.01)	G06F 3/06 301M	
G06F 12/00 (2006.01)	G06F 3/06 304N	
	G06F 3/06 301N	
	G06F 3/06 301Z	
	G06F 12/00 501B	

審査請求 未請求 請求項の数 7 O L (全 32 頁)

(21) 出願番号 特願2014-195001 (P2014-195001)
 (22) 出願日 平成26年9月25日 (2014.9.25)

(71) 出願人 00005223
 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 100092152
 弁理士 服部 毅巖
 (72) 発明者 清水 俊宏
 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 村田 美穂
 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

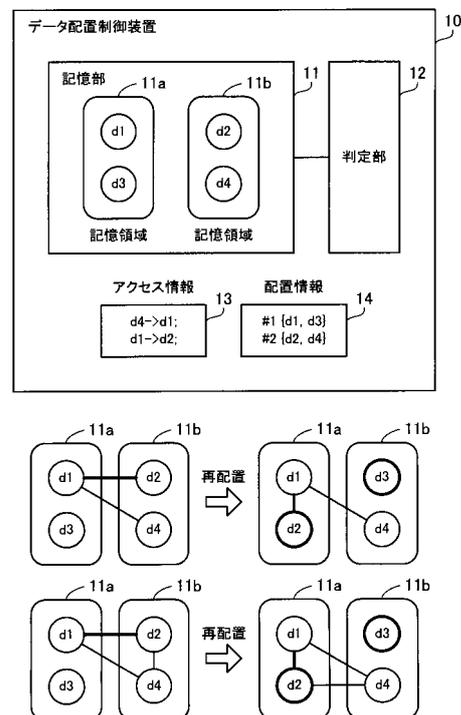
(54) 【発明の名称】 データ配置制御プログラム、データ配置制御装置およびデータ配置制御方法

(57) 【要約】

【課題】 過剰なデータ再配置を抑制する。

【解決手段】 データ配置制御装置 10 は、記憶領域 11 a に配置された単位データ d 1 に対するアクセスの直後に、記憶領域 11 b に配置された単位データ d 2 に対するアクセスが発生した場合に、単位データ d 1 の識別情報と単位データ d 2 の識別情報とに基づいて、単位データ間のアクセス順序を示すアクセス情報 13 を更新する。データ配置制御装置 10 は、アクセス情報 13 と、記憶領域 11 a , 11 b への単位データの配置状況を示す配置情報 14 とに基づいて、単位データ d 1 と関連する第 1 のデータ群および単位データ d 2 と関連する第 2 のデータ群の再配置を行うか否かを判定する。

【選択図】 図 1



【特許請求の範囲】**【請求項 1】**

コンピュータに、

記憶装置の中の複数の記憶領域に分類して配置された複数の単位データに対するアクセスを受け付け、

前記複数の記憶領域のうちの第 1 の記憶領域に配置された第 1 の単位データに対するアクセスの直後に、前記複数の記憶領域のうちの第 2 の記憶領域に配置された第 2 の単位データに対するアクセスが発生した場合に、前記第 1 の単位データの識別情報と前記第 2 の単位データの識別情報とに基づいて、前記複数の単位データの間アクセス順序を示すアクセス情報を更新し、

前記アクセス情報と、前記第 1 の記憶領域および前記第 2 の記憶領域への単位データの配置状況を示す配置情報とに基づいて、前記第 1 の単位データと関連する第 1 のデータ群および前記第 2 の単位データと関連する第 2 のデータ群の再配置を行うか否か判定する、処理を実行させるデータ配置制御プログラム。

【請求項 2】

前記アクセス情報および前記配置情報に基づいて、前記配置状況と前記再配置を行った場合の他の配置状況との差に応じた前記再配置の評価値を算出し、

前記評価値と閾値との比較に応じて、前記再配置を行うか否か判定する、

請求項 1 記載のデータ配置制御プログラム。

【請求項 3】

前記記憶装置からメモリに前記第 1 のデータ群および前記第 2 のデータ群がロードされている場合に、前記メモリ上における前記第 1 のデータ群および前記第 2 のデータ群の更新状況に基づいて、前記閾値を決定する、

請求項 2 記載のデータ配置制御プログラム。

【請求項 4】

前記アクセス情報が示す連続してアクセスされた単位データの組のうち、前記配置状況のもとで前記第 1 の記憶領域と前記第 2 の記憶領域とに分断されて配置された第 1 の組と、前記他の配置状況のもとで前記第 1 の記憶領域と前記第 2 の記憶領域とに分断されて配置される第 2 の組とを検索し、

前記第 1 の組の数と前記第 2 の組の数との差に基づいて、前記評価値を算出する、

請求項 2 または 3 記載のデータ配置制御プログラム。

【請求項 5】

前記第 1 の単位データと前記第 2 の単位データとが連続してアクセスされた回数をカウントし、前記連続してアクセスされた回数の統計情報に基づいて、前記評価値の算出に用いるパラメータの値を決定する、

請求項 2 乃至 4 の何れか一項に記載のデータ配置制御プログラム。

【請求項 6】

複数の単位データが複数の記憶領域に分類して配置された記憶部と、

前記複数の記憶領域のうちの第 1 の記憶領域に配置された第 1 の単位データに対するアクセスの直後に、前記複数の記憶領域のうちの第 2 の記憶領域に配置された第 2 の単位データに対するアクセスが発生した場合に、前記第 1 の単位データの識別情報と前記第 2 の単位データの識別情報とに基づいて、前記複数の単位データの間アクセス順序を示すアクセス情報を更新し、

前記アクセス情報と、前記第 1 の記憶領域および前記第 2 の記憶領域への単位データの配置状況を示す配置情報とに基づいて、前記第 1 の単位データと関連する第 1 のデータ群および前記第 2 の単位データと関連する第 2 のデータ群の再配置を行うか否か判定する判定部と、

を有するデータ配置制御装置。

【請求項 7】

コンピュータが実行するデータ配置制御方法であって、

10

20

30

40

50

記憶装置の中の複数の記憶領域に分類して配置された複数の単位データに対するアクセスを受け付け、

前記複数の記憶領域のうちの第1の記憶領域に配置された第1の単位データに対するアクセスの直後に、前記複数の記憶領域のうちの第2の記憶領域に配置された第2の単位データに対するアクセスが発生した場合に、前記第1の単位データの識別情報と前記第2の単位データの識別情報とに基づいて、前記複数の単位データの間のアクセス順序を示すアクセス情報を更新し、

前記アクセス情報と、前記第1の記憶領域および前記第2の記憶領域への単位データの配置状況を示す配置情報とに基づいて、前記第1の単位データと関連する第1のデータ群および前記第2の単位データと関連する第2のデータ群の再配置を行うか否か判定する、

データ配置制御方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はデータ配置制御プログラム、データ配置制御装置およびデータ配置制御方法に関する。

【背景技術】

【0002】

コンピュータが大量のデータを扱う場合、データを記憶する不揮発性の記憶装置として、HDD (Hard Disk Drive) などの低速・大容量の記憶装置が使用されることが多い。しかし、アクセス要求が発行される毎に低速の記憶装置にアクセスしていると、データアクセスがボトルネックとなってコンピュータの処理性能が低下するおそれがあるという問題がある。そこで、1つの方法として、RAM (Random Access Memory) などのランダムアクセスが高速なメモリを、キャッシュメモリとして使用することが考えられる。

【0003】

例えば、複数の単位データを「セグメント」にグループ化してHDDに格納しておき、セグメント毎にまとめてHDDからRAMにキャッシュするデータ管理装置が提案されている。このデータ管理装置は、ある単位データを指定した読み出し要求を受け付けると、指定された単位データを含むセグメント全体をHDDからRAMにロードする。RAMにロードした(キャッシュした)単位データは、すぐには破棄せずに残しておく。その後、データ管理装置は、キャッシュ中の単位データを指定した読み出し要求を受け付けると、指定された単位データをHDDから読み出す代わりにRAMから取得して提供する。

【0004】

また、データ管理装置は、読み出し要求の履歴を記録しておき、連続して読み出される可能性が高いという単位データ間の関連性を分析する。データ管理装置は、連続して読み出される可能性が高い単位データが同じセグメントに属するように、HDD上の単位データの配置を変更する。これにより、指定された単位データがRAMにキャッシュされている可能性を高めてHDDへのアクセスを減らし、アクセス性能を向上できる。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】国際公開第2013/114538号

【発明の概要】

【発明が解決しようとする課題】

【0006】

しかし、特許文献1に記載のデータ管理装置では、データ再配置が過剰に行われて、低速な記憶装置へのアクセスが削減されない可能性があるという問題がある。

特定の単位データの組が連続してアクセスされやすいという性質(ローカリティ)は、永續するわけではなく情報処理システムの運用に伴って変化し得る。ローカリティが変化すると、前回行ったデータ再配置によるアクセス削減の効果は減少してしまう。すなわち

10

20

30

40

50

、データ再配置のメリットには有効期限が存在し、メリットの大きさは有限である。特許文献1に記載のデータ管理装置では、ローカリティが変化すると、連続してアクセスされやすくなった別の単位データの組が検出され、検出された単位データの組に関して改めてデータ再配置が行われることになる。一方、データ再配置は、低速な記憶装置への書き込みを一時的に増加させることが多く、コストを生じさせる。

【0007】

よって、連続してアクセスされやすい単位データの組が新たに検出される毎に常にデータ再配置を行うと、データ再配置のコストに見合ったメリットが得られないことがあり、低速な記憶装置へのアクセスが削減されない場合がある。

【0008】

1つの側面では、本発明は、過剰なデータ再配置を抑制できるデータ配置制御プログラム、データ配置制御装置およびデータ配置制御方法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】

【0009】

1つの態様では、コンピュータに次の処理を実行させるデータ配置制御プログラムが提供される。記憶装置の中の複数の記憶領域に分類して配置された複数の単位データに対するアクセスを受け付ける。複数の記憶領域のうちの第1の記憶領域に配置された第1の単位データに対するアクセスの直後に、複数の記憶領域のうちの第2の記憶領域に配置された第2の単位データに対するアクセスが発生した場合に、第1の単位データの識別情報と第2の単位データの識別情報とに基づいて、複数の単位データの間のアクセス順序を示すアクセス情報を更新する。アクセス情報と、第1の記憶領域および第2の記憶領域への単位データの配置状況を示す配置情報とに基づいて、第1の単位データと関連する第1のデータ群および第2の単位データと関連する第2のデータ群の再配置を行うか否か判定する。

【0010】

また、1つの態様では、記憶部と判定部とを有するデータ配置制御装置が提供される。

また、1つの態様では、コンピュータが実行するデータ配置制御方法が提供される。

【発明の効果】

【0011】

1つの側面では、過剰なデータ再配置を抑制することができる。

【図面の簡単な説明】

【0012】

【図1】第1の実施の形態のデータ配置制御装置を示す図である。

【図2】第2の実施の形態の情報処理システムを示す図である。

【図3】サーバ装置のハードウェア例を示すブロック図である。

【図4】キャッシュメモリへのページのロード例を示す図である。

【図5】データ更新があったページのライトバック例を示す図である。

【図6】データ再配置があったページのライトバック例を示す図である。

【図7】データ再配置に応じたディスクコストの変化例を示す図である。

【図8】サーバ装置の機能例を示すブロック図である。

【図9】検索テーブルと逆検索テーブルの例を示す図である。

【図10】関連性情報キューと関連性集計テーブルの例を示す図である。

【図11】出現履歴テーブルの例を示す図である。

【図12】パラメータテーブルの例を示す図である。

【図13】アクセス実行の手順例を示すフローチャートである。

【図14】データ再配置の手順例を示すフローチャートである。

【図15】重心法によるデータ再配置の例を示す図である。

【図16】座標テーブルの例を示す図である。

【図17】第1の再配置案生成の手順例を示すフローチャートである。

【図18】ユニオンスプリット法によるデータ再配置の例を示す図である。

10

20

30

40

50

【図19】第2の再配置案生成の手順例を示すフローチャートである。

【図20】データ再配置前後のカット数の変化例を示す図である。

【図21】回帰変数テーブルの例を示す図である。

【図22】パラメータ算出の手順例を示すフローチャートである。

【図23】再出現予測式の変化の例を示す図である。

【図24】他の情報処理システムの例を示す図である。

【発明を実施するための形態】

【0013】

以下、本実施の形態を図面を参照して説明する。

[第1の実施の形態]

図1は、第1の実施の形態のデータ配置制御装置を示す図である。

【0014】

第1の実施の形態のデータ配置制御装置10は、記憶部11と判定部12を有する。

記憶部11は、ランダムアクセスが比較的低速な記憶装置である。記憶部11としては、例えば、HDD、テープ、書き換え可能なディスク媒体、不揮発性の半導体メモリなどを用いることができる。データ配置制御装置10は、記憶部11に対するキャッシュメモリとして、ランダムアクセスが比較的高速な記憶装置を用いるようにしてもよい。キャッシュメモリとしては、例えば、RAMやフラッシュメモリなどを用いることができる。

【0015】

判定部12は、記憶部11へのアクセス性能が向上するように、記憶部11上のデータ配置を制御する。判定部12は、例えば、プロセッサを用いて実現できる。プロセッサは、CPU (Central Processing Unit) やDSP (Digital Signal Processor) であってもよい。また、プロセッサは、ASIC (Application Specific Integrated Circuit) やFPGA (Field Programmable Gate Array) などの特定用途の電子回路を含んでもよい。プロセッサは、例えば、RAMなどのメモリに記憶されたプログラムを実行する。複数のプロセッサの集合 (マルチプロセッサ) を「プロセッサ」と呼ぶこともある。

【0016】

記憶部11には、記憶領域11a, 11bを含む複数の記憶領域が設けられている。記憶領域11a, 11bは、「ページ」や「セグメント」と呼ばれることがある。記憶領域11a, 11bそれぞれには、1または2以上の単位データが配置される。アクセス要求では、読み出しや書き込みの対象として単位データが指定される。ただし、キャッシュメモリを使用する場合、記憶部11へのアクセスを削減するため、記憶部11とキャッシュメモリとの間の転送は記憶領域単位でまとめて行うようにしてもよい。

【0017】

一例として、図1に示すように、記憶領域11aには単位データd1, d3が配置されている。また、記憶領域11bには単位データd2, d4が配置されている。単位データd1と単位データd3を連続してアクセスする場合、同じ記憶領域に配置されているため高速にアクセスすることができる。また、単位データd2と単位データd4を連続してアクセスする場合、同じ記憶領域に配置されているため高速にアクセスすることができる。

【0018】

判定部12は、記憶領域11a, 11bに配置された単位データd1, d2, d3, d4に対するアクセスを受け付ける。単位データd1, d2, d3, d4へのアクセスは、データ配置制御装置10の外部からの要求に応じて発生することもあるし、データ配置制御装置10で実行されるソフトウェアからの要求に応じて発生することもある。

【0019】

判定部12は、記憶領域11aに配置された単位データd1に対するアクセスの直後に、記憶領域11bに配置された単位データd2に対するアクセスが発生したことを検出する。すると、判定部12は、単位データd1の識別情報と単位データd2の識別情報とに基づいて、アクセス情報13を更新する。アクセス情報13は、複数の単位データの間のアクセス順序を示す。例えば、アクセス情報13には、単位データd4の直後に単位デー

10

20

30

40

50

タ d 1 がアクセスされ、単位データ d 1 の直後に単位データ d 2 がアクセスされたことが記録される。アクセス情報 1 3 は、例えば、データ配置制御装置 1 0 が備える H D D などの不揮発性の記憶装置または R A M などの揮発性の記憶装置に記憶されている。

【 0 0 2 0 】

判定部 1 2 は、更新されたアクセス情報 1 3 と配置情報 1 4 とに基づいて、単位データ d 1 と関連する第 1 のデータ群、および、単位データ d 1 の直後にアクセスされた単位データ d 2 と関連する第 2 のデータ群の再配置を行うか否か判定する。配置情報 1 4 は、記憶領域 1 1 a , 1 1 b への単位データの現在の配置状況を示す。

【 0 0 2 1 】

一例として、配置情報 1 4 は、単位データ d 1 , d 3 が記憶領域 1 1 a に配置され、単位データ d 2 , d 4 が記憶領域 1 1 b に配置されていることを示す。配置情報 1 4 は、例えば、データ配置制御装置 1 0 が備える H D D などの不揮発性の記憶装置または R A M などの揮発性の記憶装置に記憶されている。第 1 のデータ群は、例えば、記憶領域 1 1 a または記憶領域 1 1 b に配置され、単位データ d 1 の直前または直後にアクセスされた単位データを含む。第 2 のデータ群は、例えば、記憶領域 1 1 a または記憶領域 1 1 b に配置され、単位データ d 2 の直前または直後にアクセスされた単位データを含む。

10

【 0 0 2 2 】

例えば、判定部 1 2 は、アクセス情報 1 3 と配置情報 1 4 とに基づいて、再配置によるアクセス性能の改善効果（メリット）を示す評価値を算出する。評価値は、現在の配置状況と再配置を行った場合の配置状況との差に基づいて算出できる。記憶領域 1 1 a , 1 1 b をまたがる連続アクセスに着目した場合、判定部 1 2 は、アクセス情報 1 3 が示す連続してアクセスされた単位データの組のうち、現在の配置状況において、異なる記憶領域に分断された単位データの組をカウントする（分断数）。同様に、判定部 1 2 は、再配置を行った場合の配置状況について分断数をカウントする。判定部 1 2 は、再配置前の分断数と再配置後の分断数の差に比例する評価値を算出することが考えられる。

20

【 0 0 2 3 】

評価値が算出されると、例えば、判定部 1 2 は、評価値と再配置のコストを示す閾値とを比較し、評価値が閾値より大きい場合（メリットがコストより大きい場合）に再配置を行うと判定する。再配置のコストは、記憶領域 1 1 a , 1 1 b を書き換えることによって一時的に増加する記憶部 1 1 へのアクセスを示す。再配置を行うと判定した場合、判定部 1 2 は、記憶部 1 1 上で再配置を実行する。ただし、キャッシュメモリが使用されている場合、判定部 1 2 は、単位データ d 1 , d 2 , d 3 , d 4 がキャッシュメモリから追い出されるタイミングを待って記憶領域 1 1 a , 1 1 b を書き換えてもよい。

30

【 0 0 2 4 】

一例として、単位データ d 4 の直後に単位データ d 1 がアクセスされ、単位データ d 1 の直後に単位データ d 2 がアクセスされたとする。そして、単位データ d 1 , d 2 が同一の記憶領域に配置されるよう再配置を行うことを検討するものとする。この場合、再配置前の配置状況では、図 1 に示すように、単位データの組（ d 1 , d 2 ）および（ d 1 , d 4 ）が異なる記憶領域に分断されている。一方、単位データ数の偏りが小さくなるように単位データ d 1 , d 2 を入れ替える再配置を行った場合、図 1 に示すように、単位データの組（ d 1 , d 4 ）のみが異なる記憶領域に分断される。よって、再配置によって分断数は 1 だけ減少する。分断数が 1 減少することに相当するメリットがコストより大きい場合、この再配置を行うと判定される。一方、分断数が 1 減少することに相当するメリットがコスト以下である場合、この再配置は行わないと判定される。

40

【 0 0 2 5 】

また、一例として、上記に加えて、単位データ d 2 の直後に単位データ d 4 がアクセスされたとする。この場合、再配置前の配置状況では、図 1 に示すように、単位データの組（ d 1 , d 2 ）および（ d 1 , d 4 ）が異なる記憶領域に分断されている。一方、単位データ d 1 , d 2 を入れ替える再配置を行った場合、図 1 に示すように、単位データの組（ d 1 , d 4 ）に加えて（ d 2 , d 4 ）が異なる記憶領域に分断される。すなわち、ある単

50

位データの組の分断が解消されるものの、他の単位データの組が新たに分断される。よって、再配置によって分断数は減少しない。通常、この再配置は行わないと判定される。

【0026】

第1の実施の形態のデータ配置制御装置10によれば、単位データ間のアクセス順序を示すアクセス情報13と、記憶領域11a, 11bへの単位データの現在の配置状況を示す配置情報14とに基づいて、再配置を行うか否か判定される。連続してアクセスされた単位データの組とそれら単位データの組の配置状況から、再配置を行った場合のアクセス性能向上の効果を評価することができる。例えば、異なる記憶領域に分断された単位データの組の減少量を、アクセス性能向上の効果として評価することができる。よって、連続してアクセスされた単位データの組が検出される毎に常にデータ再配置を行う場合と比べて、アクセス性能向上の効果が小さい再配置を抑制することができる。

10

【0027】

[第2の実施の形態]

図2は、第2の実施の形態の情報処理システムを示す図である。

第2の実施の形態の情報処理システムは、クライアント装置21, 22およびサーバ装置100を有する。クライアント装置21, 22およびサーバ装置100は、ネットワーク20に接続されている。ネットワーク20は、LAN (Local Area Network) を含んでもよく、インターネットなどの広域ネットワークを含んでもよい。

【0028】

クライアント装置21, 22は、ユーザが操作する端末装置としてのクライアントコンピュータである。クライアント装置21, 22は、サーバ装置100によって管理されるデータを利用して情報処理を行う。このとき、クライアント装置21, 22は、ネットワーク20を介してサーバ装置100にアクセス要求を送信する。アクセス要求は、あるデータを取得するときに発行される読み出し要求(リード要求)であることもあるし、あるデータを更新するときに発行される書き込み要求(ライト要求)であることもある。

20

【0029】

サーバ装置100は、不揮発性の記憶装置に記憶したデータを管理するサーバコンピュータである。サーバ装置100は、例えば、DBMS (Database Management System) を実行している。サーバ装置100は、クライアント装置21, 22からアクセス要求を受信すると、アクセス要求で指定されたデータに対するアクセスを実行し、実行結果をアクセス要求の送信元へ返信する。リード要求を受信した場合、サーバ装置100は、指定されたデータを読み出し、読み出したデータを送信する。ライト要求を受信した場合、サーバ装置100は、指定されたデータを更新し、更新の成否を通知する。

30

【0030】

データへのアクセスを高速化するため、サーバ装置100は、低速・大容量の不揮発性の記憶装置に加えて、高速・小容量のキャッシュメモリを使用する。第2の実施の形態では、前者としてHDDを使用し、後者としてRAMを使用することとする。ただし、前者としてSSD (Solid State Drive) ・フラッシュメモリ・光ディスク・光磁気ディスク・テープなどを使用することもでき、後者としてフラッシュメモリなどを使用することもできる。

40

【0031】

サーバ装置100は、あるデータを指定したアクセス要求を初めて受信すると、指定されたデータを含むデータ集合をHDDからRAMにロードする。RAMにロードしたデータは、アクセス実行後もすぐには消去せずに残しておく。その後、RAMにロード済みのデータ(キャッシュ中のデータ)を指定したアクセス要求を受信すると、サーバ装置100は、HDDからRAMへのロードを省略してアクセスを実行することができる。なお、サーバ装置100は、第1の実施の形態のデータ配置制御装置10の一例である。

【0032】

図3は、サーバ装置のハードウェア例を示すブロック図である。

サーバ装置100は、CPU101、RAM102、HDD103、画像信号処理部1

50

04、入力信号処理部105、媒体リーダー106および通信インタフェース107を有する。上記のユニットは、サーバ装置100内においてそれぞれバス108に接続されている。CPU101は、第1の実施の形態の判定部12の一例である。また、HDD103は、第1の実施の形態の記憶部11の一例である。

【0033】

CPU101は、プログラムの命令を実行する演算回路を含むプロセッサである。CPU101は、HDD103に記憶されたプログラムやデータの少なくとも一部をRAM102にロードし、プログラムを実行する。なお、CPU101は複数のプロセッサコアを備えてもよく、サーバ装置100は複数のプロセッサを備えてもよく、以下で説明する処理を複数のプロセッサまたはプロセッサコアを用いて並列に実行してもよい。また、複数のプロセッサの集合(マルチプロセッサ)を「プロセッサ」と呼んでもよい。

10

【0034】

RAM102は、CPU101が実行するプログラムやCPU101が演算に用いるデータを一時的に記憶する揮発性の半導体メモリである。なお、サーバ装置100は、RAM以外の種類のメモリを備えてもよく、複数個のメモリを備えてもよい。

【0035】

HDD103は、OS(Operating System)やミドルウェアやアプリケーションソフトウェアなどのソフトウェアのプログラム、および、データを記憶する不揮発性の記憶装置である。プログラムには、HDD103上のデータの配置を制御するデータ配置制御プログラムが含まれる。なお、サーバ装置100は、フラッシュメモリやSSDなどの他の種類の記憶装置を備えてもよく、複数の不揮発性の記憶装置を備えてもよい。

20

【0036】

画像信号処理部104は、CPU101からの命令に従って、サーバ装置100に接続されたディスプレイ111に画像を出力する。ディスプレイ111としては、CRT(Cathode Ray Tube)ディスプレイ、液晶ディスプレイ(LCD:Liquid Crystal Display)、プラズマディスプレイ(PDP:Plasma Display Panel)、有機EL(OEL:Organic Electro-Luminescence)ディスプレイなどを用いることができる。

【0037】

入力信号処理部105は、サーバ装置100に接続された入力デバイス112から入力信号を取得し、CPU101に出力する。入力デバイス112としては、マウスやタッチパネルやタッチパッドやトラックボールなどのポインティングデバイス、キーボード、リモートコントローラ、ボタンスイッチなどを用いることができる。また、サーバ装置100に、複数の種類の入力デバイスが接続されていてもよい。

30

【0038】

媒体リーダー106は、記録媒体113に記録されたプログラムやデータを読み取る読み取り装置である。記録媒体113として、例えば、フレキシブルディスク(FD:Flexible Disk)やHDDなどの磁気ディスク、CD(Compact Disc)やDVD(Digital Versatile Disc)などの光ディスク、光磁気ディスク(MO:Magneto-Optical disk)、半導体メモリなどを使用できる。媒体リーダー106は、例えば、記録媒体113から読み取ったプログラムやデータをRAM102またはHDD103に格納する。

40

【0039】

通信インタフェース107は、ネットワーク20に接続され、ネットワーク20を介してクライアント装置21,22と通信を行う。通信インタフェース107は、スイッチなどの通信装置とケーブルで接続される有線通信インタフェースでもよいし、基地局またはアクセスポイントと無線リンクで接続される無線通信インタフェースでもよい。

【0040】

なお、サーバ装置100は、媒体リーダー106を備えていなくてもよく、ユーザが操作する端末装置から制御可能である場合には画像信号処理部104や入力信号処理部105を備えていなくてもよい。また、ディスプレイ111や入力デバイス112が、サーバ装置100の筐体と一体に形成されていてもよい。クライアント装置21,22も、サーバ

50

装置 100 と同様のハードウェア構成によって実現することができる。

【0041】

次に、データのキャッシュおよび HDD 103 上でのデータの配置について説明する。

図 4 は、キャッシュメモリへのページのロード例を示す図である。

サーバ装置 100 は、HDD 103 の記憶領域を複数のページに分割し、ページ単位で HDD 103 からのデータの読み出しや HDD 103 へのデータの書き込みを行う。1つのページは、1つの連続した物理的な記憶領域を示す。ページをセグメントと呼ぶこともある。各ページの大きさは、複数のページの間で同じであってもよいし異なってもよい。1つのページには、複数単位のデータを収容することができる。サーバ装置 100 が関係データベース管理システム (RDBMS: Relational Database Management System) を実行している場合、1単位のデータは、例えば、テーブル内の1つのタプルに相当する。各単位のデータは、例えば、主キーまたは主キー以外の連番によって識別できる。

10

【0042】

一例として、HDD 103 は、ページ 31 (ページ P)、ページ 32 (ページ Q)、ページ 33 (ページ R) およびページ 34 (ページ S) を含む。ページ 31 は、単位データとしてデータ a, b, c を含む。同様に、ページ 32 はデータ d, e, f を含み、ページ 33 はデータ g, h, i を含み、ページ 34 はデータ j, k, l を含む。上記のように、HDD 103 からのデータの読み出しや HDD 103 へのデータの書き込みは、ページ単位で行われる。よって、以下では、あるページに属する全単位のデータを読み出す / 書き込むことを、単にページのデータを読み出す / 書き込むと言うことがある。なお、1つのページに収容できるデータ単位の数には上限が設けられているものとする。

20

【0043】

ここで、サーバ装置 100 は、アクセス要求を受信すると、受信したアクセス要求が指定する単位データを含むページを検索し、検索したページのデータを HDD 103 から RAM 102 にロードする。そして、サーバ装置 100 は、RAM 102 上のデータに対して、アクセス要求が示すアクセスを実行する。サーバ装置 100 は、リード要求に対しては、RAM 102 にロードされた単位データを提供し、ライト要求に対しては、RAM 102 にロードされた単位データを更新する。RAM 102 にロードされたデータは、すぐには破棄されずに HDD 103 のキャッシュとして利用される。後に受信したアクセス要求がキャッシュ中のページに含まれる単位データを指定している場合、サーバ装置 100 は、HDD 103 からの読み出しを省略して RAM 102 上のデータを利用できる。

30

【0044】

一例として、サーバ装置 100 が、データ a, データ e, データ b, データ f, データ g を指定したアクセス要求を順に受信したとする。まず、サーバ装置 100 は、データ a を指定したアクセス要求に対して、データ a の属するページ 31 のデータ (データ a, b, c を含むページ 31 全体) を RAM 102 にロードする。次に、サーバ装置 100 は、データ e を指定したアクセス要求に対して、データ e の属するページ 32 のデータ (データ d, e, f を含むページ 32 全体) を RAM 102 にロードする。

【0045】

次に、サーバ装置 100 は、データ b を指定したアクセス要求に対して、データ b の属するページ 31 がキャッシュ中であるため、HDD 103 へのアクセスを省略し、RAM 102 に存在するデータ b を利用する。次に、サーバ装置 100 は、データ f を指定したアクセス要求に対して、データ f の属するページ 32 がキャッシュ中であるため、HDD 103 からの読み出しを省略し、RAM 102 に存在するデータ f を利用する。次に、サーバ装置 100 は、データ g を指定したアクセス要求に対して、データ g の属するページ 33 のデータ (データ g, h, i を含むページ 33 全体) を RAM 102 にロードする。

40

【0046】

図 5 は、データ更新があったページのライトバック例を示す図である。

キャッシュメモリとして利用できる RAM 102 の記憶領域 (キャッシュ領域) は、データが格納されている HDD 103 の記憶領域と比べて小さい。よって、RAM 102 の

50

キャッシュ領域が不足すると、RAM 102から何れかのページのデータを追い出すことになる。例えば、ページ34のデータをRAM 102にロードしようとしたとき、ページ31, 32, 33がキャッシュ中であり、キャッシュ領域が不足していたとする。この場合、サーバ装置100は、ページ31~33のうち少なくとも1つのデータをRAM 102から追い出して、キャッシュ領域に空きを作ることになる可能性がある。

【0047】

このとき、サーバ装置100は、更新された単位データを含まないページについては、そのページのデータをRAM 102上で破棄すればよく、HDD 103に書き戻さなくてよい。一方、サーバ装置100は、更新された単位データを含むページについては、そのページのデータをHDD 103に書き戻すことになる(ライトバック)。

10

【0048】

例えば、データa, b, c, d, e, f, g, h, iのうちデータeのみが、アクセス要求に応じて更新されたとする。データeの更新は、HDD 103への書き込みを減らすため、すぐにはHDD 103に反映されない。この場合、ページ31のデータをRAM 102から追い出すときは、単にページ31のデータを破棄すればよい。RAM 102上のデータの破棄は、明示的な消去処理を行わずに、そのデータが記憶されていた記憶領域に他のデータを上書きすることによっても実現できる。同様に、ページ33のデータをRAM 102から追い出すときは、単にページ33のデータを破棄すればよい。

【0049】

一方、ページ32のデータをRAM 102から追い出すときは、サーバ装置100は、データeの更新をHDD 103に反映させるため、ページ32のデータ(データd, e, fを含むページ32全体)をHDD 103に書き戻す。ただし、RAM 102上でのデータ更新をHDD 103に反映させるタイミングは、キャッシュ中のデータをRAM 102から追い出すときに限定しなくてもよい。例えば、サーバ装置100は、定期的に、更新された単位データを含むページを確認してライトバックを行うようにしてもよい。

20

【0050】

ところで、あるページのデータがRAM 102にキャッシュされていると、そのページに属する単位データを指定したアクセス要求に対しては、HDD 103からの読み出しを省略することができる。小容量のキャッシュ領域を活用してデータアクセスの性能を向上させるためには、連続してアクセスされる可能性の高い単位データを同じページに配置することが好ましい。そこで、サーバ装置100は、単位データの間のアクセス順序の履歴を記録しておき、履歴に基づいてHDD 103上のデータ配置(何れの単位データを何れのページに配置するか)を動的に変更することとする。

30

【0051】

図6は、データ再配置があったページのライトバック例を示す図である。

一例として、サーバ装置100がクライアント装置21から、データeを指定したアクセス要求を受信し、その直後にデータgを指定したアクセス要求を受信したとする。しかし、現在はデータeはページ32に属しており、データgはページ33に属している。このため、データgを指定したアクセス要求を受信した時点で、ページ33のデータがキャッシュされておらず(キャッシュミスヒットが発生し)、HDD 103からの読み出しが発生する可能性がある。今後もデータe, データgという順番のアクセスが出現する可能性が高い場合、データeとデータgは同じページに属していることが好ましい。

40

【0052】

そこで、サーバ装置100は、データeの属するページ32とデータgの属するページ33の間で、データe, gが同じページに属するように再配置を行うことが考えられる。例えば、サーバ装置100は、データeとデータgを入れ替える。ページ32にはデータd, e, gが含まれ、ページ33にはデータf, h, iが含まれることになる。これにより、データe, データgという順番のアクセスが今後出現した場合、データeのアクセスの時点で、データgを含むページ32のデータがRAM 102にキャッシュされ、データgのアクセスの時点ではHDD 103からの読み出しは原則として発生しない。

50

【 0 0 5 3 】

ページ間のデータ再配置は、再配置の対象となる2つのページのデータがRAM 102にキャッシュされている間に行われる。RAM 102上で行われたデータ再配置は、後でHDD 103に反映される。例えば、ページ32, 33のデータ再配置が行われると、キャッシュされたページ32, 33のデータとHDD 103に記憶されたページ32, 33のデータとは一致していない。よって、サーバ装置100は、ページ32, 33のデータをRAM 102からHDD 103に書き戻すことになる。再配置されたページのデータのライトバックは、前述の更新されたデータのライトバックの場合と同様に、そのページのデータをRAM 102から追い出すときに行うことができる。ただし、定期的に、データ再配置が行われたページを確認してライトバックを行うようにしてもよい。

10

【 0 0 5 4 】

ここで、図5に示したデータ更新と図6に示したデータ再配置の両方を考慮して、RAM 102にキャッシュされたデータをHDD 103に書き戻すコストについて検討する。ページ31, 32, 33のデータがRAM 102にキャッシュされ、その後、ページ31, 32, 33のデータが全てRAM 102から追い出されるとする。キャッシュ中、ページ31, 32, 33に含まれる単位データのうち、データeのみが更新されたとする。

【 0 0 5 5 】

ページ32, 33の間でデータ再配置を行わない場合、図5に示したように、ページ31, 32, 33のうち、更新されたデータeを含むページ32のデータをHDD 103に書き戻すことになる。これに対し、ページ32, 33の間でデータ再配置を行うと、ページ31, 32, 33のうちページ32, 33のデータをHDD 103に書き戻すことになる。ページ32のライトバックはデータ再配置の有無に関係なく発生する一方、ページ33のライトバックはデータ再配置を行う場合のみ発生する。すなわち、この例では、データ再配置を行うと、ページ1つ分だけHDD 103へのデータの書き込みが増加する。増加するデータの書き込みは、データ再配置のコストとして認識することができる。

20

【 0 0 5 6 】

一方で、ある特定の単位データと他の特定の単位データとが連続してアクセスされる可能性が高いという性質(ローカリティ)は、永續するとは限らず変化し得る。例えば、データe, gが連続してアクセスされる可能性が高いことに応じて上記のデータ再配置を行った後、ローカリティが変化して、データd, fが連続してアクセスされる可能性が高くなったとする。すると、サーバ装置100は、現在データdとデータfが異なるページに属しているため、更にページ32, 33の間でデータ再配置を行うことになる。このように、ローカリティが変化する可能性を考慮すると、データ再配置を行うことで得られるデータアクセスの性能向上というメリットは、有限の値として評価される。

30

【 0 0 5 7 】

すなわち、データアクセスのローカリティが変化する環境下では、コストがメリットを上回るためにデータ再配置を行わない方が好ましい場合が存在する。そこで、サーバ装置100は、データ再配置を実行することで発生するコストとデータ再配置を実行することで得られるメリット(再配置を実行しないことのペナルティ(機会コスト)と言うこともできる)とを比較して、データ再配置を実行するか否か判定する。

40

【 0 0 5 8 】

図7は、データ再配置に応じたディスクコストの変化例を示す図である。

ここでは、HDD 103からRAM 102にデータを読み出すコストと、RAM 102からHDD 103へデータを書き戻すコストとを合わせた、ディスクコストを考える。グラフ41は、ローカリティの持続が長い場合、すなわち、ローカリティの変化が小さい場合のディスクコストの時間変化を示す。グラフ42は、ローカリティの持続が短い場合、すなわち、ローカリティの変化が大きい場合のディスクコストの時間変化を示す。

【 0 0 5 9 】

グラフ41において、サーバ装置100がデータ再配置を全く行わない場合、ランダムにキャッシュミスヒットが発生し、HDD 103からのデータの読み出しが安定的に行わ

50

れる。一方、データ再配置を行わないため、HDD 103へのデータの書き込みは増加しない。よって、データ再配置なしの場合のディスクコストは一定のレベルに安定する。

【0060】

これに対し、グラフ41において、サーバ装置100がデータ再配置を行う場合、データ再配置によってHDD 103へのデータの書き込みが一時的に増加する。一方、データ再配置がHDD 103に反映されると、その後はキャッシュミスヒットが減少し、HDD 103からのデータの読み出しが抑制される。よって、データ再配置ありの場合のディスクコストは、一時的に増加した後に大きく減少する。ローカリティの持続が長いことから、減少したディスクコストはしばらくの間維持される。その後、ローカリティが変化すると、データ再配置の効果が徐々に消えてしまい、ディスクコストは再配置なしの場合と同じレベルまで増加する。変化後のローカリティに応じてサーバ装置100がデータ再配置を行うと、再びディスクコストは一時的に増加した後に大きく減少する。

10

【0061】

グラフ42において、サーバ装置100がデータ再配置を全く行わない場合、グラフ41の場合と同様に、ディスクコストは一定のレベルに安定する。これに対し、サーバ装置100がデータ再配置を行う場合、ディスクコストは一時的に増加した後に減少し始める。しかし、ローカリティの持続が短いことから、データ再配置の効果が早く消えてしまい、ディスクコストは十分に低下する前に増加し始める。様々な単位データの間のローカリティは、一斉に変化するわけではなく異なるタイミングに分散して変化する。そのため、複数のページ全体のディスクコストは、図7に示すように緩やかに変換する。すなわち、データ再配置ありの場合のディスクコストは、データ再配置のコストとしての増加、データ再配置のメリットとしての微減少、ローカリティの変化に伴う増加を繰り返す。

20

【0062】

ローカリティの持続が長い場合、再配置ありのディスクコストの積分値は、再配置なしのディスクコストの積分値よりも小さい。すなわち、サーバ装置100が積極的にデータ再配置を行うことで、ディスクコストが低減し、データアクセスの性能が向上する。一方、ローカリティの持続が短い場合、再配置ありのディスクコストの積分値は、再配置なしのディスクコストの積分値よりも大きい。すなわち、サーバ装置100が積極的にデータ再配置を行うことで、ディスクコストが増加し、かえってデータアクセスの性能が低下するおそれがある。このように、データ再配置は一時的にディスクコストを増加させるため、常にデータ再配置を行うことがデータアクセスの平均性能を向上させるとは限らない。

30

【0063】

そこで、サーバ装置100は、データ再配置の実行コストと不実行ペナルティとを評価し、前者が後者より小さい場合に限定してデータ再配置を実行することとする。

次に、サーバ装置100によるデータ再配置について説明する。

【0064】

図8は、サーバ装置の機能例を示すブロック図である。

サーバ装置100は、データ記憶部121、キャッシュ部122、制御情報記憶部123、アクセス実行部131、再配置制御部133およびパラメータ算出部136を有する。データ記憶部121は、HDD 103に確保した記憶領域として実現できる。キャッシュ部122および制御情報記憶部123は、RAM 102に確保した記憶領域として実現できる。アクセス実行部131、再配置制御部133およびパラメータ算出部136は、例えば、CPU 101が実行するプログラムのモジュールとして実装することができる。

40

【0065】

データ記憶部121は、連続した物理的な記憶領域として、それぞれ1または2以上の単位データを記憶することができる複数のページを含む。単位データは、識別情報によって識別されアクセス要求に応じてアクセスされるデータの単位であり、例えば、テーブルの1つのタプルに相当する。データ記憶部121からのデータの読み出しやデータ記憶部121へのデータの書き込みは、アクセス実行部131によってページ単位で行われる。

【0066】

50

キャッシュ部 1 2 2 は、データ記憶部 1 2 1 に対するキャッシュメモリである。キャッシュ部 1 2 2 の記憶容量は、データ記憶部 1 2 1 より小さい一方、キャッシュ部 1 2 2 のアクセス速度（特に、ランダムアクセス速度）は、データ記憶部 1 2 1 より速い。キャッシュ部 1 2 2 には、データ記憶部 1 2 1 に含まれる複数のページのうち一部のページのデータが、ページ単位でロードされる。アクセス要求に応じたデータ更新やデータ再配置は、キャッシュ部 1 2 2 にロードされたデータに対して行われ、キャッシュ部 1 2 2 からそのデータが追い出されるときにデータ記憶部 1 2 1 に反映される。

【 0 0 6 7 】

制御情報記憶部 1 2 3 は、データアクセス、キャッシュ管理およびデータ再配置の制御に用いられる制御情報を記憶する。制御情報には、ページと単位データの対応関係を示す検索情報、連続してアクセスされた単位データの組を示す履歴情報、データ再配置の実行コストおよびデータ再配置の不実行ペナルティを算出するときに用いるパラメータを示すパラメータ情報などが含まれる。制御情報の詳細は後述する。

10

【 0 0 6 8 】

アクセス実行部 1 3 1 は、アクセス要求を受信し、受信したアクセス要求に応じて、キャッシュ部 1 2 2 にキャッシュされたデータに対するアクセスを実行する。リード要求を受信した場合、アクセス実行部 1 3 1 は、リード要求で指定された単位データをキャッシュ部 1 2 2 から取得し、取得した単位データを返信する。ライト要求を受信した場合、アクセス実行部 1 3 1 は、ライト要求に含まれるデータを用いてキャッシュ部 1 2 2 上の単位データを更新し、更新の成否を返信する。また、アクセス実行部 1 3 1 は、受信したア

20

【 0 0 6 9 】

アクセス実行部 1 3 1 は、キャッシュ制御部 1 3 2 を有する。キャッシュ制御部 1 3 2 は、データ記憶部 1 2 1 からキャッシュ部 1 2 2 へのデータのロードを制御する。キャッシュ制御部 1 3 2 は、未キャッシュの単位データを指定したアクセス要求を受信すると、当該単位データを含むページのデータ全体をデータ記憶部 1 2 1 からキャッシュ部 1 2 2 にロードする。データをロードするにあたり、キャッシュ部 1 2 2 の空き領域が不足している場合、キャッシュ制御部 1 3 2 は、キャッシュ中の何れかのページのデータをキャッシュ部 1 2 2 から追い出す。追い出すページのデータに対してキャッシュ部 1 2 2 上で更新または再配置が行われていた場合、キャッシュ制御部 1 3 2 は、追い出すページのデータ全体をキャッシュ部 1 2 2 からデータ記憶部 1 2 1 に書き戻す。

30

【 0 0 7 0 】

再配置制御部 1 3 3 は、制御情報記憶部 1 2 3 に記憶された履歴情報を分析し、キャッシュ中のページのデータに対してキャッシュ部 1 2 2 上で再配置を実行する。再配置制御部 1 3 3 は、再配置案生成部 1 3 4 および実行可否判定部 1 3 5 を有する。

【 0 0 7 1 】

再配置案生成部 1 3 4 は、所定の開始条件が満たされると、ページと単位データの現在の対応関係、および、最近連続してアクセスされた単位データの組に基づいて、再配置案を生成する。開始条件は、例えば、データ再配置を前回検討してからの経過時間や履歴情報の蓄積量などを基準として予め決められる。再配置案は、例えば、データ再配置を行う 2 つのページの識別情報と、当該 2 つのページの間で移動する単位データの識別情報とを用いて表現される。再配置案生成部 1 3 4 は、連続してアクセスされた単位データができる限り同じページに属するように、ページ間での単位データの移動を検討する。

40

【 0 0 7 2 】

実行可否判定部 1 3 5 は、再配置案生成部 1 3 4 が再配置案を生成すると、制御情報記憶部 1 2 3 に記憶されたパラメータ情報を用いて、その再配置案を採用した場合の実行コストと不実行ペナルティとを算出する。実行コストは、データ記憶部 1 2 1 に書き戻すページの増加量、HDD 1 0 3 の書き込み速度などを考慮して算出される。不実行ペナルティは、連続してアクセスされた単位データが異なるページに分かれている状況の改善度、ある単位データの組が今後連続してアクセスされる回数の期待値、HDD 1 0 3 の読み出

50

し速度などを考慮して算出される。実行コストおよび不実行ペナルティの詳細は後述する。実行可否判定部 135 は、算出した実行コストと不実行ペナルティとを比較する。実行可否判定部 135 は、不実行ペナルティが実行コストより大きい場合は再配置案を採用し、不実行ペナルティが実行コスト以下である場合は再配置案を採用しない。

【0073】

パラメータ算出部 136 は、制御情報記憶部 123 に記憶された履歴情報を分析し、実行コストおよび不実行ペナルティの算出に用いるパラメータ情報を生成する。例えば、パラメータ算出部 136 は、ある単位データの組について、過去の連続アクセスの出現状況を分析して、その連続アクセスが今後出現する回数を予測するための予測式を求める。この予測式は、サーバ装置 100 におけるローカリティの持続度を反映している。

10

【0074】

図 9 は、検索テーブルと逆検索テーブルの例を示す図である。

検索テーブル 141 は、制御情報記憶部 123 に記憶されている。検索テーブル 141 は、データ ID およびページ ID の項目を有する。データ ID は、単位データを識別する識別情報である。データ ID として、テーブルの主キーを用いてもよいし、DBMS によって自動的に付与される連番を用いてもよい。ページ ID は、ページを識別する識別情報である。ページ ID として、HDD 103 のアドレスを用いてもよい。

【0075】

検索テーブル 141 では、1つのデータ ID に対して 1つのページ ID が対応付けられる。これは、そのデータ ID をもつ単位データがそのページ ID をもつページに属していることを示している。検索テーブル 141 を用いることで、ある単位データのデータ ID から、その単位データが属するページのページ ID を検索できる。

20

【0076】

逆検索テーブル 142 は、制御情報記憶部 123 に記憶されている。逆検索テーブル 142 は、ページ ID、データ ID、更新フラグおよび再配置フラグの項目を有する。更新フラグは、あるページに属する単位データの中に、キャッシュ部 122 上で更新された単位データが存在するか否かを示す。更新フラグ = 1 は、更新された単位データがあり、その更新がデータ記憶部 121 に未反映であることを示す。更新フラグ = 0 は、更新された単位データがないことを示す。再配置フラグは、あるページに対してキャッシュ部 122 上でデータ再配置が実行されたか否かを示す。再配置フラグ = 1 は、データ再配置が行われており、そのデータ再配置がデータ記憶部 121 に未反映であることを示す。再配置フラグ = 0 は、データ再配置が行われていないことを示す。

30

【0077】

逆検索テーブル 142 では、1つのページ ID に対して、0 または 1 以上のデータ ID と 1つの更新フラグと 1つの再配置フラグとが対応付けられる。逆検索テーブル 142 を用いることで、あるページのページ ID から、そのページに属する全ての単位データのデータ ID を検索できる。また、逆検索テーブル 142 を用いることで、あるページのページ ID から、そのページに対応する更新フラグと再配置フラグを検索できる。

【0078】

なお、キャッシュ部 122 上でのデータ再配置については、RAM 102 上で単位データを移動してその単位データの格納位置を変更してもよい。また、キャッシュ部 122 上でのデータ再配置については、RAM 102 上の単位データの格納位置は変更せずに、検索テーブル 141 および逆検索テーブル 142 の更新のみ行うようにしてもよい。

40

【0079】

図 10 は、関連性情報キューと関連性集計テーブルの例を示す図である。

関連性情報キュー 143 は、制御情報記憶部 123 上に形成されている。関連性情報キュー 143 は、先入れ先出し (FIFO: First In First Out) のリスト構造をもつ。関連性情報キュー 143 には、アクセス要求が到着する毎に関連性情報が追加される。

【0080】

関連性情報は、クライアント ID、データ ID および前データ ID を含む。クライアン

50

トIDは、アクセス要求を送信したクライアント装置を識別する識別情報である。クライアントIDとして、クライアント装置21, 22の通信アドレス(例えば、IP(Internet Protocol)アドレス)を用いてもよい。関連性情報に含まれるデータIDは、アクセス要求で指定された単位データのデータIDである。前データIDは、同じクライアント装置が前回送信したアクセス要求で指定された単位データのデータIDである。

【0081】

関連性情報は、前データIDが示す単位データの直後にデータIDが示す単位データがアクセスされたという、単位データ間の「関連性」を示している。1つ前にアクセスされた単位データは、例えば、今回のアクセスと同じクライアント装置についての直近の関連性情報を関連性情報キュー143から検索することで特定することができる。ただし、クライアント装置21, 22が、1つ前にアクセスした単位データのデータIDを、前データIDとしてアクセス要求に付加するようにしてもよい。以下では、データIDが示す単位データと前データIDが示す単位データの組を「関連データ対」と言うことがある。

10

【0082】

関連性情報キュー143に登録された関連性情報は、再配置案生成部134が再配置案を生成するとき、登録順に従って1つずつ抽出される。再配置案生成部134によって利用された関連性情報は、関連性情報キュー143から消去される。また、あるページのデータがキャッシュ部122から追い出されると、そのページに属する単位データに関する関連性情報は関連性情報キュー143から消去される。すなわち、関連性情報キュー143には、キャッシュ部122にキャッシュされているページについての関連性情報であって、データ再配置の検討にまだ利用されていないものが蓄積される。

20

【0083】

関連性集計テーブル144は、制御情報記憶部123に記憶されている。関連性集計テーブル144は、データIDおよび重みの項目を有する。重みの項目には、データIDの項目が示す単位データの直前にアクセスされた単位データを示す識別情報と、そのアクセス順序の出現回数とが登録される。例えば、データbに対して{a:2, c:2}という重み情報が登録される。これは、データaの直後にデータbがアクセスされたことが2回あり、データcの直後にデータbがアクセスされたことが2回あることを示す。

【0084】

関連性集計テーブル144は、関連性情報キュー143に関連性情報が追加される毎、すなわち、アクセス要求が到着する毎に、追加された関連性情報に従って更新される。関連性集計テーブル144を用いることで、ある単位データのデータIDから、その直前にアクセスされた単位データのデータIDと、その関連データ対の出現回数を検索できる。あるページのデータがキャッシュ部122から追い出されると、追い出されたページに属する単位データに関する重み情報は関連性集計テーブル144から消去される。すなわち、関連性集計テーブル144には、ページがキャッシュ部122にキャッシュされている一期間内に出現した、そのページに関する関連データ対の出現回数が集計される。

30

【0085】

図11は、出現履歴テーブルの例を示す図である。

出現履歴テーブル145は、制御情報記憶部123に記憶されている。出現履歴テーブル145は、関連データ対および出現回数の項目を有する。1つの関連データ対に対して、出現回数の列が対応付けられる。関連データ対の出現回数は、所定の区分基準に従って複数の期間に区分してカウントされる。例えば、今日の出現回数, 前日の出現回数, 前々日の出現回数, ...のように、日単位で出現回数がカウントされる。

40

【0086】

出現履歴テーブル145は、関連性情報キュー143に関連性情報が追加される毎、すなわち、アクセス要求が到着する毎に、追加された関連性情報に従って更新される。例えば、関連性情報キュー143に関連性情報が追加されると、追加された関連性情報が示す関連データ対に対応する今日の出現回数を1だけ加算(インクリメント)する。出現履歴テーブル145には、長期間の出現回数を蓄積することができる。関連性情報キュー14

50

3 および関連性集計テーブル 1 4 4 からは、キャッシュ部 1 2 2 から追い出されたページに関する情報が消去されるのに対し、出現履歴テーブル 1 4 5 には、キャッシュ部 1 2 2 から追い出されたページに関する情報も蓄積しておいてよい。ただし、所定期間以上古い出現回数の情報は、出現履歴テーブル 1 4 5 から消去してもよい。

【 0 0 8 7 】

図 1 2 は、パラメータテーブルの例を示す図である。

パラメータテーブル 1 4 6 は、制御情報記憶部 1 2 3 に記憶されている。パラメータテーブル 1 4 6 には、実行可否判定部 1 3 5 がデータ再配置の実行コストおよび不実行ペナルティを算出するときに用いるパラメータの名称と値が登録される。パラメータの値の少なくとも一部は、パラメータ算出部 1 3 6 によって動的に算出される。パラメータの値の中には、ユーザによって静的に設定されるものが含まれていてもよい。

10

【 0 0 8 8 】

パラメータには、書き込み速度、読み出し速度、全体の予測式、および、複数の関連データ対それぞれに対応する個別の予測式が含まれる。書き込み速度は、1 ページのデータを RAM 1 0 2 から HDD 1 0 3 に書き戻すのに要する時間を示す。読み出し速度は、1 ページのデータを HDD 1 0 3 から RAM 1 0 2 にロードするのに要する時間を示す。書き込み速度および読み出し速度それぞれの単位は、例えば、ミリ秒毎ページである。

【 0 0 8 9 】

なお、書き込み速度および読み出し速度は、ユーザが HDD 1 0 3 の物理性能とページサイズの期待値から推定しておき、予めパラメータテーブル 1 4 6 に登録しておいてもよい。また、ユーザが書き込み速度および読み出し速度を実測し、実測値の平均を予めパラメータテーブル 1 4 6 に登録しておいてもよい。また、パラメータ算出部 1 3 6 が HDD 1 0 3 の書き込み速度および読み出し速度を監視し、パラメータテーブル 1 4 6 に登録された書き込み速度および読み出し速度の値を継続的に更新してもよい。

20

【 0 0 9 0 】

予測式は、関連データ対の過去の出現状況から、今後一定期間内に同じ関連データ対が出現する回数を予測する式である。予測式は、例えば、 $y = u_1 \times x_1 + u_2 \times x_2 + u_3 \times x_3 + \dots$ という線形式である。変数 y (目的変数) は、関連データ対の将来の出現回数の期待値 (再出現期待値) を表し、変数 x_1, x_2, x_3, \dots (説明変数) は、関連データ対の過去の出現状況に応じた特徴量を表す。係数 u_1, u_2, u_3, \dots は、特徴量の重みを表す。パラメータ算出部 1 3 6 は、後述するように、出現履歴テーブル 1 4 5 を用いて回帰分析を行い、係数 u_1, u_2, u_3, \dots を算出する。

30

【 0 0 9 1 】

このとき、パラメータ算出部 1 3 6 は、様々な関連データ対についてのデータを合わせて回帰分析することで、特定の関連データ対に限定しない全体の予測式の係数 u_1, u_2, u_3, \dots を算出することができる。また、パラメータ算出部 1 3 6 は、特定の関連データ対についてのデータを回帰分析することで、その関連データ対に対応する予測式の係数 u_1, u_2, u_3, \dots を算出することができる。全体の予測式および複数の個別の予測式の間では、係数 u_1, u_2, u_3, \dots が異なることが多い。

【 0 0 9 2 】

後述する例では、基準日の前日に所望の関連データ対が出現したか否かを示す変数 x_1 と、基準日の前々日に所望の関連データ対が出現したか否かを示す変数 x_2 と、基準日から一定期間前までの出現率を示す変数 x_3 と、基準日の季節を示す変数 x_4 を用いる。変数 y は、基準日から一定期間後までに所望の関連データ対が出現する回数の期待値を示す。変数 y の出現回数をカウントする期間 (変数 y の値を定める「一定期間」) は、1 つのページが連続してキャッシュされている期間の平均に応じて決めてもよい。

40

【 0 0 9 3 】

このようにして決定した予測式は、サーバ装置 1 0 0 におけるローカリティの持続傾向を反映している。ある関連データ対の出現回数は、その関連データ対がバースト的に出現し初めてから収束するまでの間、一定に安定するわけではなく非線形に増減することがあ

50

る。そこで、直近の出現状況に関する複数の特徴量を用いることで、出現回数の分布のどの地点まで現在進んでおり、今後どの程度の出現回数が期待されるかを推定できる。ただし、パラメータ算出部 136 は、特徴量を用いた回帰分析に代えて、各関連データ対の出現回数の分布を詳細に分析して、再出現期待値を算出できるようにしてもよい。

【0094】

次に、サーバ装置 100 が実行する処理の手順について説明する。

図 13 は、アクセス実行の手順例を示すフローチャートである。

(S10) アクセス実行部 131 は、クライアント装置 21, 22 の何れかから、ネットワーク 20 を介してアクセス要求を受信する。アクセス要求は、ある単位データを読み出すリード要求またはある単位データを更新するライト要求などである。

10

【0095】

(S11) アクセス実行部 131 は、制御情報記憶部 123 に記憶された検索テーブル 141 を参照して、アクセス要求で指定された単位データを含むページ T を検索する。

(S12) キャッシュ制御部 132 は、検索されたページ T がキャッシュ中であるか、すなわち、ページ T のデータがキャッシュ部 122 に記憶されているか判断する。検索されたページ T がキャッシュ中である場合はステップ S19 に処理が進み、ページ T が未キャッシュである場合はステップ S13 に処理が進む。

【0096】

なお、各ページがキャッシュ中であるか判断するため、キャッシュ制御部 132 は、キャッシュ中のページを示すリストまたは未キャッシュのページを示すリストを保持していてもよい。また、逆検索テーブル 142 にはキャッシュ中のページに関する情報のみ登録するようにし、キャッシュ制御部 132 は、所望のページの情報に逆検索テーブル 142 に存在するか確認することで当該ページがキャッシュ中か否かを判断してもよい。また、各ページがキャッシュ中か否かを示すフラグを逆検索テーブル 142 に追加してもよい。

20

【0097】

(S13) キャッシュ制御部 132 は、キャッシュ部 122 のキャッシュ領域に、ページ T のデータを格納するだけの空きが存在するか判断する。キャッシュ領域に空きが存在するか否かは、キャッシュ中のページの数に所定の上限に達しているか否かによって判断してもよい。キャッシュ領域に空きが存在する場合はステップ S18 に処理が進み、キャッシュ領域が不足している場合はステップ S14 に処理が進む。

30

【0098】

(S14) キャッシュ制御部 132 は、キャッシュ中の複数のページのうちキャッシュ部 122 から追い出すページ U を選択する。ページ U を選択するアルゴリズム (キャッシュアルゴリズムや置換アルゴリズムなどと呼ばれることがある) としては、様々なものが考えられる。例えば、LRU (Least Recently Used)、LFU (Least Frequency Used)、FIFO などのアルゴリズムを用いることができる。キャッシュ制御部 132 は、使用するアルゴリズムに応じた情報 (例えば、ページのアクセス回数) を保持してもよい。

【0099】

(S15) キャッシュ制御部 132 は、逆検索テーブル 142 から、ステップ S14 で選択したページ U に対応する更新フラグと再配置フラグを取得する。そして、キャッシュ制御部 132 は、更新フラグ = 1 または再配置フラグ = 1 であるか、すなわち、ページ U に含まれる単位データが更新されたかまたはページ U に対してデータ再配置が行われたか判断する。更新フラグ = 1 または再配置フラグ = 1 である場合、ステップ S16 に処理が進む。更新フラグ = 0 かつ再配置フラグ = 0 である場合、ステップ S17 に処理が進む。

40

【0100】

(S16) キャッシュ制御部 132 は、キャッシュ部 122 に記憶されたページ U のデータ全体をデータ記憶部 121 に書き戻す。すなわち、RAM 102 にキャッシュされたページ U のデータが HDD 103 に書き戻される。

【0101】

(S17) キャッシュ制御部 132 は、逆検索テーブル 142 に登録されたページ U に

50

対応する更新フラグおよび再配置フラグを「0」にクリアする。また、キャッシュ制御部132は、ページUに含まれる単位データを逆検索テーブル142から検索し、検索した単位データについての情報を関連性情報キュー143および関連性集計テーブル144から消去する。なお、キャッシュ制御部132は、キャッシュ部122上のページUのデータを破棄する。そのとき、キャッシュ制御部132は、キャッシュ部122から明示的にページUのデータを消去してもよいし、ページUのデータを消去せずにページUに割り当てられていた記憶領域を上書き可能に設定するようにしてもよい。

【0102】

(S18) キャッシュ制御部132は、ステップS11で検索されたページTのデータ全体を、データ記憶部121からRAM102上のキャッシュ部122に読み出す。

(S19) アクセス実行部131は、キャッシュ部122に記憶されたデータに対して、受信したアクセス要求に応じたアクセスを実行し、アクセス要求を送信したクライアント装置に対して応答する。アクセス要求がリード要求である場合、アクセス実行部131は、アクセス要求で指定された単位データをキャッシュ部122から抽出して、アクセス要求を送信したクライアント装置に送信する。アクセス要求がライト要求である場合、アクセス実行部131は、アクセス要求で指定された単位データをキャッシュ部122上で更新し、アクセス要求を送信したクライアント装置に更新の成否を通知する。また、アクセス要求がライト要求である場合、アクセス実行部131は、逆検索テーブル142に登録されたページTの更新フラグを「1」に書き換える。

【0103】

(S20) アクセス実行部131は、受信したアクセス要求に応じて関連性情報を生成し、制御情報記憶部123に形成された関連性情報キュー143に保存する。関連性情報には、アクセス要求を送信したクライアント装置の識別情報と、アクセス要求で指定された単位データの識別情報が含まれる。また、関連性情報には、同じクライアント装置からの要求によって前回アクセスされた単位データの識別情報が含まれる。前回アクセスされた単位データは、例えば、そのクライアント装置についての直近の関連性情報を関連性情報キュー143から検索することで特定できる。また、前回アクセスされた単位データの識別情報がアクセス要求に付加されている場合、その識別情報を利用できる。

【0104】

また、アクセス実行部131は、生成した関連性情報を用いて関連性集計テーブル144を更新する。具体的には、アクセス実行部131は、関連性集計テーブル144において、今回アクセスされた単位データに対応する前回アクセスされた単位データの重みを1だけ加算する。また、アクセス実行部131は、制御情報記憶部123に記憶された出現履歴テーブル145を更新する。具体的には、アクセス実行部131は、出現履歴テーブル145において、今回アクセスされた単位データと前回アクセスされた単位データの組(関連データ対)に対応する最新の出現回数を1だけ加算する。

【0105】

図14は、データ再配置の手順例を示すフローチャートである。

(S30) 再配置制御部133は、制御情報記憶部123に形成された関連性情報キュー143に、新たな関連性情報が追加されたことを検出する。

【0106】

(S31) 再配置制御部133は、以下のステップS33~S40に示すデータ再配置の検討を前回行ってから所定時間以上経過したか判断する。前回のデータ再配置の検討から所定時間以上経過した場合、ステップS33に処理が進み、データ再配置の検討が開始される。所定時間以上経過していない場合、ステップS32に処理が進む。

【0107】

(S32) 再配置制御部133は、関連性情報キュー143に保存された関連性情報が示す関連データ対のうち、ページをまたがる関連データ対(異なるページに属する単位データの組)の数をカウントする。各単位データが属するページは、制御情報記憶部123に記憶された検索テーブル141を参照して特定することができる。そして、再配置制御

10

20

30

40

50

部 1 3 3 は、ページをまたがる関連データ対の数が所定の閾値以上であるか判断する。条件を満たす場合、ステップ S 3 3 に処理が進み、データ再配置の検討が開始される。条件を満たさない場合、データ再配置の検討は開始されない。

【 0 1 0 8 】

なお、図 1 4 では、データ再配置の検討を開始する開始条件として、ステップ S 3 1 , S 3 2 の 2 つの条件を用いることとした。ただし、ステップ S 3 1 , S 3 2 の何れか一方のみを開始条件として用いてもよい。また、ステップ S 3 1 , S 3 2 に代えて、または、ステップ S 3 1 , S 3 2 と合わせて、他の開始条件を用いてもよい。例えば、関連性情報キュー 1 4 3 に保存された関連性情報の量が閾値に達したことを開始条件としてもよい。

【 0 1 0 9 】

(S 3 3) 再配置案生成部 1 3 4 は、関連性情報キュー 1 4 3 から 1 つの関連データ対の情報を抽出する。抽出する関連データ対の情報は、例えば、関連性情報キュー 1 4 3 に記憶されているもののうち最も古いものとする。抽出した関連データ対の情報は、関連性情報キュー 1 4 3 から削除される。以下では、関連データ対が示す今回アクセスされた単位データを $m 1$ 、前回アクセスされた単位データを $n 1$ と表記することがある。

【 0 1 1 0 】

(S 3 4) 再配置案生成部 1 3 4 は、単位データ $m 1$ が属するページ M と単位データ $n 1$ が属するページ N との間の再配置案を 1 つ生成する。再配置案は、ページ M , N のページ ID と、一方のページから他方のページ (ページ M からページ N またはその逆) に移動する単位データのデータ ID とを用いて表現できる。再配置案生成の詳細は後述する。

【 0 1 1 1 】

(S 3 5) 実行可否判定部 1 3 5 は、ステップ S 3 4 で生成された再配置案に従ってデータ再配置を実行した場合の実行コストを算出する。実行コストは、ライトバックするページの増加量 \times 書き込み速度、として算出できる。

【 0 1 1 2 】

ライトバックするページの増加量は、図 5 , 6 で説明したように、キャッシュ部 1 2 2 上でのページ M , N の更新状況に基づいて算出することができ、「 0 」, 「 1 」, 「 2 」の何れかの値をとる。実行可否判定部 1 3 5 は、逆検索テーブル 1 4 2 に登録されたページ M , N の更新フラグを確認して、ページ M , N のうち更新フラグ = 1 であるページの数 (更新されたページの数) を算出する。データ再配置によって増加するライトバックのページ数は、「 2 」 - 更新されたページの数である。書き込み速度は、制御情報記憶部 1 2 3 に記憶されたパラメータテーブル 1 4 6 を参照して特定できる。

【 0 1 1 3 】

(S 3 6) 実行可否判定部 1 3 5 は、ステップ S 3 4 で生成された再配置案に従ってデータ再配置を実行した場合の不実行ペナルティを算出する。不実行ペナルティは、ページ M , N の間のカット数の減少量 \times 再出現期待値 \times 読み出し速度、として算出できる。

【 0 1 1 4 】

ページ M , N の間のカット数は、制御情報記憶部 1 2 3 に記憶された関連性集計テーブル 1 4 4 に登録されている関連データ対のうち、ページ M , N をまたがる関連データ対の重みの合計である。すなわち、ページ M , N の間のカット数は、ページ M , N の両方が今回キャッシュされている期間内に出現した関連データ対のうち、ページ M , N をまたがる関連データ対の出現回数を示す。実行可否判定部 1 3 5 は、関連性集計テーブル 1 4 4 を参照して、現在の配置状況におけるカット数とデータ再配置を実行した後の配置状況におけるカット数とを算出し、前者から後者を引いたカット数の減少量を算出する。

【 0 1 1 5 】

ページ M , N にまたがっていた関連データ対の中には、データ再配置によってページ M , N の何れか一方の中に収まり、アクセス性能が向上するものがあり得る。逆に、ページ M , N の何れか一方の中に収まっていた関連データ対の中には、データ再配置によってページ M , N にまたがるようになり、アクセス性能が低下するものがあり得る。カット数の減少量は、一部の関連データ対についてのアクセス性能の向上と一部の関連データ対につ

10

20

30

40

50

いてのアクセス性能の低下とを反映したものであり、データ再配置を行うことによる単位データの配置状況の全体的な改善度を表した指標であると言える。

【0116】

再出現期待値は、ページM, Nに属する単位データの間に関連データ対が今後一定期間の間に出現する回数の期待値を表し、パラメータテーブル146に登録された予測式を用いて算出される。例えば、実行可否判定部135は、ページM, Nに属する単位データの間に関連データ対を関連性集計テーブル144から抽出し、抽出した関連データ対それぞれについて、出現履歴テーブル145を参照して変数 x_1 , x_2 , x_3 , x_4 の値を算出する。そして、実行可否判定部135は、抽出した関連データ対それぞれについて、当該関連データ対に対応する個別の予測式を用いて個別の再出現期待値を算出する。また、実行可否判定部135は、抽出した関連データ対全体に対応する変数 x_1 , x_2 , x_3 , x_4 の平均値を算出し、全体の予測式を用いて全体の再出現期待値を算出する。個別の再出現期待値および全体の再出現期待値の平均を、不実行ペナルティの算出に用いる。

10

【0117】

なお、不実行ペナルティを算出するにあたり、個別の再出現期待値および全体の再出現期待値の一方のみを用いるようにしてもよい。また、予測式に代えて、予め算出した再出現期待値をパラメータテーブル146に登録しておくようにしてもよい。読み出し速度は、パラメータテーブル146を参照して特定できる。

【0118】

(S37) 実行可否判定部135は、ステップS35で算出した実行コストとステップS36で算出した不実行ペナルティとを比較し、不実行ペナルティが実行コストより大きいか判断する。不実行ペナルティが実行コストより大きい場合、再配置案を採用すると決定され、ステップS38に処理が進む。不実行ペナルティが実行コスト以下である場合、再配置案を採用しないと決定され、ステップS40に処理が進む。

20

【0119】

(S38) 再配置制御部133は、ステップS34で生成された再配置案に従ったデータ再配置を、キャッシュ部122上(RAM102上)で実行する。このとき、再配置制御部133は、RAM102上で単位データを移動してもよいし移動しなくてもよい。

【0120】

(S39) 再配置制御部133は、検索テーブル141および逆検索テーブル142を更新する。具体的には、再配置制御部133は、ページM, Nの間で移動する単位データの情報を検索テーブル141から検索し、その単位データに対応付けられたページIDを書き換える。また、再配置制御部133は、ページM, Nの情報を逆検索テーブル142から検索し、ページM, Nに対応付けられたデータIDを書き換える。また、再配置制御部133は、ページM, Nの再配置フラグを「1」に書き換える。

30

【0121】

(S40) 再配置案生成部134は、関連性情報キュー143から全ての関連データ対の情報を抽出したか、すなわち、関連性情報キュー143が空であるか判断する。関連性情報キュー143が空である場合、データ再配置の検討が終了する。関連性情報キュー143が空でない場合、ステップS33に処理が進む。

40

【0122】

次に、ステップS34で行われる再配置案の生成について説明する。以下では、再配置案の生成方法の例として、重心法とユニオンスプリット法を挙げる。

図15は、重心法によるデータ再配置の例を示す図である。

【0123】

重心法では、単位データの間に関連性の強さ(連続してアクセスされる可能性の高さ)を、N次元空間(Nは2以上の整数)上の距離として表現し、N次元空間上で単位データをグルーピングする。ここでは、一例として2次元空間を用いる。グラフ43は、関連性情報キュー143から抽出した関連データ対の情報を適用する前の関連性を表す。グラフ44は、抽出した関連データ対の情報を適用した後の単位データの間に関連性を表す。

50

【 0 1 2 4 】

重心法では、ページおよび単位データそれぞれに対して座標を付与する。ページの座標は、互いに十分に離れるように予め付与しておく。単位データの座標の初期値は、その単位データが属するページの座標の近傍になるように付与しておく。グラフ 4 3 では、ページ Q, R (ページ 3 2, 3 3) および単位データ e, f, g, h が配置されている。

【 0 1 2 5 】

初期状態では、所定のグルーピング方法を用いると、単位データ e, f はページ 3 2 と同じグループに振り分けられ、単位データ g, h はページ 3 3 と同じグループに振り分けられるようにしておく。グルーピング方法としては、例えば、各ページが順に、グループが決定していない単位データのうちそのページから座標が最も近い単位データを自グループに取り込むという方法が考えられる。グラフ 4 3 の場合、1 巡目でページ Q が単位データ f を選択し、ページ R が単位データ g を選択する。2 巡目でページ Q が単位データ e を選択し、ページ R が単位データ h を選択する。これにより、単位データ e, f はページ Q に属し、単位データ g, h はページ R に属するというグルーピングを行うことができる。

10

【 0 1 2 6 】

ここで、再配置案生成部 1 3 4 が関連性情報キュー 1 4 3 から関連データ対の情報を抽出すると、その関連データ対に応じて単位データの座標を変更する。具体的には、一方の単位データの座標を、他方の単位データが属するページの座標に近付ける。単位データ f の直後に単位データ g がアクセスされた場合、グラフ 4 3 では、単位データ f の座標がページ R の座標に近付き、単位データ g の座標がページ Q の座標に近づく。これは、単位データ f とページ R の関連性が現在よりも強くなり、単位データ g とページ Q の関連性が現在よりも強くなったことを表す。座標の移動量は、一定量としてもよい。また、座標の移動量は、単位データの座標と近づく先のページの座標との間の距離 (例えば、単位データ f の座標とページ R の座標の距離) に対する一定割合 (例えば、10%) としてもよい。

20

【 0 1 2 7 】

2 次元空間上で単位データの座標が変更されると、上記のグルーピング方法を用いて単位データのグループが再計算される。例えば、グラフ 4 4 の場合、1 巡目でページ Q が単位データ f を選択し、ページ R が単位データ h を選択する。2 巡目でページ Q が単位データ g を選択し、ページ R が単位データ e を選択する。これにより、単位データ f, g はページ Q のグループに振り分けられ、単位データ e, h はページ R のグループに振り分けられることになる。これは、単位データ e がページ Q からページ R に移動し、単位データ g がページ R からページ Q に移動するという再配置案を表す。

30

【 0 1 2 8 】

図 1 6 は、座標テーブルの例を示す図である。

再配置案の生成に重心法を用いる場合、座標テーブル 1 4 7 が制御情報記憶部 1 2 3 に記憶される。座標テーブル 1 4 7 は、ノード ID および座標の項目を有する。ノード ID は、N 次元空間上に配置するノードの識別情報である。ノード ID として、ページについてはページ ID を用い、単位データについてはデータ ID を用いる。ノード ID に対して、N 次元空間上の現在の座標が対応付けられる。単位データに対応する座標は、上記のように再配置案生成部 1 3 4 によって更新され得る。あるページがキャッシュ部 1 2 2 から追い出されても、そのページに関する情報を座標テーブル 1 4 7 から消去しなくてよい。

40

【 0 1 2 9 】

図 1 7 は、第 1 の再配置案生成の手順例を示すフローチャートである。

第 1 の再配置案生成は、上記のステップ S 3 4 で実行される。

(S 5 0) 再配置案生成部 1 3 4 は、制御情報記憶部 1 2 3 に記憶された検索テーブル 1 4 1 から単位データ m 1 を含むページ M と単位データ n 1 を含むページ N を検索する。

【 0 1 3 0 】

(S 5 1) 再配置案生成部 1 3 4 は、制御情報記憶部 1 2 3 に記憶された座標テーブル 1 4 7 から、単位データ m 1, n 1 およびページ M, N に対応する座標を検索する。

(S 5 2) 再配置案生成部 1 3 4 は、単位データ m 1 の座標をページ N の座標に向かっ

50

て近付ける。例えば、再配置案生成部 134 は、座標テーブル 147 で、単位データ m 1 の座標を、単位データ m 1 の座標とページ N の座標の距離が 10% 縮まるように変更する。また、再配置案生成部 134 は、単位データ n 1 の座標をページ M の座標に向かって近付ける。例えば、再配置案生成部 134 は、座標テーブル 147 で、単位データ n 1 の座標を、単位データ n 1 の座標とページ M の座標の距離が 10% 縮まるように変更する。

【0131】

(S53) 再配置案生成部 134 は、制御情報記憶部 123 に記憶された逆検索テーブル 142 から、ページ M, N に含まれる全ての単位データを検索する。再配置案生成部 134 は、座標テーブル 147 から、検索した単位データそれぞれの座標を検索する。

【0132】

(S54) 再配置案生成部 134 は、ステップ S54 で検索された単位データを、それら単位データの座標とページ M, N の座標を用いてグルーピングする。グルーピングでは、ページ M, N の座標と単位データそれぞれの座標との間の距離が考慮される。ページ M との距離が近い単位データはページ M に配置されることが好ましく、ページ N との距離が近い単位データはページ N に配置されることが好ましい。例えば、ページ M, N が交互に、未選択の単位データのうち距離が最も近い単位データを 1 つずつ選択していく。

【0133】

(S55) 再配置案生成部 134 は、現在のページ M, N のデータ配置とステップ S54 で求めたページ M, N のデータ配置とを比較し、ページ M, N の間で移動する単位データを特定する。これにより、ページ M, N の再配置案が生成される。

【0134】

図 18 は、ユニオンスプリット法によるデータ再配置の例を示す図である。

ユニオンスプリット法では、再配置案生成部 134 が関連性情報キュー 143 から関連データ対の情報を抽出すると、関連データ対が示す 2 つのページが統合される。ページの統合では、一方のページに属する全ての単位データを、他方のページに移動させる。統合後の一方のページは、単位データを含まない空のページとなる。

【0135】

ただし、統合後の他方のページに含まれる単位データの量が上限を超えてしまうことがある。その場合、単位データそれぞれのアクセス状況に応じて、統合後の他方のページを分割する。ページの分割では、他方のページに集められた単位データを、キャッシュ部 122 に今回キャッシュされている間にアクセスされたものとアクセスされなかったものとにグルーピングする。そして、何れか一方のグループの単位データを移動させる。

【0136】

例えば、単位データ d, e, f を含むページ 32 (ページ Q) と、単位データ g, h, i を含むページ 33 (ページ R) がキャッシュ部 122 にキャッシュされているとする。また、今回のキャッシュ中、単位データ e の直後に単位データ f がアクセスされ、単位データ f の直後に単位データ g がアクセスされたとする。すると、ページ Q とページ R が統合される。例えば、ページ R に含まれる単位データ g, h, i がページ Q に移動する。その結果、ページ Q は単位データ d, e, f, g, h, i を含み、ページ R は空となる。

【0137】

しかし、このように移動するとページ Q に含まれる単位データの量が所定の上限を超えてしまう場合、単位データ d, e, f, g, h, i が、キャッシュ中にアクセスされた単位データ e, f, g とアクセスされなかった単位データ d, h, i とに分けられる。そして、ページ Q が分割される。例えば、キャッシュ中にアクセスされなかった単位データ d, h, i がページ Q からページ R に移動する。その結果、ページ Q は単位データ e, f, g を含み、ページ R は単位データ d, h, i を含むこととなる。

【0138】

図 19 は、第 2 の再配置案生成の手順例を示すフローチャートである。

第 2 の再配置案生成は、上記のステップ S34 で実行される。

(S60) 再配置案生成部 134 は、制御情報記憶部 123 に記憶された検索テーブル

10

20

30

40

50

141から単位データm1を含むページMと単位データn1を含むページNを検索する。

【0139】

(S61)再配置案生成部134は、制御情報記憶部123に記憶された逆検索テーブル142から、ページM、Nに含まれる全ての単位データを検索する。

(S62)再配置案生成部134は、ページMとページNを統合する再配置案を生成する。具体的には、再配置案生成部134は、ページNに含まれる全ての単位データをページMに移動する再配置案を生成する。この再配置案によれば、ページNは空となる。

【0140】

(S63)再配置案生成部134は、ステップS62で生成した再配置案を採用した場合に、ページMのデータ量(例えば、単位データの個数)が所定の上限を超えるか判断する。ページMのデータ量が上限を超える場合、ステップS64に処理が進む。ページMのデータ量が以下である場合、ステップS66に処理が進む。

10

【0141】

(S64)再配置案生成部134は、ページMに集められた単位データそれぞれが、キャッシュ部122に今回キャッシュされている間にアクセスされたか判定する。各単位データのアクセスの有無は、例えば、制御情報記憶部123に記憶された関連性集計テーブル144に、その単位データに関する情報が登録されているか否かで判定できる。

【0142】

(S65)再配置案生成部134は、ステップS64で判定したアクセスの有無に応じてページMを分割するように、ステップS62で生成した再配置案を修正する。具体的には、再配置案生成部134は、ページMに集められた単位データのうち、アクセスされなかった単位データがページNに移動するように再配置案を修正する。

20

【0143】

(S66)再配置案生成部134は、ステップS62で生成した再配置案またはステップS65で修正した再配置案に基づいて、現在のページM、Nのデータ配置から移動する単位データを特定する。これにより、ページM、Nの再配置案が確定される。

【0144】

なお、再配置案生成部134は、重心法およびユニオンスプリット法を含む複数の再配置案の生成方法のうち、何れか1つを使用すればよい。使用する再配置案の生成方法は、例えば、ユーザが予め再配置案生成部134に設定しておく。重心法は、関連データ対の出現回数の増加に応じて徐々にデータ配置を変更していくことが可能な方法であり、データ配置の長期的な最適化に適しているという利点がある。ユニオンスプリット法は、新たな関連データ対の出現に反応して、データ配置を迅速に修正できるという利点がある。

30

【0145】

次に、上記のステップS36で算出するカット数の減少量について補足する。

図20は、データ再配置前後のカット数の変化例を示す図である。

ここでは、ページ32(ページQ)に単位データd、e、fが含まれ、ページ33(ページR)に単位データg、h、iが含まれているとする。また、単位データdと単位データg、単位データeと単位データf、単位データeと単位データg、単位データhと単位データiが、連続してアクセスされたとする。また、単位データfをページRに移動し、単位データgをページQに移動するという再配置案が生成されたとする。

40

【0146】

データ再配置前は、単位データd、gの関連データ対および単位データe、gの関連データ対が、ページQ、Rをまたがっている。よって、データ再配置前のカット数は「2」である。一方、生成された再配置案によれば、データ再配置後は、単位データd、gの関連データ対および単位データe、gの関連データ対がページQ、Rをまたがっておらず、単位データe、fの関連データ対がページQ、Rをまたがっている。よって、データ再配置後のカット数は「1」であり、カット数の減少量(Cut)が「1」と算出される。

Cutは、生成された再配置案の良否を反映していると言いうことができる。

【0147】

50

次に、再出現期待値の予測式を算出する方法について説明する。

図 2 1 は、回帰変数テーブルの例を示す図である。

パラメータ算出部 1 3 6 は、定期的またはユーザからの指示に応じて、出現履歴テーブル 1 4 5 を用いて個別の予測式および全体の予測式を算出し、パラメータテーブル 1 4 6 に登録する。予測式を算出するとき、パラメータ算出部 1 3 6 によって回帰変数テーブル 1 4 8 が制御情報記憶部 1 2 3 上に生成される。回帰変数テーブル 1 4 8 は、前日フラグ、前々日フラグ、過去出現率、季節および将来出現回数の項目を有する。

【 0 1 4 8 】

前日フラグは、回帰分析の説明変数であり、前述の変数 $\times 1$ に相当する。前日フラグは、ある関連データ対が基準日の前日に 1 回以上出現したか否かを示す。1 回以上出現した場合は前日フラグ = 1 となり、1 回も出現していない場合は前日フラグ = 0 となる。前々日フラグは、回帰分析の説明変数であり、前述の変数 $\times 2$ に相当する。前々日フラグは、ある関連データ対が基準日の前々日に 1 回以上出現したか否かを示す。1 回以上出現した場合は前々日フラグ = 1 となり、1 回も出現していない場合は前々日フラグ = 0 となる。

10

【 0 1 4 9 】

過去出現率は、回帰分析の説明変数であり、前述の変数 $\times 3$ に相当する。過去出現率は、基準日から所定日数前までのうち、ある関連データ対が 1 回以上出現した日の割合を示す。例えば、基準日から 10 日前までのうち、ある関連データ対が 3 日出現して 7 日出現しなかった場合、過去出現率は 0.3 となる。季節は、回帰分析の説明変数であり、前述の変数 $\times 4$ に相当する。春は「0」、夏は「1」、秋は「2」、冬は「3」と表記される。将来出現回数は、回帰分析の目的変数であり、前述の変数 y に相当する。将来出現回数は、基準日から所定日数後までの間に、ある関連データ対が出現した回数を示す。

20

【 0 1 5 0 】

パラメータ算出部 1 3 6 は、基準日を 1 つ選択し、出現履歴テーブル 1 4 5 に登録された関連データ対それぞれについて、基準日の前後の出現回数を用いて、前日フラグ・前々日フラグ・過去出現率・季節・将来出現回数を算出する。パラメータ算出部 1 3 6 は、複数の基準日について、前日フラグ・前々日フラグ・過去出現率・季節・将来出現回数を算出する。複数の基準日は、互いに一定日数以上離れるようにする。そして、パラメータ算出部 1 3 6 は、各関連データ対について、基準日の異なる説明変数および目的変数の値を用いて回帰分析を行い、当該関連データ対の個別の予測式の係数を算出する。また、パラメータ算出部 1 3 6 は、全ての関連データ対についての説明変数および目的変数の値をまとめて使用して回帰分析を行い、全体の予測式の係数を算出する。

30

【 0 1 5 1 】

図 2 2 は、パラメータ算出の手順例を示すフローチャートである。

(S 7 0) パラメータ算出部 1 3 6 は、複数の基準日を選択する。

(S 7 1) パラメータ算出部 1 3 6 は、説明変数と目的変数を決定する。例えば、説明変数として前日フラグ (変数 $\times 1$) ・前々日フラグ (変数 $\times 2$) ・過去出現率 (変数 $\times 3$) ・季節 (変数 $\times 4$) を用い、目的変数として将来出現回数 (変数 y) を用いる。過去出現率および将来出現回数については、集計期間の長さも決定する。ただし、関連データ対の過去の出現状況を表す説明変数として、上記以外の特徴量を用いることも可能である。

40

【 0 1 5 2 】

(S 7 2) パラメータ算出部 1 3 6 は、制御情報記憶部 1 2 3 に記憶された出現履歴テーブル 1 4 5 を用いて、異なる関連データ対と基準日の組み合わせ毎に、説明変数の値と目的変数の値を算出し、回帰変数テーブル 1 4 8 に登録する。

【 0 1 5 3 】

(S 7 3) パラメータ算出部 1 3 6 は、回帰変数テーブル 1 4 8 に登録された値を用いて回帰分析を行い、説明変数の係数 (重み) を算出する。例えば、パラメータ算出部 1 3 6 は、前日フラグの係数 u_1 、前々日フラグの係数 u_2 、過去出現率の係数 u_3 、季節の係数 u_4 を算出する。このとき、関連データ対毎に値を分けて回帰分析を行うことで、関連データ対毎の個別の係数が算出される。また、全ての関連データ対の値をまとめて用い

50

て回帰分析を行うことで、全体の係数が算出される。

【0154】

(S74)パラメータ算出部136は、ステップS73で算出した係数を含む予測式をパラメータテーブル146に保存する。予測式には、関連データ対毎の再出現期待値を求める個別の予測式と、平均の再出現期待値を求める全体の予測式とが含まれる。一例として、再出現期待値 $(y) = -3 \times \text{前日フラグ}(x1) - 1 \times \text{前々日フラグ}(x2) + 2 \times \text{過去出現率}(x3) + 0.03 \times \text{季節}(x4)$ という予測式が得られる。

【0155】

図23は、再出現予測式の変化の例を示す図である。

サーバ装置100において、同じ関連データ対の出現回数の時間分布やローカリティの持続傾向が変わると、予測式によって算出される再出現期待値と実際の出現回数とのずれが大きくなるおそれがある。その場合には、予測式を更新することが好ましい。例えば、図23において、回帰変数テーブル148に登録された上3つのサンプルは、将来出現回数 = 前日フラグ + 10 × 過去出現率という予測式にフィットする。一方、回帰変数テーブル148に登録された下2つのサンプルは、将来出現回数 = 前々日フラグ + 10 × 過去出現率という予測式にフィットする。これは、サーバ装置100における出現回数の時間分布やローカリティの持続傾向が変化した可能性を示している。

10

【0156】

次に、第2の実施の形態の情報処理システムの構成の変形例について説明する。上記では、サーバ装置100が集中的にデータを管理することとした。これに対し、複数のサーバ装置が分散してデータを管理することも可能である。

20

【0157】

図24は、他の情報処理システムの例を示す図である。

変形例に係る情報処理システムは、クライアント装置21a, 22aおよびサーバ装置100a, 100b, 100cを有する。クライアント装置21a, 22aおよびサーバ装置100a, 100b, 100cは、ネットワーク20に接続されている。

【0158】

サーバ装置100a, 100b, 100cは、複数のページのデータを分散して記憶する。例えば、サーバ装置100aがページ31のデータを記憶し、サーバ装置100bがページ32のデータを記憶し、サーバ装置100cがページ33のデータを記憶する。

30

【0159】

クライアント装置21a, 22aは、アクセスしたい単位データを記憶するサーバ装置を知っている場合、当該サーバ装置に対してアクセス要求を送信する。一方、クライアント装置21a, 22aは、アクセスしたい単位データを記憶するサーバ装置を知らない場合、サーバ装置100a, 100b, 100cの全てにアクセス要求を送信するようにしてもよいし、任意の1つのサーバ装置に対してアクセス要求を送信してもよい。前者の場合、アクセス要求で指定された単位データをもつサーバ装置のみ、アクセス要求の送信元に応答すればよい。後者の場合、アクセス要求を受信したサーバ装置は、アクセス要求で指定された単位データをもつサーバ装置にアクセス要求を転送する。サーバ装置100a, 100b, 100cは、ページとサーバ装置との対応関係の情報を保持している。

40

【0160】

アクセスの連続性を検出するため、サーバ装置100a, 100b, 100cは、受信したアクセス要求が指定する単位データのデータIDを相互に通知し合う。または、クライアント装置21a, 22aが、前回アクセスした単位データのデータIDをアクセス要求に付加する。これにより、サーバ装置100a, 100b, 100cそれぞれは、自装置に記憶された単位データについての関連性情報を収集できる。サーバ装置100a, 100b, 100cそれぞれは、収集した関連性情報を用いて、自装置が管理するページに関する再配置案を生成しデータ再配置を実行すればよい。データ再配置の相手ページが他のサーバ装置に存在する場合、サーバ装置間で単位データが移動される。

【0161】

50

第2の実施の形態の情報処理システムによれば、連続してアクセスされた単位データができる限り同じページに配置されるように、HDD103上の単位データの格納位置が動的に変更される。これにより、キャッシュミスヒットが減少して、ランダムアクセスの低速なHDD103からのデータの読み出しが削減され、データアクセスの性能を向上させることができる。また、RAM102にデータがキャッシュされている間にデータ再配置を行うことで、HDD103へのデータの書き込みを削減することができる。

【0162】

また、連続アクセスの発生に応じて再配置案が生成されたとき、生成された再配置案の実行コストと不実行ペナルティとが算出され、不実行ペナルティが実行コストより大きい場合のみ再配置案が採用される。実行コストには、HDD103への書き込みの増加量が反映される。不実行ペナルティには、再配置案の良否やサーバ装置100におけるローカリティの持続傾向が反映される。これにより、データ再配置の実行によってデータアクセスの性能がかえって低下することを抑制できる。例えば、データ再配置を行ってもページをまたがる連続アクセスの減少効果が小さい場合や、同じパターンの連続アクセスの出現回数が少ないと予想される場合、再配置案を採用しないと判定され得る。また、HDD103への書き込みの増加量が大きい場合、再配置案を採用しないと判定され得る。

10

【0163】

なお、前述のように、第1の実施の形態の情報処理は、コンピュータにプログラムを実行させることで実現できる。また、第2の実施の形態の情報処理は、クライアント装置21, 22（または、クライアント装置21a, 22a）やサーバ装置100（または、サーバ装置100a, 100b, 100c）にプログラムを実行させることで実現できる。

20

【0164】

プログラムは、コンピュータ読み取り可能な記録媒体（例えば、記録媒体113）に記録しておくことができる。記録媒体としては、例えば、磁気ディスク、光ディスク、光磁気ディスク、半導体メモリなどを使用できる。磁気ディスクには、FDおよびHDDが含まれる。光ディスクには、CD、CD-R（Recordable）/RW（Rewritable）、DVDおよびDVD-R/RWが含まれる。プログラムは、可搬型の記録媒体に記録されて配布されることがある。その場合、可搬型の記録媒体からHDDなどの他の記録媒体（例えば、HDD103）にプログラムをコピーして実行してもよい。

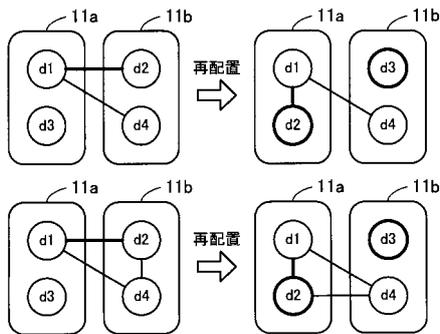
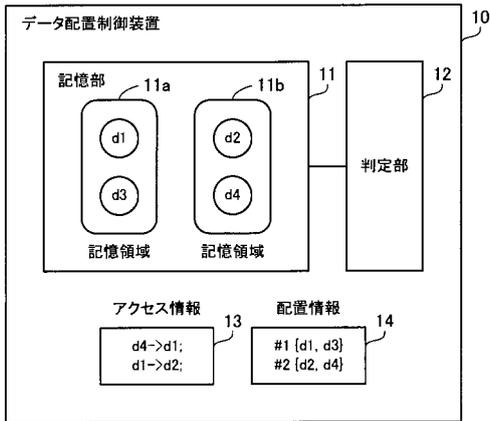
30

【符号の説明】

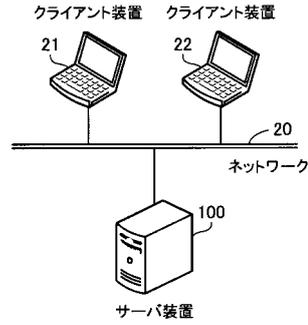
【0165】

- 10 データ配置制御装置
- 11 記憶部
- 11a, 11b 記憶領域
- 12 判定部
- 13 アクセス情報
- 14 配置情報
- d1, d2, d3, d4 単位データ

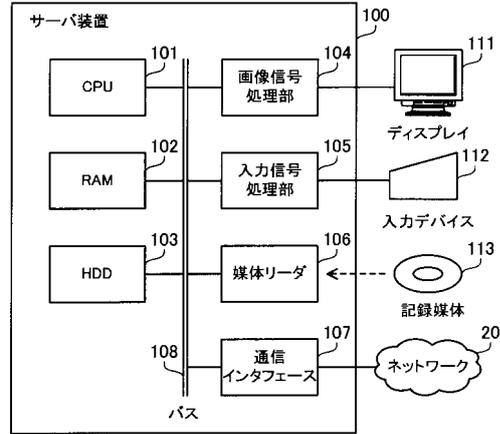
【 図 1 】



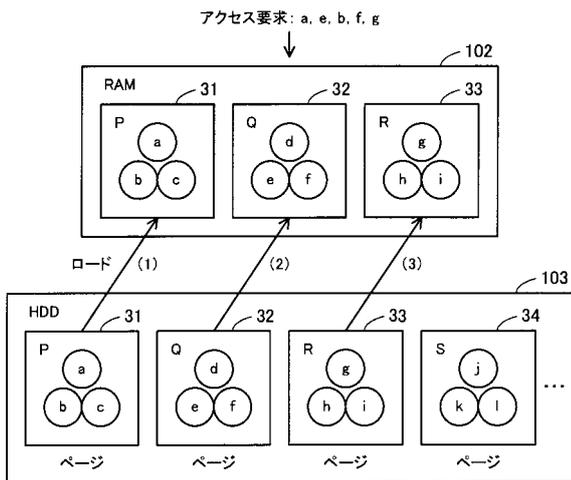
【 図 2 】



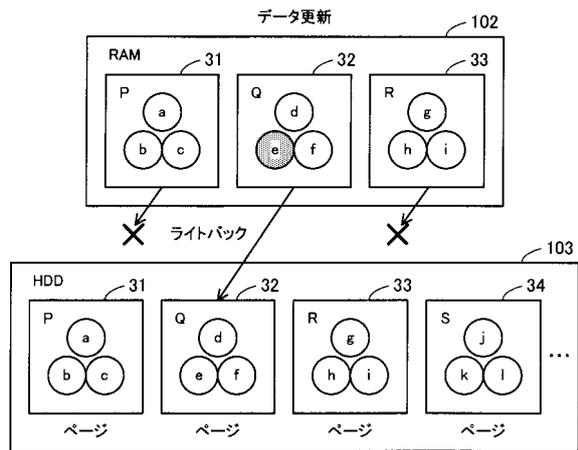
【 図 3 】



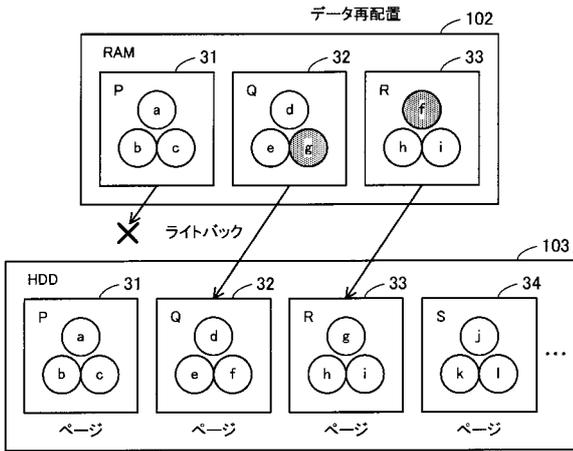
【 図 4 】



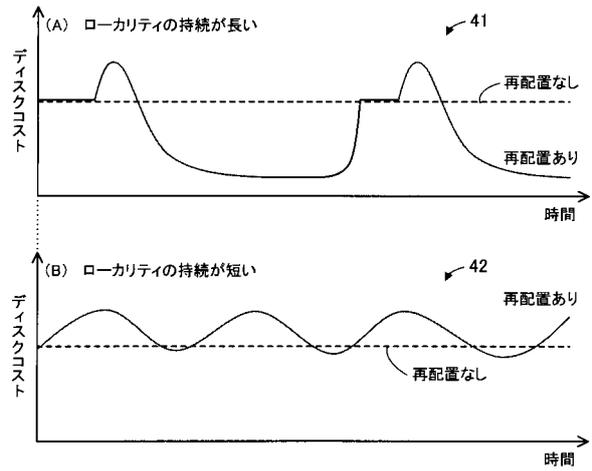
【 図 5 】



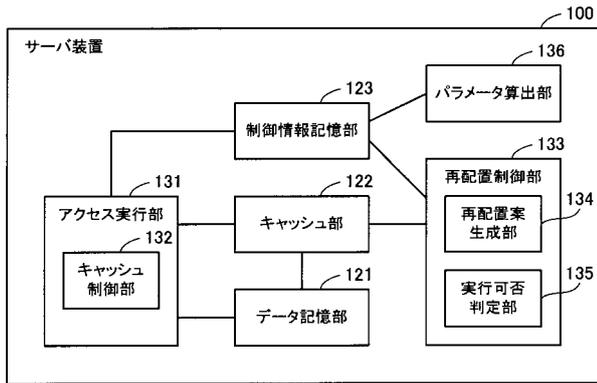
【 図 6 】



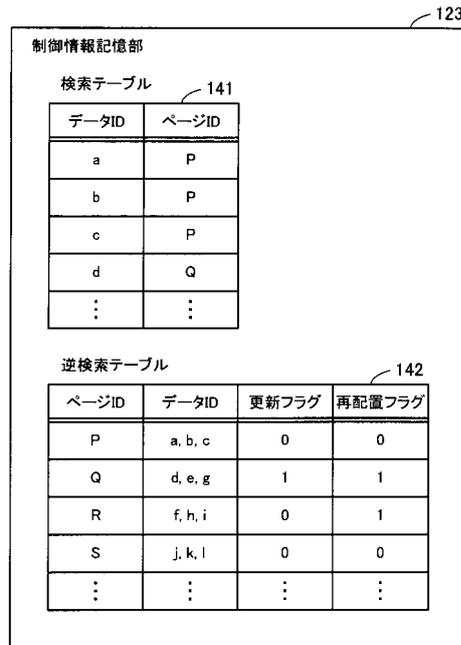
【 図 7 】



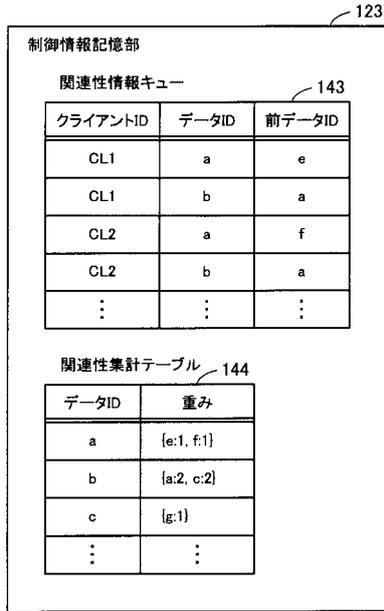
【 図 8 】



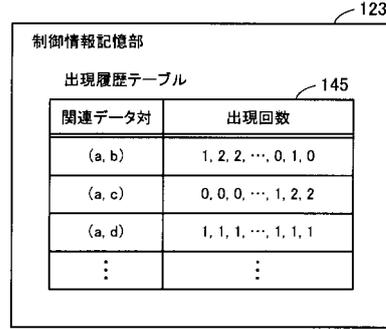
【 図 9 】



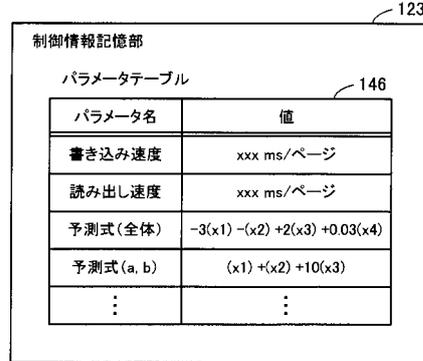
【図10】



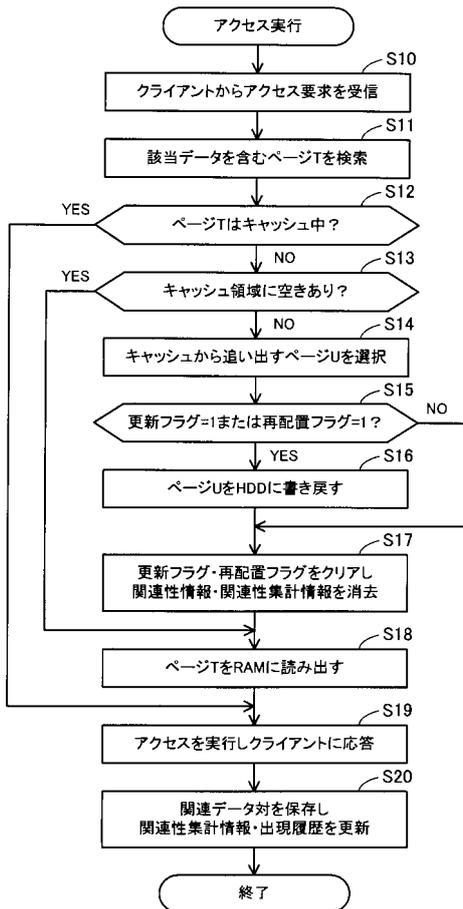
【図11】



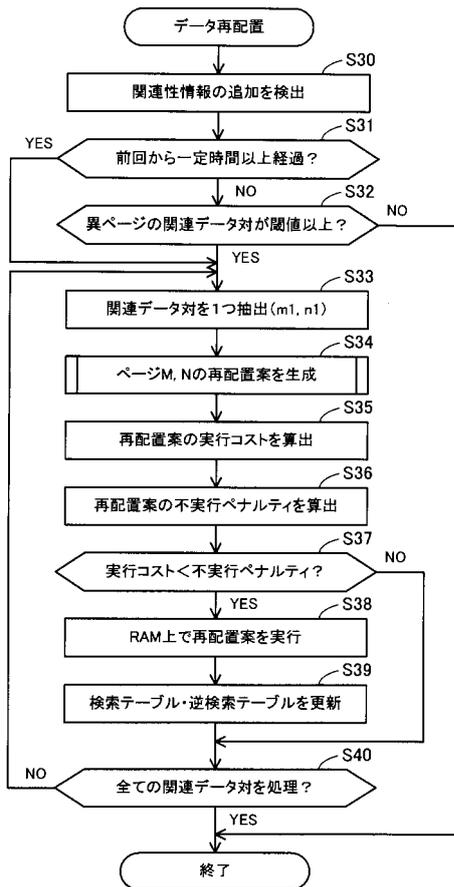
【図12】



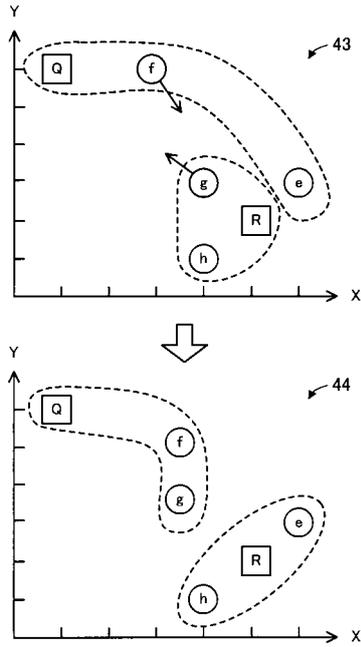
【図13】



【図14】



【図15】



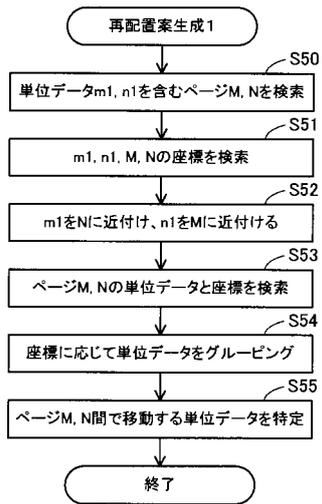
【図16】

制御情報記憶部 123

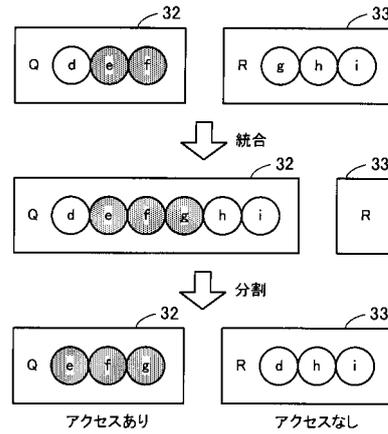
座標テーブル 147

ノードID	座標
Q	(1, 6)
R	(5, 2)
e	(6, 3)
f	(3, 6)
g	(4, 3)
h	(4, 1)
⋮	⋮

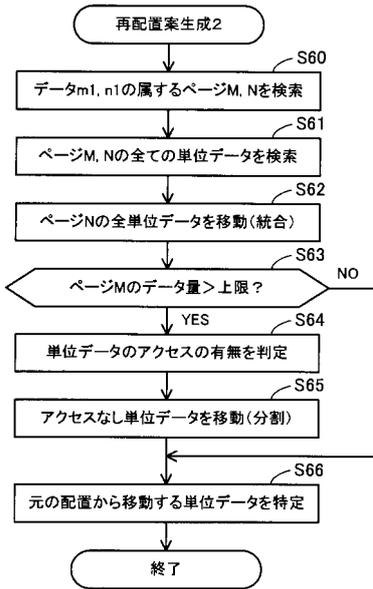
【図17】



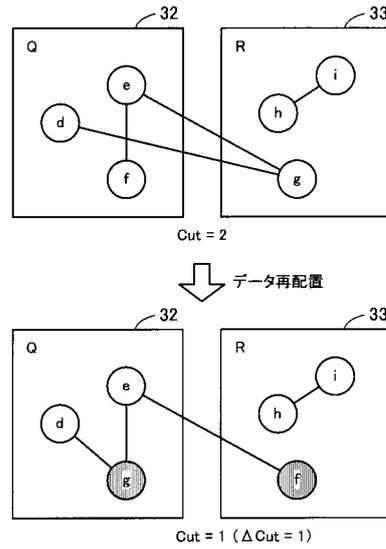
【図18】



【図19】



【図20】



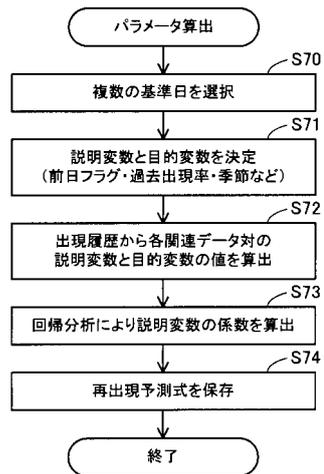
【図21】

制御情報記憶部 123

回帰変数テーブル 148

前日フラグ	前々日フラグ	過去出現率	季節	将来出現回数
0	1	0.3	3	5
1	0	0.5	2	7
1	1	0.7	0	10
0	1	0.2	0	3
1	1	0.8	1	12
⋮	⋮	⋮	⋮	⋮

【図22】



【 図 2 3 】

123

制御情報記憶部

148

回帰変数テーブル

前日フラグ	前々日フラグ	過去出現率	季節	将来出現回数
1	0	0.3	1	4
0	1	0.5	1	5
1	0	0.7	1	8
0	1	0.2	2	3
1	0	0.8	2	8
⋮	⋮	⋮	⋮	⋮

【 図 2 4 】

