

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3697149号

(P3697149)

(45) 発行日 平成17年9月21日(2005.9.21)

(24) 登録日 平成17年7月8日(2005.7.8)

(51) Int. Cl.⁷

F I

G06F 12/08

G06F 12/08 501F

G06F 3/06

G06F 12/08 509Z

G06F 12/12

G06F 12/08 551Z

G06F 12/08 557

G06F 3/06 301U

請求項の数 2 (全 14 頁) 最終頁に続く

(21) 出願番号 特願2000-294735 (P2000-294735)

(22) 出願日 平成12年9月27日(2000.9.27)

(65) 公開番号 特開2001-142778 (P2001-142778A)

(43) 公開日 平成13年5月25日(2001.5.25)

審査請求日 平成12年9月27日(2000.9.27)

(31) 優先権主張番号 09/410499

(32) 優先日 平成11年10月1日(1999.10.1)

(33) 優先権主張国 米国(US)

前置審査

(73) 特許権者 390009531

インターナショナル・ビジネス・マシー
ズ・コーポレーションINTERNATIONAL BUSIN
ESS MACHINES CORPO
RATIONアメリカ合衆国10504 ニューヨーク
州 アーモンク ニュー オーチャード
ロード

(74) 代理人 100086243

弁理士 坂口 博

(74) 代理人 100091568

弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 キャッシュ・メモリを管理する方法

(57) 【特許請求の範囲】

【請求項1】

ホスト・コンピュータと記憶装置との間でデータ・パスを介して接続され、かつ第1のキャッシュと第2のキャッシュとに論理的に分割されているキャッシュ・メモリを管理する方法であって、

前記第1のキャッシュがデータを含むセグメントを記憶し、前記第2のキャッシュが複数のセグメントからなるグループを記憶し、前記第1のキャッシュのためのセグメントLRU(最低使用頻度)リスト及び前記第2のキャッシュのためのグループLRUリストが設けられており、

前記ホスト・コンピュータへの読み取り時に、

(a) 前記ホスト・コンピュータからのデータにアクセスする要求に応じて、要求されたデータが前記セグメントLRUリストにあるかどうかを判断するステップと、

(b) 前記ステップ(a)において、もしも前記要求されたデータが前記セグメントLRUリストにあるならば、前記要求されたデータのコピーを前記第1のキャッシュから前記ホスト・コンピュータに転送するステップと、

(c) 前記ステップ(a)において、もしも前記要求されたデータが前記セグメントLRUリストにないならば、前記要求されたデータが前記グループLRUリストにあるかどうかを判断するステップと、

(d) 前記ステップ(c)において、もしも前記要求されたデータが前記グループLRUリストにあるならば、前記要求されたデータを含むセグメントのコピーを前記第2のキ

10

20

キャッシュから前記第1のキャッシュへ転送すると共に、前記セグメントLRUリストを更新して、前記転送されたセグメントのアドレスを前記セグメントLRUリストのMRU(最後に使用された)部分に入れるステップと、

(e)前記ステップ(d)に続いて、前記要求されたデータのコピーを前記第1のキャッシュから前記ホスト・コンピュータに転送するステップと、

(f)前記ステップ(c)において、もしも前記要求されたデータが前記グループLRUリストにないならば、前記要求されたデータを含むグループのコピーを前記記憶装置から前記第2のキャッシュへ転送すると共に、前記グループLRUリストを更新して、前記転送されたグループのアドレスを前記グループLRUリストのMRU部分に入れるステップと、

10

(g)前記ステップ(f)に続いて、前記要求されたデータを含むセグメントのコピーを前記第2のキャッシュから前記第1のキャッシュに転送するステップと、

(h)前記ステップ(g)に続いて、前記要求されたデータのコピーを前記第1のキャッシュから前記ホスト・コンピュータに転送するステップと

を行い、

前記ホスト・コンピュータからのデータの書き込み時に、

(イ)通常はスリープ・モードになっており、前記ホスト・コンピュータにより修正され且つ前記第2のキャッシュ及び前記記憶装置へ転送する必要があるセグメントが前記第1のキャッシュ内に存在することに応答して起動するステップと、

(ロ)前記ステップ(イ)に続いて、前記セグメントLRUリスト内の前記修正されたセグメントを見出すステップと、

20

(ハ)前記ステップ(ロ)に続いて、前記修正されたセグメントが前記グループLRUリストにあるかどうかを判断するステップと、

(ニ)前記ステップ(ハ)において、前記修正されたセグメントが前記グループLRUリストにあるならば、前記修正されたセグメントを前記第2のキャッシュに転送して既存のグループを修正するステップと、

(ホ)前記ステップ(ハ)において、前記修正されたセグメントが前記グループLRUリストにないならば、前記修正されたセグメントのセグメント待ちフラグが0に設定されているかどうかを判断するステップと、

(ヘ)前記ステップ(ホ)において、前記セグメント待ちフラグが0に設定されているならば、該セグメント待ちフラグを1に設定するステップと、

30

(ト)前記ステップ(ヘ)に続いて、前記ステップ(イ)、(ロ)、(ハ)及び(ホ)を行い、該ステップ(ホ)において前記セグメント待ちフラグが1であることを検出して、前記修正されたセグメントを前記記憶装置に転送するステップと

を行うことを特徴とする、キャッシュ・メモリ管理方法。

【請求項2】

もしも前記ステップ(a)で要求されたデータが前記第1のキャッシュに格納されていないと判断したならば、最低使用頻度手順にもとづいて最低使用頻度セグメントを前記第1のキャッシュからデステージするステップと、

もしも前記ステップ(a)で要求されたデータが前記第2のキャッシュに格納されていないと判断したならば、最小使用頻度手順にもとづいて前記セグメントの最小使用頻度グループを前記第2のキャッシュからデステージするステップと、

40

をさらに有することを特徴とする請求項1に記載のキャッシュ・メモリ管理方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

この発明は、ディスクメモリ等の直接アクセス記憶装置とキャッシュメモリとを有する記憶システムに関する。特に、記憶システムへのデータ・アクセスに対するホスト・コンピュータ及びチャンネルの待ち時間を減らすことによってパフォーマンスを向上させる方法及びシステムに関する。

50

【 0 0 0 2 】

【 従来 の 技 術 】

典型的なコンピュータ・システムは、ホスト・コンピュータと、直接アクセス記憶装置（D A S D）、例えば1枚以上の磁気ディスクとを含む。ホスト・コンピュータ上で実行されているアプリケーションは、データの読み取り及び書き出しのためにD A S D内のアドレス位置にアクセスする。そのようなアクセスは、ディスク入出力（I / O）オペレーションとして知られている。ホスト・コンピュータは、ディスクI / O操作が該I / O操作の完了を実行中のアプリケーションに待たせるように、D A S Dよりもかなり速い速度で稼動する。結果は、ホスト・コンピュータの処理能力が妨げられるということである。このことを避けるため、別の高速キャッシュ・メモリが用いられ、アプリケーションによつて最も頻繁に使用されるデータが格納される。

10

【 0 0 0 3 】

記憶システム・マネージャがディスクI / O操作の制御及びキャッシュ・メモリに対するアクセスの制御に使われる。一般に、データはレコード等のデータ対象物に編成される。実行中のアプリケーションがレコードを要求すると、記憶システム・マネージャは最初にキャッシュ・メモリ内のレコードを探索する。もし要求されたレコードがキャッシュ・メモリで見いだされるならば（すなわち、「ヒット」）記憶システム・マネージャは素早くそのレコードにアクセスするので時間を浪費するディスクI / O操作を実行する必要性がない。もし要求したレコードがキャッシュ・メモリで見いだされなければ（すなわち、「ミス」）、記憶システム・マネージャはディスクI / O操作を実行してD A S Dから要求されたレコードを取得し、該要求されたレコードをキャッシュ・メモリに書き込む。

20

【 0 0 0 4 】

一般に、記憶システム・マネージャは、最低使用頻度（L R U）手法によってキャッシュ・メモリ内のレコード保存を管理する。L R U手法は、制御ブロックの連鎖又は待ち行列を使用する。それぞれの制御ブロックが（a）レコードのアドレス、（b）連鎖内の次に続くレコードのアドレスを識別する順方向連鎖ポインタ、及び（c）その連鎖内の先のレコードのアドレスを識別する逆方向連鎖ポインタを識別する。記憶システム・マネージャは、L R Uレコード、例えば連鎖の先端を識別する第1のアンカー・ポインタを保持する。また、記憶システム・マネージャは最後に使用された（M R U）レコード、例えば連鎖の末端を識別する第2のアンカー・ポインタを保持する。

30

【 0 0 0 5 】

キャッシュ・ヒットが起こるたびに、ヒット・レコードの制御ブロックがデキューされ、続いてL R U連鎖の末端でM R Uレコードとして待ち行列に入れられる。キャッシュ・ミスが起こるたびに、L R U制御ブロックは連鎖の先端からデキューされる。新たに要求されたレコードはD A S Dからキャッシュ・メモリ内の割り当てられたアドレスへステージされる。デキューされた制御ブロックは、ステージされたレコードと割り当てられたアドレスとの一致によってアップデートされ、さらにM R U制御ブロックとしてL R U連鎖の末端で待ち行列に入れられる。

【 0 0 0 6 】

記憶システムのキャッシュ・メモリを設計する上で、要求されたデータ・レコードがキャッシュ・メモリ内で見つかる可能性を増やすことによりかなりの配慮が払われる。例えば、米国特許第5,717,893号は、グローバル・キャッシュと、各々が特定のタイプのデータ・レコードに割り当てられている複数のデステージング・ローカル・キャッシュとに分割されたキャッシュ・メモリを開示している。全てのタイプのデータ・レコードがローカル・キャッシュ又はディスク記憶システムからグローバル・キャッシュにデステージされる。L R Uアルゴリズムにもとづいて、L R Uデータ・レコードがグローバル・キャッシュからローカル・キャッシュへ降格される。この際、ローカル・キャッシュのデータ・タイプは降格されたL R Uレコードのデータ・タイプと一致する。ローカル・キャッシュが一杯になると、L R Uは記憶システムにデステージされる。より頻繁に使用されるデータ・レコード・タイプに対してより多くのキャッシュを割り当てるように分割スキームが

40

50

設計されるので、キャッシュ・ヒット率が増加する。より頻繁に使用されるデータ・タイプに対してキャッシュが増加し、使用頻度の低いデータ・タイプに対してはキャッシュが同時に減少するように、区画の論理的かつダイナミックなサイズ変更を可能とすることも特徴の一つである。

【0007】

他の従来スキームは、キャッシュ内のデータ・レコードの複製を削除することによって、キャッシュ・ヒット率を増やす。このタイプの典型的なスキームは、米国特許第5,802,572号及び第5,627,990号に開示されている。

【0008】

DASDシステムは、障害が生じた際にデータ回復を保証する幾何学的配置で構成された多数の小さな記憶装置モジュールを使用することで改善されている。改善がなされたそれらのシステムは、比較的安価な(又は独立した)複数のディスクからなる重複アレイ(RAID)としてしばしば呼ばれる。そのような幾何学的配置のいくつかでは、データ・オブジェクトが複数のデータ部分に分割され、かつ各データ部分が複数のディスクの異なるディスク上に格納される。RAIDレベル4として知られている一つの幾何学的配置では、ディスクの一つがデータ部分に対するパリティの保存専用となる。パリティは、障害が生じた際にデータ部分を再構築するために使われる。書き出し操作のために、この幾何学的配置は別々の2つの書き出しアクセスを必要とするもので、一つはデータ部分が格納されるディスクに対するアクセスであり、もう一つはパリティが格納されるディスクに対するアクセスである。

10

20

【0009】

RAIDレベル5として知られている別の幾何学的配置では、アレイ内の複数のアクティブなディスクをまたがってデータ及びパリティ情報の分配を行うためにディスクが分割される。各区画は一般にストライプと呼ばれる。ストライプに対するパリティ情報は通常一つのディスクに置かれ、データはストライプの残りのディスク上に置かれる。パリティ情報を含んでいるディスクはストライプ間で異なる。このことによって多数のストライプが並列処理可能となり、それによってステージ又はデステージされるデータのどちらかと言えば大きいチャンクを可能とする。

【0010】

キャッシュ・ヒット率を高めるための上記スキームは、どちらかと言うと小さなデータ・オブジェクト、例えばページ、テーブル等に関係している。また上記スキームは、より小さなページ・オブジェクトを大量に含むかなり大きなデータ・オブジェクト(例えばストライプ)を扱うためにRAIDシステムの能力を利用することはない。

30

【0011】

【発明が解決しようとする課題】

したがって、キャッシュ・ヒットの確率が改善されたキャッシュ・メモリが求められている。特に、RAID記憶装置が持つストライプ・アクセス能力を利用するキャッシュ・メモリが求められている。

【00012】

【課題を解決するための手段】

本発明は、記憶システムからデータ・オブジェクトを必要とするアプリケーションを実行するホスト・コンピュータを使用する。記憶システムは、記憶装置(例えば、ディスク記憶装置)とキャッシュ・メモリとを有する。ホスト・コンピュータによって頻繁に使われるデータ・オブジェクトがキャッシュ・メモリに格納される。データ・オブジェクトは、データ・オブジェクトが論理的に配列された記憶装置にも格納される。キャッシュ・メモリは、第1のキャッシュと第2のキャッシュとに論理的に分割されている。

40

【0013】

本発明の方法は、わずかな細分性(granularity)を持つセグメントを格納するための第1のキャッシュと、大きな細分性を持つセグメントからなるグループを格納するための第2のキャッシュとを使用する。ホスト・コンピュータがデータに対するアクセスを要求する

50

と、本発明の方法は要求されたデータが第1のキャッシュに格納されているかどうかを判断する。もし要求されたデータが第1のキャッシュに格納されていないならば、方法は要求されたデータが第2のキャッシュに格納されているかどうかを判断する。もし要求されたデータが第2のキャッシュに格納されていないならば、記憶装置に格納された複数のセグメントからなる一グループにアクセスする。この際、それらのセグメントの一つに要求されたデータが含まれている。次に、セグメントからなるグループを第2のキャッシュに格納し、要求されたデータが含まれる第1のセグメントを第1のキャッシュに格納する。次に、要求されたデータが第1のキャッシュからアクセスされる。

【0014】

もし方法が第2のキャッシュではなく第1のキャッシュに要求されたデータが格納されていると判断するならば、要求されたデータを含むセグメントのコピーが第1のキャッシュに転送される。

10

【0015】

本方法は、別々のLRU手順を用いて最小使用頻度セグメントを第1のキャッシュから、またセグメントからなるグループを第2のキャッシュからデステージして第1及び第2のキャッシュに格納されていない要求データに対して記憶域を割り当てる。

【0016】

第1のキャッシュにセグメントを格納し、第2のキャッシュにグループを格納するとともに、キャッシュ・メモリを第1のキャッシュ及び第2のキャッシュへ論理的に分割することは、本発明の重要な特徴である。この特徴は、一グループ内の一つのデータ・オブジェクトを要求するアプリケーションが同一グループの別のデータ・オブジェクトも要求するが、同一のセグメントである必要性はない可能性を利用する。

20

【0017】

本発明のキャッシュ・メモリ・システムは上記した方法の手順を含む多重細分性(multi-granular)キャッシュ・マネージャ・プログラムを使用する。

【0018】

本発明にもとづくメモリ媒体は、キャッシュ・メモリを制御して上記した発明の方法の手順を実行する。

【0019】

本発明の別の目的、利点、及び特徴は、図面と共に以下の説明を参照することによって理解することができよう。なお、同一の符号は同一の構成要素を示している。

30

【0020】

【発明の実施の形態】

図1は、本発明に適用されるコンピュータ・システム10の概略的構成を示すブロック図である。コンピュータ・システム10は、ホスト・コンピュータ12と記憶システム14とを含む。記憶システム14は、ホスト・アダプタ16、記憶装置18、及び本発明にもとづく多重細分性キャッシュ・メモリ30を備える。多重細分性キャッシュ・メモリ30はホスト・アダプタ16とともに、ホスト・コンピュータ12と記憶装置18との間のデータ・パスに結合している。

【0021】

40

ホスト・コンピュータ12は、アプリケーションを実行する1種類以上のプロセッサを一般に有する。プロセッサは、並行して単一のアプリケーション、又は異なる時間又は同時に別々のアプリケーション、あるいはそれらの任意の組み合わせを実行してもよい。アプリケーションは、論理データ構造を構成するデータ・オブジェクトを使う。論理データ構造は一般にページ、テーブル等のレコードであってもよい。格納を目的として、複数のデータ・オブジェクトが一つのセグメントに論理的に配置され、また複数のセグメントが論理的に一つのグループに配置される。

【0022】

RAID-5アレイに関する好ましい実施形態例では、グループはストライプ(アレイ・シリンダとして当業者にしばしば知られている)に対応し、またセグメントはストライプ

50

(また、アレイ・トラックとしてしばしば知られている)の一つのディスクに対応する。しかし、本発明は、オブジェクト、セグメント、及びグループの論理データ構造をフレキシブルな方法で規定することが可能であると考え。また、当業者はそれらの論理データ構造に対応する代わりに物理的レイアウトが可能であることを理解するであろう。

【0023】

ホスト・アダプタ16は、好ましくは記憶システム14とホスト・コンピュータ12との間でデータ・レコードを交換する送信及び受信能力を持つ従来のファシリティであってもよい。

【0024】

記憶装置18は、好ましくは複数のディスク19を有するディスク記憶装置であってもよい。好ましくは、記憶装置18は多重細分性キャッシュ・メモリ30がどちらかと言えば大きいチャンク(1つ以上のシリンダ)を操作するRAID能力を利用することができるように、RAID記憶スキームでディスク19を使用する。

10

【0025】

本発明によれば、キャッシュ・メモリ30は多重細分性マネージャ36を含む。

本発明の方法によれば、マネージャ36はキャッシュ・メモリ30を論理的に第1のキャッシュ32と第2のキャッシュ34とに分割する。マネージャ36は、ホスト・コンピュータ12からのデータ・アクセス要求に応答して、より小さな細分性のデータ構造、例えば第1のキャッシュ32に格納されるセグメントを生成する。マネージャ36は、ホスト・コンピュータ12からのデータ・アクセス要求に応答して、より大きな細分性のデータ構造、例えば第1のキャッシュ34に格納されるセグメントを生成する。第1及び第2のキャッシュ32、34は、キャッシュ・メモリ30の指定物理領域、又は論理領域であってもよい。

20

【0026】

セグメントを第1のキャッシュに格納し、グループを第2のキャッシュ34に格納するとともにキャッシュ・メモリ30を第1のキャッシュ32と第2のキャッシュ34とに論理分割することは、本発明の重要な特徴の一つである。この特徴は、一グループ内の一つのデータ・オブジェクトを要求するアプリケーションが同一グループの別のデータ・オブジェクトも要求するが、同一のセグメントである必要性はない可能性を利用する。キャッシュ・ミス率は約4分の1に減少すると推測される。

30

【0027】

マネージャ36は、セグメントLRUリスト35と、第1のキャッシュ32のためにセグメントLRUリスト35を管理する第1のキャッシュ・マネージャ37とを含むもので、従来のLRU手順を利用する。また、マネージャ36はグループLRUリスト38と、第2のキャッシュ34のためにグループLRUリスト38を管理するグループ・マネージャ39とを含むもので、従来のLRU手順を利用する。マネージャ36は、さらに読み込み手順50、セグメント書き出し手順60、及びグループ書き出し手順80を有する。

【0028】

ここで、図2を参照する。読み込み手順50は、ホスト・コンピュータ12が読み込み操作のためにデータ・アクセスを要求する時に開示される。ステップ51は、要求されたデータがセグメントLRUリスト35にあるかどうかを判断する。もし要求されたデータがセグメントLRUリスト35にあるならば、ステップ57は要求されたデータのコピーを第1のキャッシュ32からホスト・コンピュータ12に転送する。

40

【0029】

もし要求されたデータがセグメントLRUリスト35に含まれていない場合(第1のキャッシュ・ミス)、ステップ52は要求されたデータがグループLRUリスト38にあるかどうかを判断する。もし要求されたデータがグループLRUリスト38に含まれている場合、ステップ54は要求されたデータを含むセグメントのコピーを第2のキャッシュ34から第1のキャッシュ32へステージする。ステップ54は、また第1のキャッシュ・マネージャと協力してセグメントLRUリスト35をアップデートし、そのようなステージ

50

されたセグメントのアドレスをリストのMRU部分に入れる。次に、ステップ57は要求されたデータのコピーをホスト・コンピュータ12に転送する。

【0030】

もし要求されたデータがグループLRUリスト38にないならば(第2のキャッシュ・ミス)、ステップ55が要求されたデータを含むグループのコピーをディスク記憶装置18から第2のキャッシュ34へステージする。また、ステップ55は、第2のキャッシュ・マネージャと協力してグループLRUリスト38をアップデートし、そのようなステージされたグループのアドレスをリストのMRU部分に入れる。次に、ステップ54は要求されたデータを含むセグメントを第2のキャッシュ34から第1のキャッシュ32へ上記したようにステージする。次に、ステップ57は要求されたレコードのコピーをホスト・コンピュータ12に転送する。

10

【0031】

ここで、図3及び図4を参照しながら説明する。セグメント書き出し手順60は、図3に示す同期プロセス61と、図4に示す非同期プロセス70とを含む。最初に図3を参照すると、同期プロセス61はホスト・コンピュータ12が書き出し操作のためにデータ・アクセスを要求する時に開始される。ステップ62は、要求されたデータがすでにセグメントLRUリスト35に存在するかどうかを判断する。もしそうであるならば、ステップ64は既存のデータを修飾するためにセグメントLRUキャッシュ32に書き込まれるデータを転送する。また、ステップ64はセグメントLRUリスト35も更新する。さらにステップ64はこのセグメントのセグメント待ちフラグを0にセットする。

20

【0032】

もし要求データがセグメントLRUリスト35に存在しないとステップ62によって判断されるならば、ステップ63は第1のキャッシュ32に位置を割り当てる。次に、ステップ64が上記したようにして実行される。この際、書き込まれるデータは既に存在する版を修正しない。

【0033】

図4を参照する。非同期プロセス70は通常はスリープ・モードになっており、第2のキャッシュ34及び/又はディスク記憶装置18にデステージされる必要がある第1のキャッシュ32内に修飾されたセグメントが存在する場合にステップ71によって起こされる。ステップ71がデステージをする作業が存在すると判断した場合、ステップ72はセグメントLRUリスト35に次の修飾されたセグメントを見つける。ステップ73はこのセグメントがグループLRUリスト38にあるかどうかを判断する。もしなければ、ステップ74はセグメント待ちフラグが0に設定されたかどうかを判断する。もしそうであるならば、ステップ77はセグメント待ちフラグを1に設定する。次に、非同期プロセス70は、セグメントLRUリスト35の次の修飾されたセグメントに進むか、又はもし作業のための現在の操作が完了するならば短時間の間スリープ・モードに入る。このことは、ディスク記憶装置18にステージされる前に作られるセグメントに対する別の修飾を可能とする待ち機構を提供する。

30

【0034】

次にステップ71は再び非同期プロセス70を起こし、またステップ72からステップ74までが繰り返される。それによって、ステップ74はセグメント待ちフラグが0に設定されていないと判断する。次に、ステップ76は修飾されたセグメントをディスク記憶装置18にデステージする。そして、非同期プロセス70はセグメントLRUリスト35の次の修飾されたセグメントに進むか、又は作業に対する現行の操作が完了したならば短時間の間スリープ・モードに入る。

40

【0035】

もし修飾されたセグメントがグループLRUリスト28にあるとステップ73によって判断されるならば、ステップ73は既にあるグループを修飾するように、第2のキャッシュ34へ修飾されたセグメントを転送する。また、ステップ75はグループLRUリスト38を更新する。さらに、ステップ75はこのグループに対するグループ待ちフラグを0に

50

設定する。次に、非同期プロセス 70 がスリープ・モードに入る。

【0036】

ここで、図 5 を参照する。グループ書き出し手順 80 もまた非同期プロセスであり、グループをデステージする作業が存在する時にステップ 81 で起こされる。グループ書き出し手順 80 がステップ 81 によって起こされると、ステップ 82 はグループ LRU リスト 38 に次の修飾されたグループを見つける。ステップ 83 は、グループ待ちフラグが 0 に設定されているかどうかを判断する。もしそうであるならば、ステップ 84 はグループ待ちフラグを 1 に設定する。次に、グループ書き出し手順 80 は短時間の間、スリープ・モードに入る。このことは、ディスク記憶装置 18 にステージされる前に作られるセグメントに対する別の修飾を可能とする待ち機構を提供する。

10

【0037】

次に、ステップ 81 は再びグループ書き出し手順 80 を起こし、ステップ 82 及び 83 が繰り返される。それによって、やっとステップ 83 はグループ待ちフラグが 0 に設定されていないと判断する。次に、ステップ 85 は修飾されたセグメントをディスク記憶装置 18 にデステージする。次に、グループ書き出し手順 80 はスリープ・モードに入る。

【0038】

以上、本発明をその好ましい実施の形態にもとづいて説明したが、本発明は特許請求の範囲で定義されるように、種々の変更及び修飾が本発明の精神及び範囲から逸脱することなく可能であることは当業者に容易に理解されよう。

【0039】

まとめとして、本発明の構成に関して以下の事項を開示する。

20

(1) ホスト・コンピュータと記憶装置との間のデータ・パスに位置し、かつ第 1 のキャッシュと第 2 のキャッシュとに論理的に分割されているキャッシュ・メモリを管理する方法であって、

(a) 前記ホスト・コンピュータからのデータにアクセスする要求に応じて、要求されたデータが前記第 1 のキャッシュに格納されているかどうかを判断し、もし前記要求されたデータが第 1 のキャッシュに格納されていなければ、前記要求されたデータが第 2 のキャッシュに格納されているかどうかを判断するステップと、

(b) 前記要求されたデータが前記第 2 のキャッシュに格納されていないならば、前記記憶装置に格納されたセグメントからなるグループにアクセスし、前記セグメントのグループに含まれた第 1 のセグメントに前記要求されたデータが含まれるステップと、

30

(c) 前記第 2 のキャッシュに前記セグメントのグループを格納し、また前記第 1 のキャッシュに前記要求されたデータが含まれる前記第 1 のセグメントを格納するステップと、

(d) 前記ステップ (a) で前記要求されたデータが前記第 1 のキャッシュに格納されていると判断した場合、又は前記ステップ (c) で前記第 1 のセグメントを前記第 1 のキャッシュに格納した時に前記第 1 のキャッシュからの前記要求されたデータにアクセスするステップと、

を有することを特徴とするキャッシュ・メモリ管理方法。

(2) (e) もしステップ (a) が要求されたデータは第 1 のキャッシュではなく第 2 のキャッシュに格納されていると判断するならば、前記第 2 のキャッシュにおいて前記要求されたデータを含む第 2 のセグメントにアクセスし、前記第 2 のセグメントを前記第 1 のキャッシュに格納するステップと、

40

(f) 前記ステップ (e) が前記第 1 のキャッシュを前記第 2 のセグメントに格納した時に、前記第 1 のキャッシュから前記要求されたデータにアクセスするステップと、

をさらに有することを特徴とする上記 (1) に記載のキャッシュ・メモリ管理方法。

(3) (g) もし前記ステップ (a) で要求されたデータが前記第 1 のキャッシュに格納されていないと判断したならば、最低使用頻度手順にもとづいて最低使用頻度セグメントを前記第 1 のキャッシュからデステージするステップと、

(h) もし前記ステップ (a) で要求されたデータが前記第 2 のキャッシュに格納されていないと判断したならば、最小使用頻度手順にもとづいて前記セグメントの最小使用頻度

50

グループを前記第2のキャッシュからデステージするステップと、
をさらに有することを特徴とする上記(2)に記載のキャッシュ・メモリ管理方法。

(4)(i)前記ステップ(g)で使用される前記第1のキャッシュに対する第1の最小使用頻度リストを保持するステップと、

(j)前記ステップ(h)で使用される前記第2のキャッシュに対する第2の最小使用頻度リストを保持するステップと、

をさらに有することを特徴とする上記(2)に記載のキャッシュ・メモリ管理方法。

(5)(g)もし前記アクセス要求が前記要求されたデータの修復バージョンを書き出す要求であるならば、前記修復バージョンを前記第1のキャッシュに転送するステップと、

(h)前記記憶装置に対して前記第1のキャッシュから前記修復バージョンのコピーを直接的に、又は前記第2のキャッシュを介して間接的に転送するステップと、

10

をさらに有することを特徴とする上記(2)に記載のキャッシュ・メモリ管理方法。

(6)前記ステップ(h)は、もし前記要求されたデータが前記第2のキャッシュになれば直接的に転送を行い、またもし要求されたデータが前記第2のキャッシュにあれば間接的に転送を行うことを特徴とする上記(5)に記載のキャッシュ・メモリ管理方法。

(7)前記記憶装置は、RAID記憶装置であり、前記セグメントの各々は論理アレイ・トラックであり、また前記グループの各々は論理アレイ・シリンダであることを特徴とする上記(6)に記載のキャッシュ・メモリ管理方法。

(8)ホスト・コンピュータと記憶装置との間のデータ・パスに位置した多重細分性キャッシュ・メモリ・システムであって、

20

第1のキャッシュと第2のキャッシュとに論理的に分割されるキャッシュ・メモリと、
データにアクセスするために前記ホスト・コンピュータからの要求を処理するための多重細分性マネージャとを備え、さらに、
前記多重細分性マネージャは、

(a)前記ホスト・コンピュータからのデータにアクセスする要求に応じて、要求されたデータが前記第1のキャッシュに格納されているかどうかを判断し、もし前記要求されたデータが第1のキャッシュに格納されていないならば、前記要求されたデータが第2のキャッシュに格納されているかどうかを判断するステップと、

(b)前記要求されたデータが前記第2のキャッシュに格納されていないならば、前記記憶装置に格納されたセグメントからなるグループにアクセスし、前記セグメントのグループに含まれた第1のセグメントに前記要求されたデータが含まれるステップと、

30

(c)前記第2のキャッシュに前記セグメントのグループを格納し、また前記第1のキャッシュに前記要求されたデータが含まれる前記第1のセグメントを格納するステップと、

(d)前記ステップ(a)で前記要求されたデータが前記第1のキャッシュに格納されていると判断した場合、又は前記ステップ(c)で前記第1のセグメントを前記第1のキャッシュに格納した時に前記第1のキャッシュからの前記要求されたデータにアクセスするステップと、

を有する手順を実行することを特徴とする多重細分性キャッシュ・メモリ・システム。

(9)前記手順は、

(e)もしステップ(a)が要求されたデータは第1のキャッシュではなく第2のキャッシュに格納されていると判断するならば、前記第2のキャッシュにおいて前記要求されたデータを含む第2のセグメントにアクセスし、前記第2のセグメントを前記第1のキャッシュに格納するステップと、

40

(f)前記ステップ(e)が前記第1のキャッシュを前記第2のセグメントに格納した時に、前記第1のキャッシュから前記要求されたデータにアクセスするステップと、

をさらに有することを特徴とする上記(8)に記載の多重細分性キャッシュ・メモリ・システム。

(10)前記手順は、

(g)もし前記ステップ(a)で要求されたデータが前記第1のキャッシュに格納されていないと判断したならば、最低使用頻度手順にもとづいて最低使用頻度セグメントを前記

50

第1のキャッシュからデステージするステップと、

(h) もし前記ステップ(a)で要求されたデータが前記第2のキャッシュに格納されていないと判断したならば、最小使用頻度手順にもとづいて前記セグメントの最小使用頻度グループを前記第2のキャッシュからデステージするステップと、
をさらに有することを特徴とする上記(9)に記載の多重細分性キャッシュ・メモリ・システム。

(11) 前記手順は、

(i) 前記ステップ(g)で使用される前記第1のキャッシュに対する第1の最小使用頻度リストを保持するステップと、

(j) 前記ステップ(h)で使用される前記第2のキャッシュに対する第2の最小使用頻度リストを保持するステップと、

をさらに有することを特徴とする上記(10)に記載の多重細分性キャッシュ・メモリ・システム。

(12) 前記手順は、

(g) もし前記アクセス要求が前記要求されたデータの修復バージョンを書き出す要求であるならば、前記修復バージョンを前記第1のキャッシュに転送するステップと、

(h) 前記記憶装置に対して前記第1のキャッシュから前記修復バージョンのコピーを直接的に、又は前記第2のキャッシュを介して間接的に転送するステップと、

をさらに有することを特徴とする上記(9)に記載の多重細分性キャッシュ・メモリ・システム。

(13) 前記ステップ(h)は、もし前記要求されたデータが前記第2のキャッシュになれば直接的に転送を行い、またもし要求されたデータが前記第2のキャッシュにあれば間接的に転送を行うことを特徴とする上記(12)に記載の多重細分性キャッシュ・メモリ・システム。

(14) 前記記憶装置は、RAID記憶装置であり、前記セグメントの各々は論理アレイ・トラックであり、また前記グループの各々は論理アレイ・シリンダであることを特徴とする上記(6)に記載の多重細分性キャッシュ・メモリ・システム。

(15) ホスト・コンピュータと記憶装置との間のデータ・パスに位置した多重細分性キャッシュ・メモリ・システムを制御するためのメモリ媒体であって、

前記キャッシュ・メモリ・システムは、キャッシュ・マネージャと、第1のキャッシュと第2のキャッシュとに論理的に分割されるキャッシュ・メモリとを有し、

前記メモリ媒体は、

前記多重細分性キャッシュ・メモリを制御するための手段を有し、該制御は、

(a) 前記ホスト・コンピュータからのデータにアクセスする要求に応じて、要求されたデータが前記第1のキャッシュに格納されているかどうかを判断し、もし前記要求されたデータが第1のキャッシュに格納されていないならば、前記要求されたデータが第2のキャッシュに格納されているかどうかを判断するステップと、

(b) 前記要求されたデータが前記第2のキャッシュに格納されていないならば、前記記憶装置に格納されたセグメントからなるグループにアクセスし、前記セグメントのグループに含まれた第1のセグメントに前記要求されたデータが含まれるステップと、

(c) 前記第2のキャッシュに前記セグメントのグループを格納し、また前記第1のキャッシュに前記要求されたデータが含まれる前記第1のセグメントを格納するステップと、

(d) 前記ステップ(a)で前記要求されたデータが前記第1のキャッシュに格納されていると判断した場合、又は前記ステップ(c)で前記第1のセグメントを前記第1のキャッシュに格納した時に前記第1のキャッシュからの前記要求されたデータにアクセスするステップと、

を有する手順によることを特徴とするメモリ媒体。

(16) 前記手順は、

(e) もしステップ(a)が要求されたデータは第1のキャッシュではなく第2のキャッシュに格納されていると判断するならば、前記第2のキャッシュにおいて前記要求された

10

20

30

40

50

データを含む第2のセグメントにアクセスし、前記第2のセグメントを前記第1のキャッシュに格納するステップと、

(f) 前記ステップ(e)が前記第1のキャッシュを前記第2のセグメントに格納した時に、前記第1のキャッシュから前記要求されたデータにアクセスするステップと、
をさらに有することを特徴とする上記(15)に記載のメモリ媒体。

(17) 前記手順は、

(g) もし前記ステップ(a)で要求されたデータが前記第1のキャッシュに格納されていないと判断したならば、最低使用頻度手順にもとづいて最低使用頻度セグメントを前記第1のキャッシュからデステージするステップと、

(h) もし前記ステップ(a)で要求されたデータが前記第2のキャッシュに格納されていないと判断したならば、最小使用頻度手順にもとづいて前記セグメントの最小使用頻度グループを前記第2のキャッシュからデステージするステップと、
をさらに有することを特徴とする上記(16)に記載のメモリ媒体。

(18) 前記手順は、

(i) 前記ステップ(g)で使用される前記第1のキャッシュに対する第1の最小使用頻度リストを保持するステップと、

(j) 前記ステップ(h)で使用される前記第2のキャッシュに対する第2の最小使用頻度リストを保持するステップと、
をさらに有することを特徴とする上記(17)に記載のメモリ媒体。

(19) 前記手順は、

(g) もし前記アクセス要求が前記要求されたデータの修復バージョンを書き出す要求であるならば、前記修復バージョンを前記第1のキャッシュに転送するステップと、

(h) 前記記憶装置に対して前記第1のキャッシュから前記修復バージョンのコピーを直接的に、又は前記第2のキャッシュを介して間接的に転送するステップと、
をさらに有することを特徴とする上記(16)に記載のメモリ媒体。

(20) 前記ステップ(h)は、もし前記要求されたデータが前記第2のキャッシュになれば直接的に転送を行い、またもし要求されたデータが前記第2のキャッシュにあれば間接的に転送を行うことを特徴とする上記(19)に記載のメモリ媒体。

(21) 前記記憶装置は、RAID記憶装置であり、前記セグメントの各々は論理アレイ・トラックであり、また前記グループの各々は論理アレイ・シリンダであることを特徴とする上記(20)に記載のメモリ媒体。

【図面の簡単な説明】

【図1】 本発明にもとづく多重細分性キャッシュ・メモリを有するコンピュータ・システムの概略的構成を説明するためのブロック図である。

【図2】 本発明にもとづく多重細分性プログラムの読み込み手順を示すフローチャートである。

【図3】 本発明にもとづく多重細分性プログラムの書き出し手順を示すフローチャートである。

【図4】 本発明にもとづく多重細分性プログラムの書き出し手順を示すフローチャートである。

【図5】 本発明にもとづく多重細分性プログラムの書き出し手順を示すフローチャートである。

【符号の説明】

- 10 コンピュータ・システム
- 12 ホスト・コンピュータ
- 14 記憶システム
- 16 ホスト・アダプタ
- 18 記憶装置
- 19 ディスク
- 30 多重細分性キャッシュ・メモリ

10

20

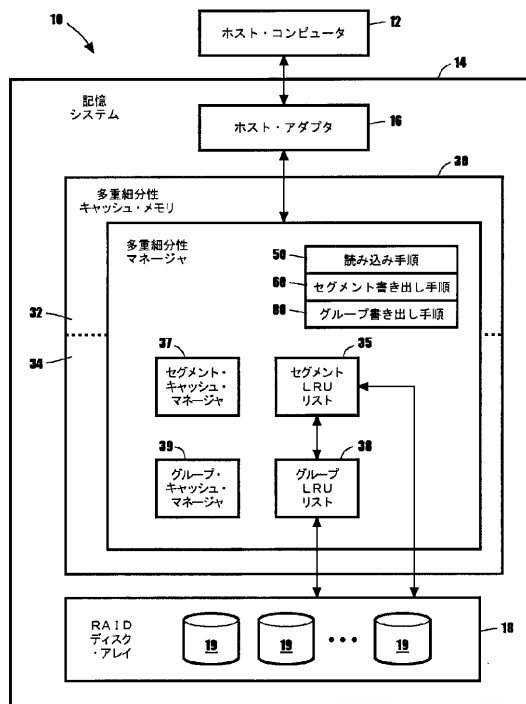
30

40

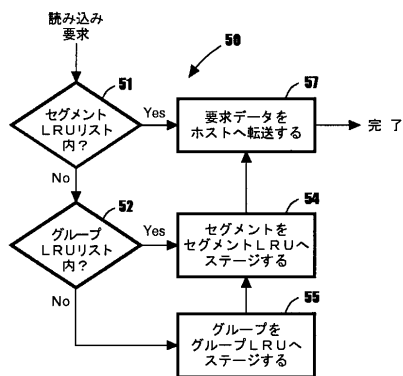
50

- 3 2 第 1 のキャッシュ
- 3 4 第 2 のキャッシュ
- 3 5 セグメントLRUリスト
- 3 6 多重細分性マネージャ
- 3 7 第 1 のキャッシュ・マネージャ
- 3 8 グループLRUリスト
- 3 9 グループ・マネージャ
- 5 0 読み込み手順
- 6 0 セグメント書き出し手順
- 6 1 同期プロセス
- 7 0 非同期プロセス
- 8 0 グループ書き出し手順

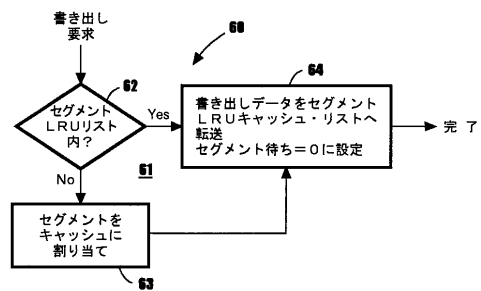
【 図 1 】



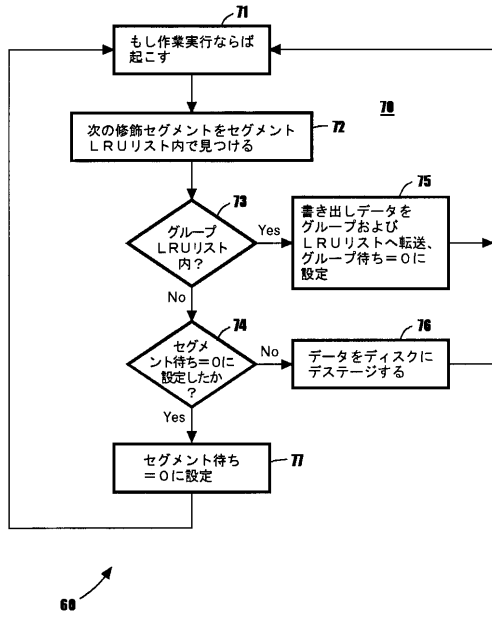
【 図 2 】



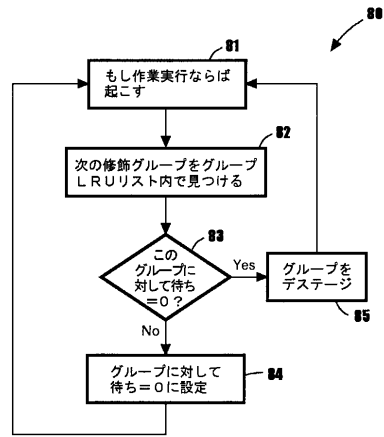
【 図 3 】



【 図 4 】



【 図 5 】



フロントページの続き

(51) Int.Cl.⁷

F I

G 0 6 F 3/06 3 0 2 A

G 0 6 F 3/06 5 4 0

G 0 6 F 12/12 5 5 7 Z

(72)発明者 ブルース・マクナット

アメリカ合衆国 9 5 0 2 0 カリフォルニア州、 ギルロイ、 デニオ・アベニュー 9 5

審査官 清木 泰

(56)参考文献 特開平 4 - 1 1 2 2 5 3 (J P , A)

特開平 6 - 2 4 3 0 0 3 (J P , A)

特開平 2 - 2 8 1 3 5 0 (J P , A)

(58)調査した分野(Int.Cl.⁷, D B名)

G06F12/08-12/12

G06F 3/06- 3/08