

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4246922号
(P4246922)

(45) 発行日 平成21年4月2日(2009.4.2)

(24) 登録日 平成21年1月16日(2009.1.16)

(51) Int.Cl.	F I
G06F 12/12 (2006.01)	G06F 12/12 551
G06F 12/08 (2006.01)	G06F 12/12 557Z
	G06F 12/08 507G
	G06F 12/08 543B
	G06F 12/08 523B

請求項の数 21 (全 15 頁)

(21) 出願番号	特願2000-618844 (P2000-618844)	(73) 特許権者	500271546 インテル・コーポレーション INTEL CORPORATION アメリカ合衆国 95051 カリフォル ニア、サンタ・クララ、ウォルシュ・アベ ニュー 2625
(86) (22) 出願日	平成12年5月17日(2000.5.17)	(74) 代理人	100088454 弁理士 加藤 紘一郎
(65) 公表番号	特表2002-544624 (P2002-544624A)	(72) 発明者	ウイルカーソン, クリストファー, ビイ アメリカ合衆国 オレゴン州 97209 ポートランド ホーリー・ストリート ノースウエスト 1807
(43) 公表日	平成14年12月24日(2002.12.24)		
(86) 国際出願番号	PCT/US2000/013569		
(87) 国際公開番号	W02000/070468		
(87) 国際公開日	平成12年11月23日(2000.11.23)		
審査請求日	平成17年1月12日(2005.1.12)		
(31) 優先権主張番号	09/314, 233		
(32) 優先日	平成11年5月18日(1999.5.18)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 保護機能付き最近最低使用頻度置換方法

(57) 【特許請求の範囲】

【請求項 1】

キャッシュメモリシステムにデータを記憶させる方法であって、

キャッシュメモリに記憶させたデータのリストを含むキャッシュデータエントリリストを作成し、

キャッシュデータエントリリストを再使用リストとフィルタリストに分割し、

新しいキャッシュデータエントリをフィルタリストに割り当てることによりフィルタリストにエントリを記入し、

昇進プロセスにより再使用リストにエントリを記入し、

選択されたキャッシュデータエントリを再使用リストとフィルタリストのいずれかから捨てるように動作する保護プロセスの作動によりフィルタリスト及び再使用リストからキャッシュデータエントリを排除するステップより成る、キャッシュメモリシステムへのデータ記憶方法。

【請求項 2】

キャッシュデータエントリの再使用リストとフィルタリストへの分割は、各キャッシュデータエントリに関連する再使用ビットを所定の値にセットすることにより行われる請求項 1 の方法。

【請求項 3】

再使用リストへのエントリの記入は、フィルタリストの一部である間にアクセスされたキャッシュデータエントリを再使用リストに昇進させることにより行われる請求項 1

の方法。

【請求項 4】

再使用リストへのエントリーの記入は、指名プロセスにより選択されたキャッシュデータエントリーを再使用リストに昇進させることにより行われる請求項 1 の方法。

【請求項 5】

再使用リストへのエントリーの記入は、キャッシュデータエントリーをフィルタリストから再使用リストへランダムに昇進させることにより行われる請求項 1 の方法。

【請求項 6】

キャッシュデータエントリーは、キャッシュメモリシステムに現在記憶されていないデータに対するリクエストであるキャッシュミスに限り応答して、フィルタリストから再使用リストへランダムに昇進される請求項 5 の方法。

10

【請求項 7】

再使用リストへのエントリーの記入は、さらに、フィルタリストから、新しく割り当てられたキャッシュラインをランダムに選択して再使用リストへ含まれるようにすることである請求項 1 の方法。

【請求項 8】

再使用リストまたはフィルタリストからのキャッシュデータエントリーの排除は、再使用リスト及びフィルタリストの各々の長さを計算し、再使用リスト及びフィルタリストの長さを所望の長さと比較し、キャッシュメモリシステムに記憶させる新しいキャッシュデータエントリーを受け、キャッシュデータエントリーを何れかのリストから排除して、新しいキャッシュデータエントリーの記憶場所を提供することより成り、キャッシュデータエントリーは、計算した長さと所望の長さとの比較結果に応じて選択されたリストから排除される請求項 1 の方法。

20

【請求項 9】

データを電子的に記憶させ、取り出すように作動可能で、キャッシュメモリを実現するコンピュータシステムであって、キャッシュメモリへのデータの記憶は、キャッシュメモリに記憶させたデータのリストを含むキャッシュデータエントリーリストを作成し、キャッシュデータエントリーリストを再使用リストとフィルタリストに分割し、新しいキャッシュデータエントリーをフィルタリストに割り当てることによりフィルタリストにエントリーを記入し、フィルタリストからキャッシュデータエントリーを昇進させることにより再使用リストにエントリーを記入し、選択されたキャッシュデータエントリーを再使用リストとフィルタリストのいずれかから捨てるように動作する保護プロセスの作動によりフィルタリスト及び再使用リストからキャッシュデータエントリーを排除するように行われるコンピュータシステム。

30

【請求項 10】

キャッシュメモリにデータを記憶させる方法であって、キャッシュメモリに記憶させたデータのリストを含むキャッシュデータエントリーリストを作成し、キャッシュデータエントリーリストを再使用リストとフィルタリストに分割し、キャッシュデータエントリーリストの各キャッシュデータエントリーの使用順序を追跡し、新しいキャッシュデータエントリーをフィルタリストに割り当てることによりフィルタリストにエントリーを記入し、フィルタリストからキャッシュデータエントリーを昇進させることにより再使用リストにエントリーを記入し、選択されたキャッシュデータエントリーを再使用リストとフィルタリストのいずれかか

40

50

ら捨てるように動作する保護プロセスの作動によりフィルタリスト及び再使用リストからキャッシュデータエントリーを排除するステップより成る、キャッシュメモリへのデータ記憶方法。

【請求項 1 1】

再使用リストへのエントリーの記入は、さらに、フィルタリストの一部である間にアクセスされたキャッシュデータエントリーを再使用リストに昇進させることにより行われる請求項 1 0 の方法。

【請求項 1 2】

再使用リストへのエントリーの記入は、さらに、キャッシュデータエントリーがフィルタリストにあるが、フィルタリストにおける最近使用頻度が最も高いキャッシュデータエントリーではない間に使用されると、フィルタリストから再使用リストに昇進させることにより行われる請求項 1 0 の方法。

10

【請求項 1 3】

再使用リストへのエントリーの記入は、さらに、キャッシュデータエントリーがフィルタリストにあるが、キャッシュデータエントリーリストにおける最近使用頻度が最も高いキャッシュデータエントリーではない間に使用されると、フィルタリストから再使用リストに昇進させることにより行われる請求項 1 0 の方法。

【請求項 1 4】

再使用リストへのエントリーの記入は、さらに、フィルタリストから、新しく割り当てられたキャッシュラインをランダムに選択して再使用リストへ含まれるようにすることである請求項 1 0 の方法。

20

【請求項 1 5】

フィルタリストから再使用リストへの昇進のための新しく割り当てられたキャッシュラインのランダムな選択は、再使用リストのサイズが大きくなるにつれてますます可能性が低くなる請求項 1 4 の方法。

【請求項 1 6】

保護プロセスは、さらに、

キャッシュメモリに現在記憶されていないデータに対するリクエストであるキャッシュミス、再使用リストの最近使用頻度が最も低いキャッシュデータエントリーが不変である間に、ストライクカウンタにおいてカウントし、

30

ストライクカウンタの値を最大許容ストライクしきい値と比較し、

ストライクカウンタの値が最大許容ストライクしきい値よりも大きい場合、キャッシュデータエントリーリストの最近使用頻度が最も低いキャッシュデータエントリーを排除するステップより成る請求項 1 0 の方法。

【請求項 1 7】

保護プロセスは、さらに、

再使用リストの最近使用頻度が最も低いキャッシュデータエントリーを参照する毎にカウンタを増分し、

フィルタリストの最近使用頻度が最も低いキャッシュデータエントリーを参照する毎にカウンタを減分し、

40

カウンタの値に基づき再使用リストの最大許容サイズを計算し、

再使用リストの実際のサイズを計算し、

再使用リストの実際のサイズが再使用リストの最大許容サイズより大きい場合、キャッシュデータエントリーリストの最近使用頻度が最も低いキャッシュデータエントリーを排除し、

フィルタリストの最近使用頻度が最も低いキャッシュデータエントリーが存在し、且つ再使用リストの実際のサイズが再使用リストの最大許容サイズより大きくない場合、フィルタリストのかかるキャッシュデータエントリーを排除し、

フィルタリストが空であり、再使用リストの実際のサイズが再使用リストの最大許容サイズより大きくない場合、キャッシュデータエントリーリストの最近使用頻度が最も低い

50

キャッシュデータエントリーを排除するステップより成る請求項 10 の方法。

【請求項 18】

保護プロセスはさらに、

カウンタの値をしきい値と比較し、

しきい値に対するカウンタの値に基づき再使用リストの最大許容サイズを発生させるステップより成る請求項 17 の方法。

【請求項 19】

保護プロセスは、さらに、

フィルタリストの最近使用頻度が最も低いキャッシュデータエントリーへの中間的な参照なしに再使用リストの最近使用頻度が最も低いキャッシュデータエントリーを連続して参照すると、1つおきの参照毎にカウンタを増分し、

再使用リストの最近使用頻度が最も低いキャッシュデータエントリーへの中間的な参照なしにフィルタリストの最近使用頻度が最も低いキャッシュデータエントリーを連続して参照すると、1つおきの参照毎にカウンタを減分し、

カウンタの値に基づき再使用リストの最大許容サイズを計算し、

再使用リストの実際のサイズを計算し、

再使用リストの実際のサイズが再使用リストの最大許容サイズより大きい場合、キャッシュデータエントリーリストの最近使用頻度が最も低いキャッシュデータエントリーを排除し、

フィルタリストの最近使用頻度が最も低いキャッシュデータエントリーが存在し、再使用リストの実際のサイズが再使用リストの最大許容サイズよりも大きくない場合、フィルタリストのかかるキャッシュデータエントリーを排除し、

フィルタリストが空であり、再使用リストの実際のサイズが再使用リストの最大許容サイズより大きくない場合、キャッシュデータエントリーリストの最近使用頻度が最も低いキャッシュデータエントリーを排除するステップより成る請求項 10 の方法。

【請求項 20】

キャッシュメモリにデータを記憶させるように作動可能な回路を有する集積回路であって、

キャッシュデータエントリーリストを再使用リストとフィルタリストに分割し、

キャッシュデータエントリーリストの各キャッシュデータエントリーの使用順序を追跡し、

新しいキャッシュデータエントリーをフィルタリストに割り当てることによりフィルタリストにエントリーを記入し、

フィルタリストからキャッシュデータエントリーを昇進させることにより再使用リストにエントリーを記入し、

選択されたキャッシュデータエントリーを再使用リストとフィルタリストのいずれかから捨てるように動作する保護プロセスの作動によりフィルタリスト及び再使用リストからキャッシュデータエントリーを排除するように動作する集積回路。

【請求項 21】

データを電子的に記憶させ、取り出すように作動可能で、キャッシュメモリより成るコンピュータシステムであって、

キャッシュメモリへのデータの記憶は、

キャッシュデータエントリーリストを再使用リストとフィルタリストに分割し、

キャッシュデータエントリーリストの各キャッシュデータエントリーの使用順序を追跡し、

新しいキャッシュデータエントリーをフィルタリストに割り当てることによりフィルタリストにエントリーを記入し、

フィルタリストからキャッシュデータエントリーを昇進させることにより再使用リストにエントリーを記入し、

選択されたキャッシュデータエントリーを再使用リストとフィルタリストのいずれかか

10

20

30

40

50

ら捨てるように動作する保護プロセスの作動によりフィルタリスト及び再使用リストからキャッシュデータエントリを排除するように行われるコンピュータシステム。

【発明の詳細な説明】

【0001】

【発明の分野】

本発明は、一般的に、電子コンピュータシステムのキャッシュメモリに記憶されたデータの置換に関し、さらに詳細には、望ましくない置換が起こらないように保護する最近最小使用頻度置換方法に関する。

【0002】

【発明の背景】

キャッシュメモリは、データへの高速アクセスのため、高いシステム速度を必要とする高速で高価なメモリを備えた高性能コンピュータシステムに使用されることが多い。プロセッサの動作速度が速ければ速いほど、プロセッサはメモリからデータ及び命令を高速で取り出す必要がある。このため、高速でアクセス可能なメモリが必要となる。このメモリは、非常に高いクロック速度が使用されることが多いため、プロセッサに近いところに配置するのが一般的である。しかしながら、キャッシュとして動作可能な高速メモリは高価であり、大きなまたは特注のメモリをプロセッサに近いところに配置するのは実際的でない場合が多い。従って、限られた数のキャッシュメモリのバンクを、大きい主システムメモリとは別個に、プロセッサのコア近くに配置することが多い。

【0003】

キャッシュメモリは、通常、スタティックランダムアクセスメモリ（SRAM）のような高速メモリと、キャッシュコントローラとより成る。コントローラは、プロセッサがどのデータを早期に必要とする可能性があるかの予測に応じて比較的低速の主メモリからキャッシュメモリへコピーされるデータを管理する。キャッシュメモリの容量は、システムの全メモリ容量の10%と1%の間であることが多いが、部分的にコンピュータシステムのメモリアクセス特性の予測可能性に応じて大きな範囲で変化する。

【0004】

通常、メモリへのアクセスはメモリアドレスの比較的小さな領域で次々に起こるため、アクセス頻度の最も高いデータをキャッシュメモリに記憶させると、システム性能を有意に改善することができる。使用頻度の最も高いデータが比較的高速のキャッシュメモリからアクセスされるようにすると、低速の主メモリからデータがアクセスされる間プロセッサが待機を余儀なくされることがなく、これは「キャッシュヒット」と呼ばれる。プロセッサが必要とするデータがキャッシュに存在せず、主メモリから取り出さなければならない場合、そのリクエストは同様に、「キャッシュミス」と呼ばれる。

【0005】

キャッシュがメモリへのアクセス速度を事実上増加させる度合は、キャッシュヒットとキャッシュミスとのメモリリクエスト数により測定可能である。キャッシュコントローラ設計者の目標は、プロセッサが必要とする可能性が最も高いデータをキャッシュに配置して、キャッシュミスに対するキャッシュヒットの比率を最大限にすることである。かかる方式を用いると、システムは、高速メモリを備える利点を最大限に活用すると共に大部分のデータを比較的 low 価格の低速メモリに記憶させてシステムの全体コストを減少させることができる。

【0006】

使用頻度の低いデータを使用頻度の高いデータに置換するかかる方式の1つは、キャッシュに記憶されたデータのうち最近の使用頻度が最低（LRU）であるデータに着目してそのデータをキャッシュにロードされる新しいデータと置換することである。かかる方式は通常、データが最後にアクセスされた時からの時間を示すデータを含む、キャッシュに記憶された各データブロックに関連する多数のLRUビットを使用する。かかる方式は、キャッシュのサイズに限りがあるため、再び必要になる前にデータを大きなデータセットから変位させる。さらに、LRUアルゴリズムは、最近の使用頻度の高いデータを、そのデ

10

20

30

40

50

ータの使用頻度が格段に低く、たとえただ1回使用されたものであったとしても、どちらかと言えば使用頻度の高いデータと置換する。従って、再びリクエストされる可能性の高いデータが、最近使用されたが使用頻度が低いデータに置換されないようにする方式が必要とされる。

【0007】

【発明の概要】

本発明は、キャッシュメモリシステムにデータを記憶させる方法を提供する。キャッシュデータエントリリストは、キャッシュメモリに記憶されるキャッシュデータエントリを記録するが、このリストは、フィルタリストと、再使用リストに分割されている。フィルタリストは、フィルタリストに新しいキャッシュデータエントリを全て記憶させることによりエントリが行われ、また再使用リストは、キャッシュデータエントリをフィルタリストから選択的に昇進させることによりエントリが行われる。選択されたキャッシュデータエントリを再使用リストとフィルタリストのいずれかから捨てるように動作する保護プロセスが作動すると、フィルタリストと再使用リストとから要素が排除される。

10

【0008】

【実施例の詳細な説明】

本発明の以下の詳細な説明において、本願の一部であり、本発明の特定の実施例を例示する添付図面を参照する。これらの実施例は、当業者が本発明を実施できるように十分に詳しく記載されている。他の実施例も可能であり、本発明の範囲から逸脱することなく構造的、論理的及び電氣的な設計変更を行うことができる。従って、以下の詳細な説明は限定的な意味に解すべきでなく、本発明の範囲は、頭書の特許請求の範囲と、特許請求の範囲が当然受けるべき均等物の全範囲とにより規定される。

20

【0009】

L R Uキャッシュメモリシステムを改良して、新たなキャッシュスペースが必要になると最近最低使用頻度のキャッシュデータエントリをただ排除する場合に生じる問題の影響を減少する方法が必要とされている。本発明は、単純なL R U置換方式を改良するものであり、キャッシュサイズに限りがあるため、再使用の前に大きなデータセットからデータを排除する問題に対処するものである。また、使用頻度の高いデータを使用頻度は低い最近使用されたデータと置換する問題に対処するものである。本発明は、L R Uキャッシュのこれら共通の問題に対する解決法を提供する。

30

【0010】

典型的なキャッシュメモリは、種々のバンクで同じ短縮アドレスを有する多数のメモリセルを備えており、これらのバンクは「ウェイ」と呼ばれる。例えば、アドレス長が16ビットのシステムのメモリは、4ウェイのキャッシュを備えており、各々が8ビットのアドレスを有する4つのバンクのメモリより成る。主メモリからのデータをキャッシュに記憶させる場合、データの主メモリアドレスの8つの最上位ビットが捨てられて、残りの8ビットがデータが記憶されるキャッシュアドレスとなる。所望のアドレスを有する4つのウェイ、またはメモリバンクが存在するため、最近最低使用頻度のデータを有するウェイは、通常、データが記憶されるウェイとして選択される。これは、他のメモリのウェイに記憶された使用頻度の高いデータの前に、最近最低使用頻度のデータが必要とされる可能性が低いことによる。この仮定は常に正しい訳ではなく、L R Uデータが実際、使用頻度が高い場合がある。多数の異なるキャッシュアーキテクチャ及びキャッシュ記入/置換アルゴリズムが存在するが、上記説明は、最近最低使用頻度(L R U)置換アルゴリズムを用いる1つのかかるマルチウェイキャッシュの一例にすぎない。

40

【0011】

本発明は、キャッシュメモリから排除されるデータが早期使用の可能性が最低ではないデータを含むウェイから選択される普通の方式である、容量による変位法及び関連性による変位法の問題を克服するように設計されたL R Uアルゴリズムを含む。

【0012】

50

容量による変位は、データの使用頻度が低いため、最近最低使用頻度のデータの早期アクセス可能性が、最近使用頻度の高いウェイのデータよりも実際に大きい場合に起こる。例えば、容量による変位は、16kBのキャッシュの一部が記憶される32kBの大きさのデータより成る画像に対する画像編集プロセスに、その画像を反復してアクセスする必要がある場合に、影響を与える。画像を構成するデータバイトが次々にアクセスされ、アクセスと共にキャッシュに記憶される場合、キャッシュには常に16kBの大きさの画像情報が存在するが、そのプロセスは、キャッシュに実際に記憶された情報を再び必要とする前に、もう1つの16kBの画像情報をアクセスする必要がある。画像ファイルの第1の部分の部分を再びアクセスする時、キャッシュは画像の第2の部分からの最近使用頻度の最も高いデータだけを含んでいる。同様に、プロセスが画像の第2の部分からのデータを必要とする時、そのデータは、画像の第1の部分からより最近アクセスされたデータと置換されている。かかるプロセスにはキャッシュミスが必ず大きな割合で起こるが、その理由は、キャッシュが反復使用される全てのデータを記憶する容量をただ備えていないためである。

10

【0013】

データがキャッシュから排除されるもう1つの普通の方式は、関連性による変位である。これは、ある特定のデータが、繰り返し必要とされるが、使用頻度の低い多量のデータがアクセスされるため、排除される時に起こる。例えば、データ構造Aと、データ要素Bとを取り扱うただ4つのエントリーを有するキャッシュを考える。データがA1、A2、A3、A4、A5、B、A6、A7、A8、A9、B、A10、A11、A12、A13、A14、Bなどの順序でアクセスされる場合、キャッシュはBを決して含まず、データ構造Aの最近使用頻度の高い要素を含む。Bはデータ構造Aの何れの単一要素より使用頻度は格段に高いが、従来のLRUアルゴリズムはBのようなアクセス頻度の高いデータは優先しない。

20

【0014】

LRUアルゴリズムを実現するためには、キャッシュコントローラは、ウェイ内の特定のアドレスについて何れのウェイが最近使用頻度が最も高いかを示すデータを処理しなければならない。キャッシュコントローラは、何れのウェイが最近使用頻度が最も高いかに関するデータを記憶させると共に、そのデータを復号して、何れのウェイが、最近使用頻度が最低であり、新しく使用されるデータと置換しなければならないかを決定する役割を有する。図1は、4ウェイキャッシュの特定のアドレスについて何れのウェイの最近使用頻度が最低かを追跡する1つの方式を部分的に示す表である。

30

【0015】

図1は、各ウェイ内の特定のアドレスの使用順序を記録するために使用可能な要素(i , j)の二次元マトリックスより成る表であり、この要素(i , j)は、ウェイ(i)のデータがウェイ(j)のデータより使用頻度が高い場合は1であり、ウェイ(i)のデータがウェイ(j)のデータより使用頻度が低い場合は0である。マトリックスの対角線上要素であるアドレス(0, 0)乃至(3, 3)は、(i)及び(j)が共に同じウェイを参照するため、無視できるから、意味のある情報を含まないことに注意されたい。また、この対角線の一方の側の要素は他方の側の要素の逆数であるため、(3, 2)が1で、ウェイ3がウェイ2よりも使用頻度が高いことを意味する場合、ウェイ2の最近使用頻度はウェイ3より高いことはありえず、(2, 3)は論理0でなければならない。従って、対角線の一方の側の要素は対角線と共に捨てることが可能であり、使用順序の記録に重要なマトリックスの数で表される要素だけを残す。従って、これら5個の数字を付したマトリックス要素が、5個のマトリックス要素(i , j)のアレイLで表される5ビットの列により表され、 $L = \{ (0, 1), (0, 2), (1, 2), (0, 3), (1, 3), (2, 3) \}$ であり、各ビットは、要素(i)が要素(j)よりもより最近使用頻度が高ければ1を含む。アレイLのこれら5つのビットは、単一の数で容易に参照することができるため、 $L = (0, 1, 2, 3, 4, 5)$ であり、各ビットはアレイLの前のマトリックス表示で与えられる順序のマトリックス要素に対応する。

40

50

【 0 0 1 6 】

図2は、新しいデータにアクセスし、新しくアクセスしたデータをキャッシュに書き込む時、アレイLの5つのビットを更新するためにキャッシュのこの例の動作に用いる更新表を示す。この更新マトリックスは、新しくアクセスするデータを記憶させるために何れのウェイを使用するかに応じて、Lの何れのビットを特定のアドレスに対して変化させるかを示す。キャッシュは4ウェイキャッシュであり、4つのウェイのうちの任意の1つがデータを記憶させるために選択可能であるため、マトリックスの1つの軸はキャッシュシステムの4つのウェイを表す。もう1つの軸は、4つのウェイのうち最近使用頻度が最も高いウェイにデータを記憶させ、もう1つの軸に表されるようにアクセスされるウェイに対応する、アレイLの5つのビットを表す。例えば、ウェイ2がワードLのチェックにより特定のアドレスについて最近使用頻度が最も低いと判定される動作では、ウェイ2は入来データと置換される。ウェイ軸のウェイ2を参照して、マトリックスは、ビット0、3及び4は不変であるが、ビット1及び2は0にセットされ、ビット5は1に変化していることを示す。アレイLのマトリックス要素表示法を再び参照して、その式は、ビット1及び2を0にセットすると、ウェイ0がウェイ2よりも最近使用頻度が低く、ウェイ1がウェイ2よりも最近使用頻度が低いことを示す。同様に、ビット5を1にセットすると、ウェイ2がウェイ3よりも最近使用頻度が高いことを示す。アレイLのマトリックス表示法をさらに検討すると、ビット0、3及び4はウェイ2についての情報を含まず、そのため不変であることがわかる。キャッシュへの書き込みの毎にこの方法に従ってビットを更新することにより、アレイLは、キャッシュのウェイの使用順序についての正確な情報を常に含むことになるため、キャッシュコントローラが新しいデータをキャッシュの最近使用頻度が最も低いウェイに割り当てるために使用することができる。

10

20

【 0 0 1 7 】

キャッシュコントローラは、図2の更新表の補足物に過ぎない犠牲者選択表を用いてこれを行う。図3は、マトリックスの全ての1が0で置換され、全ての0が1で置換され、全ての不変ビットが不変のままであるような、この例のキャッシュの犠牲者選択表を示す。マトリックスから読み取られる要素の補足物の発見は、エレクトロニクスまたは他の手段により全く容易に行うことが可能であるため、犠牲者選択表を維持する必要はないが、図2の更新表からデータを取り出しながら変換することによって行うことができる。図3の各要素は動作時にマスクとして使用されるため、マスクの内容がアレイLに記憶されたデータと比較される。4つのマスクのうちの任意のもののデータがアレイLのデータとマッチする場合、そのマスクに関連するウェイが、最近使用頻度が最も低く新しいデータの目標となるウェイとして選択される。図3のマスクが、ウェイ間で互いに排他的なビットデータを含むため、4つのマスクのうち1つだけがアレイLの任意可能な組み合わせのビットにマッチし、特定のマスクを最近使用頻度が最も低いものとして指示することに注意されたい。これは、マスクとアレイLとは同時に比較可能であることを意味し、ただ1つのウェイだけがキャッシュの書き込み動作のための目標として選択される。

30

【 0 0 1 8 】

例えば、特定のアドレスのウェイが2、3、0、1の順序で使用されており、ウェイ1が最近使用頻度が最も高く、ウェイ2が最近使用頻度が最も低い場合を考える。図2の更新表を用いると、全ての更新を行った後、アレイLはデータL = (0, 1, 1, 1, 1, 0)を含む。このワードを図3のマスクと比較すると、ウェイ0では、ビット0、1及び3が不正確であり、ウェイ1では、ビット1及び3が不正確であり、ウェイ2では、ビット1、2、5が不正確であることがわかる。ウェイ2では、犠牲者選択マトリックスのマスクのビット1、2及び5はアレイLのデータとマッチするため、ウェイ2は最近使用頻度が最も低いウェイとして適正に選択される。

40

【 0 0 1 9 】

前の段落で説明した例は、キャッシュシステムにおけるLRU置換の一例であり、前述したような普通のLRU置換方式の短所を有する。本発明は、上述した方式を改良するもので、関連性及び容量による変位のありふれた問題に対処するものである。

50

【 0 0 2 0 】

保護機能付きLRU、即ち、LRUPは、従来のLRUウェイリストを少なくとも2つのサブリストに分割したものである。1つの実施例では、従来のLRUウェイリストを改良するために、各ウェイが、そのウェイに記憶されるデータが再使用サブリストまたはフィルタサブリストの一部であることを示す再使用/フィルタビットを含むようにする。LRUの順序は大域的に維持されるため、LRUウェイから最近使用頻度が最も高いウェイまでの任意のウェイは、大域的リストに加えて、再利用リストまたはフィルタリストのいずれか1つの一部である。再使用リスト及びフィルタリストのデータを別個に管理する種々のプロセスは、その後、変位の問題に対処し、LRUPキャッシュの性能を最適化するように調整される。

10

【 0 0 2 1 】

本発明の1つの実施例において、キャッシュ要素を、フィルタリストの一部として最初に割り当てる。フィルタリストの選択要素を、指名プロセスにより再使用リストに昇進させ、再使用リストへエントリとして記入する。指名によりキャッシュラインがフィルタリストから再使用リストへ移動するため、指名を引き続き行くと、再使用リストが無限に大きくなって、最終的にキャッシュ要素がフィルタリストから押し出されることになる。従って、保護プロセスは、入来データを記憶させるために新しいキャッシュラインを割り当てると、再使用リストまたはフィルタリストの間で排除するメンバーの選択を行うことにより、フィルタリストと再使用リストの間で要素が適正に分配されるようにする。

20

【 0 0 2 2 】

フィルタリストから再使用リストへの要素の指名は、望ましくない変位を減少させるように調整された多数のプロセスにより行われる。関連性による変位を減少させるトリガーによる指名は、記憶されるデータがフィルタリストにある間に再使用されるなら、キャッシュラインをフィルタリストから再使用リストへ移動させることを含む。容量による変位を減少させるランダムな指名では、再使用リストに含まれるようにするために、フィルタリストから、新しく割り当てられたキャッシュラインがランダムに選択される。トリガーによる指名とランダムな指名との組み合わせ及び変更は、本発明の範囲内に含まれると考えられ、かかるものの幾つかの実施例を以下において説明する。

【 0 0 2 3 】

トリガーによる指名は、ある特定の要素を再使用リストに昇進させるのがそれほど望ましくない特定の状況を評価することにより、フィルタリストにある間使用されるキャッシュラインを再使用リストに昇進させるという簡単な規則を変更するものである。例えば、大域リストの最近使用頻度の最も高い位置にある間に参照される要素は、単一の演算または計算の一部として再アクセスされるため、2つの異なる時間にアクセスされる要素より再アクセスの可能性が低い。これらの要素は、その後、昇進から除外されるが、後で再使用または排除されるまでは、フィルタリストの一部として留まる。同様に、フィルタリストの最近使用頻度が最も高い位置にある間に参照される要素は、トリガーによる指名の間、再使用リストへの昇進から除外されるか、または、大域リストまたはフィルタリストの何れかの最近使用頻度の高い多数の位置にある要素が指名から除外される。

30

【 0 0 2 4 】

ランダムな指名は、ランダムに選択したキャッシュラインを再使用リストへ指名することを含む。通常ランダムな指名の変形例は、大域リストに対する再使用リストのサイズをモニターし、再使用リストが大きくなるにつれてランダムな指名をますます可能性の低いものとするランダムな選択アルゴリズムを用いることを含む。上述の例に加えて、データをランダムに選択するための多数の他の方法が存在するが、これらは本発明の範囲内に含まれるものと考えられる。

40

【 0 0 2 5 】

トリガーによる指名とランダムな指名との組み合わせを用いて、関連性による変位及び容量による変位の両方に対処することが可能であり、これら両方のアルゴリズム及びそれらの変形例を、特定のキャッシュ構成についてキャッシュヒットとキャッシュミスの比率

50

を最適化するように選択する。例えば、本発明との間に一貫性がある1つの指名プロセスの実施例は、再使用リストのサイズが増大するにつれてランダムな指名の可能性がますます低くなるランダムな指名アルゴリズムと共に、フィルタリストの最近使用頻度の最も高い要素を排除するトリガーによる指名方式を使用する。かかる組み合わせによると、トリガーによる指名アルゴリズムによる指名と、ランダムな指名アルゴリズムによる指名とが自動的に平衡するという利点があるが、その理由は、ランダムな指名アルゴリズムが再使用リストのサイズに依存するからである。これにより、再使用リストのサイズの安定化が得られると共に、関連性による変位と容量による変位の両方への対処を依然として効果的に行うことができる。

【0026】

再使用リストへのエントリーの記入は、データの特定の特性に応じてキャッシュデータのエントリーを再使用リストに昇進させる分類方式などのような、他の指名プロセスにより行うことができる。種々の指名プロセスにより昇進されるキャッシュラインは、再使用リストからデータを昇進させるか、または新しく割り当てられたキャッシュデータエントリーを昇進させるか、もしくは別の方法により選択されたキャッシュデータエントリーを昇進させる。指名プロセスを組み合わせるにより、容量及び関連性による変位、あるいはキャッシュを非効率にする他の問題に対処することができる。再使用リストへの昇進のためにキャッシュラインを選択する他の指名プロセスが存在し、また展開可能であると予想されるが、これらは本発明の範囲内であると考えべきである。

【0027】

保護アルゴリズムは、大域リストまたはフィルタリストから排除される要素を選択することにより、再使用リストのサイズをさらに制御する。大域リストはフィルタリストと再使用リストの両方から構成されるため、大域リストからの排除は、再使用リストまたはフィルタリストの何れかからの排除であることに注意されたい。2つの保護アルゴリズム、即ち、ストライクアルゴリズムと、LRUカウンタアルゴリズムについて説明する。

【0028】

ストライクアルゴリズムは、ストライクカウンタを用いて、再使用リストの最近使用頻度が最も低い要素が不変である間にメモリへのアクセスがキャッシュミスとなった回数を追跡する。最大許容ストライクしきい値を予めセットし、ストライクカウンタの値と比較して、ストライクカウンタの値が最大許容ストライクしきい値を超えると、大域リストのLRU要素が排除されるようにする。再使用リストにおける要素の寿命は、再使用リストを大きくするか、または最大許容ストライクしきい値を増加させることにより延長可能であるが、最大許容ストライクしきい値を大きくすると、大域リストの要素の寿命が増加する。従って、比較的大きな再使用リストでは、最大許容ストライクしきい値を低くするのが望ましく、比較的小さい再使用リストでは、このしきい値を高くするのが望ましい。従って、ストライクアルゴリズムの1つの実施例は、再使用リストのサイズに応じてストライクしきい値を変化させる方法を使用する。

【0029】

LRUカウンタ保護アルゴリズムは、フィルタリストからのキャッシュラインの排除に対する再使用リストからのキャッシュラインの排除コストを予想しようとする。カウンタは、フィルタリスト及び再使用リストの両方のLRU要素への参照を追跡するため、各リストのLRU要素がアクセスされる可能性の見積りを比較することができる。かかる方式の1つの実施例は、0でスタートし、再使用リストのLRU要素が参照される毎にカウントアップし、フィルタリストのLRU要素が参照される毎にカウントダウンするカウンタを使用する。別法として、ヒステリシスビットをカウンタに追加して、フィルタリストまたは再使用リストのいずれか一方のリストのLRUビットが連続して(もう一方のリストのLRUへの中間的な参照なしに)参照されると、1つおきの参照のたびにカウンタが増分または減分されるようにする。

【0030】

排除が必要な場合、LRUのカウンタの値を用いて、再使用リストの最大許容サイズを

10

20

30

40

50

求め、再使用リストまたはフィルタリストのいずれからキャッシュラインを排除すべきであるかを決定する。一例として、1つの実施例で排除が必要であり、カウンタが再使用リストへの頻繁なアクセスを示す正の数をホールドするケースを考える。その場合、LRUカウンタの値を、再使用リストのしきい値と比較する。カウンタ値がしきい値を越える場合、所定の大きなサイズのリストが示唆されており、カウンタ値がしきい値よりも小さい場合、所定の小さいサイズのリストが必要とされる。その後、選択された所定のリストサイズを再使用リストの実際のサイズと比較し、所定のサイズが再使用リストの実際のサイズよりも小さい場合に限り、再使用リストからラインを排除する。そうでなければ、フィルタリストからラインを排除する。表は、多数の所定しきい値と、それらに対応する再使用リストのサイズとを含み、その目的は再使用リストのサイズ及びそのリストからの排除の最適化にあることに注意されたい。かかる方式は、再使用リストのサイズを制御するだけでなく、フィルタリストまたは再使用リストの何れにあるかに関係なく、使用される可能性が最も低いデータを排除する作用を有する。

10

【0031】

動作時、再使用リストとフィルタリストとを管理するには、各キャッシュ要素が何れのリストの一部であることを示すデータビットを使用する必要がある。高速キャッシュを最適動作させて、大きなキャッシュヒット対ミス比率から得られる利点を得るためには、各リストの順序付けを、それぞれ独立に、且つそれらを組み合わせさせた大域リストとして、迅速且つ効率良く行わなければならない。

【0032】

再使用リスト指示ビットを使用するLRUPキャッシュの実施例を、かかるキャッシュの再使用リスト及びフィルタリストを管理する1つの方法の一例として、ここで説明する。

20

【0033】

再使用リストまたはフィルタリストのメンバーであるキャッシュ要素間の識別は、かかるキャッシュにおいて、各キャッシュ要素に関連する再使用ビットをチェックすることにより容易に行うことができる。例えば、図4において、キャッシュ要素が再使用リストのメンバーである場合、再使用ビットは、ウェイ1及び2のように、1にセットされている。要素がフィルタリストのメンバーである場合、再使用ビットは、ウェイ0及び3のように、0にセットされている。キャッシュコントローラが、LRUPキャッシュの一部として説明するプロセスの1つを実行するために、フィルタリストまたは再使用リストの何れのLRUウェイかを判定する必要がある場合、キャッシュコントローラが求める順序の少なくともそのリストについて使用順序を維持して、大域リストを2つのリストに分離しなければならない。このプロセスを、サブリスト順序抽出と呼び、かかるプロセスの一例を以下において説明する。

30

【0034】

この例の方法におけるサブリスト順序抽出の目的は、大域リストのLRU要素が抽出を受けるサブリストのLRU要素であるように、順序を変更した大域リストを作成することである。大域リストのMRU位置へ、抽出されないサブリストの要素を再び挿入する方式についてここで説明するが、この方法は、抽出されるリストの順序を保持して、リストのLRUビットが正しく識別されるようにするが、抽出されないリストの要素順序情報は失われる。

40

【0035】

図4を再び参照して、再使用リストのLRU要素の抽出について説明する。抽出されるサブリスト(再使用リスト)は、関連の再使用ビットにより示されるように、要素1及び2より成り、要素0及び3が抽出中でないサブリストのメンバーであることも示す。図4の要素0-4を表すために数字を付した各行は、リストのメンバーが新しいアレイを引き出すために抽出中でない場合、アレイLに適用されるマスクを表す。この新しいアレイは、LRUが抽出中のサブリストのLRUである順序を変更した大域リストを表す。

【0036】

50

例えば、「古い」と記した行は、元のLRUリストを表し、3、2、1、0の順序を示すが、3は最近使用頻度が最も低く、0は最近使用頻度が最も高いキャッシュ要素である。再使用リストを抽出するため、再使用リストの一部でない要素0及び3のマスクがアレイLに適用される。これにより、「新しい」と記した行に示す新しいアレイが生まれるが、これは抽出中のリストの要素の順序を反映する一方、抽出中でないリストの要素がMRU位置にあることをアレイが示すようにする。しかしながら、抽出中のリストの一部でない要素のマスクが、古いアレイLに、アレイLのそれらの古い位置を表さない順序で適用されるため、抽出中でないリストの順序情報が失われる。

【0037】

この特定の例では、新しいアレイ(0、0、0、0、1、1)は2、1、0、3の順序を表すが、2はLRU要素、3はMRU要素である。ウェイ0及び3のマスクがランダムな順序で適用されたため、新しいアレイは、0が実際に3より最近使用頻度が高い時、3が0より最近使用頻度が高いことを示す。しかしながら、この情報はアレイLに保管されず、また抽出中のリストの順序の維持だけに興味があるため、このタイプのエラーは本願にとっては重要ではない。

10

【0038】

あたらしいアレイの作成によるサブリストの抽出は、1つの例において組み合わせ論理を用いることにより行うことが可能であり、従って、非常に迅速に完了することができる。サブリストの順序を異なる方式で取り出す他の実施例が存在するが、それは、抽出中でないサブリストの順序を保持することを含む。サブリストを抽出するについて新しいアレイを作成し使用した後、新しいアレイは古いアレイLに書き戻さない。従って、アレイLはそのままであり、実際の要素の順序が保持され、必要に応じて組み合わせ論理のような方法により、サブリストの非破壊的抽出を繰り返し行うことができる。

20

【0039】

ここで説明するようなキャッシュシステムは、図5に示すシステムにおけるようなコンピュータに使用されるプロセッサの一部である。キャッシュシステムはプロセッサコアの一部である必要はないが、大きなプロセッサモジュールの一部として、別個のデバイスとして、またはプロセッサへ接続されるものとして実現可能である。システムは、システムバス502に取り付けられたプロセッサコア501を有する。システムバスとプロセッサは、キャッシュコントローラ503によりリンクされ、このコントローラはキャッシュ504にも接続されている。バスは、ハードディスクドライバ505のようなデータ記憶装置に接続され、またポート506のような1またはそれ以上の入出力ポートに接続されている。キャッシュコントローラは、プロセッサからメモリフェッチリクエストを受信し、リクエストされたデータがそこに記憶されておれば、キャッシュからリクエストに対応し、そうでなければ、主メモリからそのリクエストに対応する。キャッシュコントローラは、前述したように、本発明に従って、キャッシュ要素置換のLRUP方式により、キャッシュ504におけるデータを維持する。

30

【0040】

キャッシュコントローラ503のようなキャッシュコントローラは、集積回路上の論理回路として物理的に実現して、集積回路の少なくとも一部が、前述したように、LRUP方式によりキャッシュにおいてデータを維持するように作動可能である。例えば、上述したように、キャッシュを再使用リストと、フィルタリストに分割し、再使用リストと、フィルタリストにエンタリーを記入し、再使用リスト及びフィルタリストからデータエンタリーを排除する機能を実行可能なモジュールを備えた集積回路は、本発明の実施例と一貫性を有するものである。特に、本発明は、プロセッサまたはプロセッサモジュール、マザーボードチップセット、キャッシュコントローラチップ、または他の集積回路の少なくとも一部より成る回路で実施されることを企図されており、これらは全て本発明の範囲内にあると考えるべきである。

40

【0041】

図6は、データ記憶サブシステムを示す。このシステムは、コンピュータシステムのシス

50

テムバス602に記憶コントローラ601が接続されたものである。記憶コントローラは、多量のデータの非揮発性記憶装置として一般的に使用される、1またはそれ以上のハードディスクドライブ603にも接続されている。記憶コントローラはさらに、キャッシュ604に接続され、このキャッシュは、ハードディスクドライブ603よりも迅速にデータの記憶及び取り出しを行なうことができるが、一般的に高価で小容量である。記憶コントローラ601は、ハードディスクドライブからのデータを記憶し、取り出すリクエストに回答し、キャッシュ604に使用頻度の高いデータを維持する。キャッシュコントローラは、本願のLRUPキャッシュの説明に従って、キャッシュのデータを管理する。

【0042】

上述したような保護機能付きLRUは、従来のLRU置換方式の共通の問題に対する解決法を与えることにより、従来のLRU置換方式に対して実質的な性能上の利点を提供する。キャッシュエントリリストをフィルタサブリストと、再使用リストに分割することにより、エントリ記入及び保護プロセスを実行して、関連性及び容量による変位を減少させることができる。新しいキャッシュエントリは最初にフィルタリストに記憶され、フィルタリストから昇進したエントリが再使用リストに記入される。フィルタリスト及び再使用リストからの排除は、フィルタリスト、再使用リストまたは大域リストからデータエントリを排除する保護プロセスにより行われる。保護及び排除プロセスの多くの変形例を、それぞれが従来のLRU置換方式に存在する望ましくない変位の問題の影響を軽減することにより得られる利点と共に説明した。

10

【0043】

特定の実施例を図示説明したが、当業者にとっては、同じ目的を達成する任意の構成を図示の特定の実施例に置き換えることが可能であることがわかるであろう。本願は、本発明の任意の変形例または設計変更を包含するように意図されている。本発明は、特許請求の範囲及びその均等物の全範囲によってのみ限定されることを意図されている。

20

【図面の簡単な説明】

【図1】 図1は、本発明の実施例と一貫する、最近最小使用頻度(LRU)順序マトリックスを表す表である。

【図2】 図2は、本発明の実施例と一貫する、特定のウェイのデータへのアクセスに回答して変化するマルチウェイのLRU置換キャッシュワードのビットを表す表である。

【図3】 図3は、本発明の実施例と一貫する、マルチウェイLRUキャッシュ置換における犠牲者選択に使用する復号マスクを表す表である。

30

【図4】 図4は、本発明の実施例と一貫する、サブリスト順序付けマスクと、関連のLRUウェイ保護ビットを備えたLRUワードへのかかるマスクの適用例とを表す表である。

【図5】 図5は、本発明の実施例と一貫する、プロセッサ、メモリ及びLRUPキャッシュを備えた汎用コンピュータを示す。

【図6】 図6は、本発明の実施例と一貫する、LRUPキャッシュを組み込んだデータ記憶システムを示す。

【図1】

LRU ウェイ表

i \ j	0	1	2	3
0	-	0	1	3
1		-	2	4
2			-	5
3				-

L={0,1),(0,2),(1,2),(0,3),(1,3),(2,3)}
OR L={0,1,2,3,4,5}

【図3】

犠牲者選択復号表

ビットマスク 検出されたウェイ	0	1	2	3	4	5
0	0	0	-	0	-	-
1	1	-	0	-	0	-
2	-	1	1	-	-	0
3	-	-	-	1	1	1

【図2】

更新表

変化したビット 使用されるウェイ	0	1	2	3	4	5
0	1	1	-	1	-	-
1	0	-	1	-	1	-
2	-	0	0	-	-	1
3	-	-	-	0	0	0

1= 最近使用頻度が高い

【図4】

LRUP 大域リストウェイ表

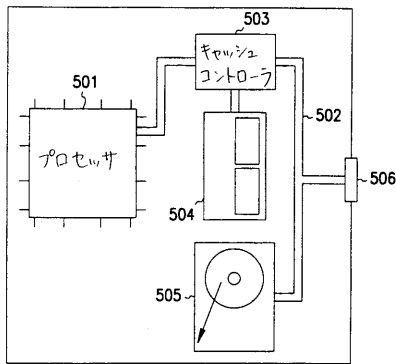
ビットウェイ	0	1	2	3	4	5	再使用
0	0	0	0	-	-	-	0
1	1	-	-	0	0	-	1
2	-	1	-	1	-	0	1
3	-	-	1	-	1	1	0

例:(マスク0.3を適用)

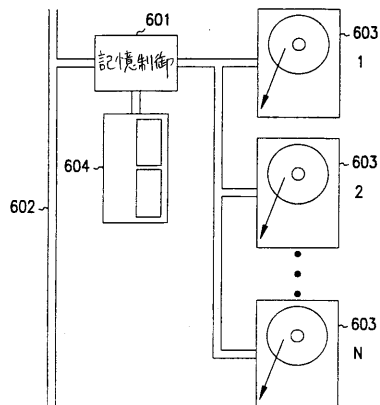
古い	0	0	0	0	0	0	0
新しい	0	0	0	0	0	1	1

XX=0
 X1=1
 01=0
 0X=0
 00=0
 11=1

【図5】



【図6】



フロントページの続き

(72)発明者 ウェード, ニコラス, デイ
アメリカ合衆国 オレゴン州 97229 ポートランド エイティファースト・ストリート ノ
ースウエスト 2730

審査官 清木 泰

(56)参考文献 特開平04-092941(JP, A)
特開平11-338775(JP, A)
特開平10-307756(JP, A)
特開平06-161898(JP, A)
特開平05-189316(JP, A)
特開平03-071347(JP, A)
特開平02-281350(JP, A)
米国特許第05778442(US, A)
Ramakrishna Karedla, J. Spencer Love, Bradley G. Wherry, HDDの性能を改善するための
キャッシュ方式を提案, 日経エレクトロニクス, 日本, 日経BP社, 1994年 9月12日,
No:617, Pages:107-121
Method for Improving Least Recently Used Buffer Performance Using Multiple Stack Inser
tion Points Based on Data Types, IBM Technical Disclosure Bulletin, 米国, Internationa
l Business Machines Corporation, 1993年 8月, Vol:36, No:08, Pages:479-480

(58)調査した分野(Int.Cl., DB名)
G06F12/08-12/12