

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5583274号
(P5583274)

(45) 発行日 平成26年9月3日(2014.9.3)

(24) 登録日 平成26年7月25日(2014.7.25)

(51) Int. Cl. F 1
G 0 6 F 12/10 (2006.01) G O 6 F 12/10 5 O 1 Z
G 0 6 F 12/08 (2006.01) G O 6 F 12/08 5 O 7 Z

請求項の数 10 (全 11 頁)

(21) 出願番号	特願2013-523548 (P2013-523548)	(73) 特許権者	391030332
(86) (22) 出願日	平成23年7月7日(2011.7.7)		アルカテルルーセント
(65) 公表番号	特表2013-544380 (P2013-544380A)		フランス国、9 2 1 0 0 ・ブローニュービ
(43) 公表日	平成25年12月12日(2013.12.12)		ヤンクール、ルート・ドゥ・ラ・レーヌ・
(86) 国際出願番号	PCT/EP2011/061513		1 4 8 / 1 5 2
(87) 国際公開番号	W02012/016783	(74) 代理人	100094112
(87) 国際公開日	平成24年2月9日(2012.2.9)		弁理士 岡部 譲
審査請求日	平成25年3月21日(2013.3.21)	(74) 代理人	100106183
(31) 優先権主張番号	10305868.1		弁理士 吉澤 弘司
(32) 優先日	平成22年8月6日(2010.8.6)	(74) 代理人	100128657
(33) 優先権主張国	欧州特許庁 (EP)		弁理士 三山 勝巳
		(74) 代理人	100160967
			弁理士 ▲濱▼口 岳久
		(74) 代理人	100170601
			弁理士 川崎 孝

最終頁に続く

(54) 【発明の名称】 コンピュータ・メモリを管理する方法、対応するコンピュータ・プログラム製品、およびそのためのデータ・ストレージ・デバイス

(57) 【特許請求の範囲】

【請求項 1】

コンピュータ・メモリを管理する方法であって、
 仮想アドレスを物理アドレスにマップするためのページ・テーブル・エントリー、および複数のデータ・ブロックを含むキャッシュを保持するステップ(101)と、

前記仮想アドレスへの参照にตอบสนองして、前記ページ・テーブル・エントリーを用いて前記仮想アドレスを前記物理アドレスへと変換するステップ(102)と、データを前記物理アドレスから前記キャッシュ内へフェッチするステップ(103)と

を含む方法において、前記ページ・テーブル・エントリーが、複数のインジケータを含み、それぞれのデータ・ブロックが、前記複数のインジケータのうちの1つに対応し、前記データを前記キャッシュ内へフェッチするステップ(103)が開始すると、当該方法が

前記複数のインジケータから選択された1つのインジケータ(110)がセットされていることにตอบสนองして、前記対応するデータ・ブロックをゼロ化するさらなるステップ(104)を含むことを特徴とする方法。

【請求項 2】

前記ページ・テーブル・エントリーがビットマスクを含み、前記インジケータが、前記ビットマスク内に含まれているビットであることを特徴とする、請求項1に記載の方法。

【請求項 3】

前記ページ・テーブル・エントリーが、前記複数のデータ・ブロックを含むメモリ・ペ

ージに関連付けられ、前記方法が、

データ・ブロックを受け取る中間ステップと、

前記データ・ブロックを前記メモリ・ページ内に格納する中間ステップと、

前記データ・ブロックが対応する前記インジケータをクリアする中間ステップとを含むことを特徴とする、請求項 1 に記載の方法。

【請求項 4】

前記データ・ブロックが対応していない、前記複数のインジケータのうちの残りのインジケータをセットするさらなるステップ

を含むことを特徴とする、請求項 1 に記載の方法。

【請求項 5】

さらなるデータ・ブロックを上書きする後続ステップと、

前記さらなるデータ・ブロックが対応する前記インジケータをクリアする後続ステップと

を含むことを特徴とする、請求項 3 または 4 に記載の方法。

【請求項 6】

コンピュータ・プログラムであって、プログラムがコンピュータ上で実行されたときに請求項 1 に記載の方法を実行するためのコンピュータ実行可能命令を含む、コンピュータ・プログラム。

【請求項 7】

オペレーティング・システムを含むことを特徴とする、請求項 6 に記載のコンピュータ・プログラム。

【請求項 8】

仮想アドレスを物理アドレスにマップするためのページ・テーブル・エントリー、および複数のデータ・ブロックを含むキャッシュを保持するステップ (1 0 1) を行うための手段と、

前記仮想アドレスへの参照に回答して、前記ページ・テーブル・エントリーを用いて前記仮想アドレスを前記物理アドレスへと変換するステップ (1 0 2)、およびデータを前記物理アドレスから前記キャッシュ内へフェッチするステップ (1 0 3) を行うための手段と

を含むデバイスにおいて、前記ページ・テーブル・エントリーが、複数のインジケータを含み、それぞれのデータ・ブロックが、前記複数のインジケータのうちの 1 つに対応し、当該デバイスが、前記データを前記キャッシュ内へフェッチするステップ (1 0 3) が開始されると、

前記複数のインジケータから選択された 1 つのインジケータ (1 1 0) がセットされていることに回答して、前記対応するデータ・ブロックをゼロ化するステップ (1 0 4) を行うためのさらなる手段を含むことを特徴とするデバイス。

【請求項 9】

中央処理装置、

メモリ・マネジメント・ユニット、および

データ・キャッシュ

のうちの少なくとも 1 つを含むことを特徴とする、請求項 8 に記載のデバイス。

【請求項 10】

前記ページ・テーブル・エントリーを格納するように構成されている変換ルックアサイド・バッファをさらに含み、前記仮想アドレスが、前記変換ルックアサイド・バッファを用いて変換される (1 0 2) ことを特徴とする、請求項 8 または 9 に記載のデバイス。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、仮想アドレスを物理アドレスにマップするためのページ・テーブル・エントリー、および複数のデータ・ブロックを含むキャッシュを保持するステップと、仮想アド

10

20

30

40

50

レスへの参照に応答して、ページ・テーブル・エントリを用いて仮想アドレスを物理アドレスへと変換するステップと、データを物理アドレスからキャッシュ内へフェッチするステップとを含む、コンピュータ・メモリを管理する方法に関する。本発明はさらに、コンピュータ・プログラム製品であって、プログラムがコンピュータ上で実行されるときに前記方法を実行するためのコンピュータ実行可能命令を含む、コンピュータ・プログラム製品に関し、また、前記方法を実行するようにプログラムまたは構成されているデバイスに関する。

【背景技術】

【0002】

コンピューティングにおいては、メモリ・マネージメントは、コンピュータ・メモリを管理する行為である。これは、そのよりシンプルな形態においては、プログラムの要求に応じてそれらのプログラムにメモリの一部分を割り当てるための方法を提供することと、もはや必要なくなったときに再利用のためにメモリを解放することとを含む。メモリのマネージメントは、いかなるコンピュータ・システムにとっても重要である。

10

【0003】

仮想メモリ・システムは、プロセスによって使用されるメモリ・アドレスを物理アドレスから分離し、プロセスの分離を可能にし、メモリのうちの有効に利用できる量を増やす。仮想メモリ・マネージャーの質は、全体的なシステム・パフォーマンスに著しい影響を与える。

【0004】

20

Ludmila CherkasovaおよびRob Gardnerは、「Measuring CPU Overhead for I/O Processing in the Xen Virtual Machine Monitor」、Proceedings of the USENIX Annual Technical Conference 2005、アメリカ合衆国カリフォルニア州アナハイムにおいて、仮想メモリ・システムの概要を示している。オペレーティング・システム(OS)が、メイン・メモリにおいて使用するために二次ストレージからデータを取り出せるようにするために、この従来のシステムは、ページングとして知られているメモリ・マネージメント・スキームを利用し、このページングにおいては、データは、メモリ・ページと呼ばれる固定サイズのブロックで取り出される。本明細書においては、二次ストレージ(外部メモリと呼ばれる場合もある)とは、ハード・ディスク・ドライブ、光ストレージ、フラッシュ・メモリ、フロッピー・ディスク、磁気テープ、紙テープ、せん孔カード、またはZipドライブなど、中央処理装置(CPU)によって直接アクセス可能ではない任意のストレージを意味する。

30

【0005】

一般に入力/出力またはI/Oと呼ばれる、他のシステムとの効率のよい通信を可能にするために、現況技術の仮想メモリ・システムは、ページ・フリッピングとして知られているメカニズムを実装している。この技術によれば、アプリケーションが、入力データを受け取るための1つまたは複数のメモリ・ページを指定し、それによってOSは、その内蔵されたページング・スキームを用いて、そのデータを供給することができる。このために、OSは、アプリケーションのワーキング・メモリ(一般には、そのアドレス空間と呼ばれる)内の指定されたメモリ・ページを、要求された入力と「交換」する。

40

【0006】

この知られているシステムの主なマイナス面は、供給される入力データがターゲットのメモリ・ページに正確にフィットすることがめったにないという事実にある。そのようにして交換されるそれぞれのメモリ・ページの残りは、メモリの初期化に依存するアプリケーションに対応するために、ならびに要求元のアプリケーションが外部の潜在的に機密性のデータへの無許可のアクセスを得ることを防止するために、ゼロ化する必要がある。ソフトウェア開発において、ゼロ化するとは、固定された無意味な値、たとえばゼロでデータを上書きして、そのデータの不測の開示を防止することを意味し、そのような開示は、

50

要求元のアプリケーションによるセキュリティ侵害を潜在的に可能にする。

【0007】

米国特許第5920895A号によれば、マップされたファイルI/Oを使用してキャッシュされるファイルを書き込む効率は、ファイルのキャッシュされたページ内の初期化されていないデータをゼロ化することを、そのファイルがユーザ・モード・スレッドによってマップされるまで控えることによって、改善される。ページング・オペレーションが仮想メモリ・マネージャーによってコントロールされ、マップされたファイルI/Oを使用したメモリ・ベースのキャッシングがキャッシュ・マネージャーによって管理されるオペレーティング・システムにおいては、書き込みの際に、マップされたファイルをゼロ化することを控えることは、仮想メモリ・マネージャーとキャッシュ・マネージャーの間における通信のための一式の内部オペレーティング・システム・インターフェースによって実施される。キャッシュされているファイルが、ユーザ・モード・スレッドによってまだマップされていないときには、キャッシュ・マネージャーは、そのファイルのキャッシュ・ページがどの程度まで書き込まれているかを追跡把握し、それによって、そのキャッシュ・ページ内の初期化されていないデータはすべて、後でそのファイルがユーザ・モード・スレッドによってマップされたときにゼロ化することができる。

10

【0008】

米国特許出願第2008/229117A1号による、コンピューティング環境におけるデジタル上の著作権侵害を防止するための方法は、アプリケーションをコンピューティング環境内にロードするステップであって、前記アプリケーションが、暗号化キーを使用して暗号化される、ステップと、仮想アドレス空間を前記アプリケーションに割り当てるステップと、前記アプリケーションのための前記暗号化キーを、中央処理装置によってのみアクセス可能であるレジスタ内にロードするステップと、前記キーのためのインデックス値を、前記アプリケーションのための前記仮想アドレス空間に対応するページ・テーブル・エントリー内の前記レジスタ内に格納し、それによって、前記仮想アドレス空間を前記アプリケーションのための前記キーにリンクさせるステップとを含む。

20

【0009】

MENON A.らによる「Optimizing Network Virtualization in Xen」、PROCEEDINGS OF THE 2006 USENIX ANNUAL TECHNICAL CONFERENCE、2006年5月29日(2006-05-29)、XP002623699において、筆者らは、Xen仮想化環境においてネットワーク・パフォーマンスを最適化するための3つの技術を提案および評価している。彼らの技術は、デバイス・ドライバを、I/Oデバイスへのアクセスを有する特権のある「ドライバ」ドメイン内に配置し、仮想化されたネットワーク・インターフェースを通じて、特権のない「ゲスト」ドメインにネットワーク・アクセスを提供するという基本的なXenアーキテクチャーを保持している。

30

【先行技術文献】

【特許文献】

【0010】

【特許文献1】米国特許第5920895A号

40

【特許文献2】米国特許出願第2008/229117A1号

【非特許文献】

【0011】

【非特許文献1】Ludmila CherkasovaおよびRob Gardner、「Measuring CPU Overhead for I/O Processing in the Xen Virtual Machine Monitor」、Proceedings of the USENIX Annual Technical Conference 2005、アメリカ合衆国カリフォルニア州アナハイム

【非特許文献2】MENON A.らによる「Optimizing Network Virtualization in Xen」、PROCEEDINGS OF TH

50

E 2006 USENIX ANNUAL TECHNICAL CONFERENCE、2006年5月29日(2006-05-29)、XP002623699

【非特許文献3】IEEE 802.3-2008標準

【発明の概要】

【発明が解決しようとする課題】

【0012】

本発明の目的は、ページ・フリッピングの際にゼロ化を行う必要性をなくすかまたは少なくする、仮想メモリ・マネージメントに対する改善されたアプローチを提示することである。さらなる目的は、限られている帯域幅のメモリに適した方法を提供することであり、すなわち、その速度でデータがCPUによってメモリから読み出され、またはメモリに格納され得る。他のさらなる目的は、メモリ・バス、すなわち、メイン・メモリを、そのメイン・メモリに入るデータおよびそのメイン・メモリから出るデータの流れを管理するメモリ・コントローラに接続するコンピュータ・サブシステムに課される負荷を少なくすることにある。本明細書においては、メイン・メモリ(プライマリー・ストレージまたは内部メモリとも呼ばれる)とは、ランダムアクセス・メモリ(RAM)、読み取り専用メモリ(ROM)、プロセッサ・レジスタ、またはプロセッサ・キャッシュなど、CPUによって直接アクセス可能な任意のストレージを意味する。

10

【課題を解決するための手段】

【0013】

この目的は、コンピュータ・メモリを管理する方法によって達成され、本方法は、仮想アドレスを物理アドレスにマップするためのページ・テーブル・エントリー、および複数のデータ・ブロックを含むキャッシュを保持するステップと、仮想アドレスへの参照に回答して、ページ・テーブル・エントリーを用いて仮想アドレスを物理アドレスへと変換するステップと、データを物理アドレスからキャッシュ内へフェッチするステップとを含み、ページ・テーブル・エントリーは、複数のインジケータを含み、それぞれのデータ・ブロックは、複数のインジケータのうちの一つに対応し、データをキャッシュ内へフェッチするステップが開始すると、本方法は、前記複数のインジケータから選択された一つのインジケータがセットされていることに回答して、対応するデータ・ブロックをゼロ化するさらなるステップを含む。この目的はさらに、コンピュータ・プログラム製品であって、プログラムがコンピュータ上で実行されたときに前記方法を実行するためのコンピュータ実行可能命令を含む、コンピュータ・プログラム製品によって、または前記方法を実行するようにプログラムもしくは構成されているデバイスによって達成される。

20

30

【0014】

本発明の主要なアイデアは、仮想アドレス、すなわち、アプリケーションのワーキング・メモリ内のあるロケーションを表すインデックスと、対応する物理アドレスとの間におけるマッピングを格納するためにOSの仮想メモリ・システムによって使用されるデータ構造であるページ・テーブルを増強することである。仮想メモリのコンテキストにおいては、仮想アドレスは、アクセス・プロセスにとってのみ一意であり、その一方で物理アドレスは、メイン・メモリの実際のストレージ・セルを指す。

【0015】

本発明のさらなる発展形態は、従属請求項および以降の説明から得ることができる。

40

【0016】

以降では、添付の図面を参照して、本発明についてさらに説明する。

【0017】

本発明の一実施形態に従ってコンピュータ・メモリを管理するために、仮想アドレスを物理アドレスにマップするためのページ・テーブル・エントリーが保持される。さらに、複数のデータ・ブロックを含むキャッシュが保持される。ページ・テーブル・エントリーは、複数のインジケータを含み、それぞれのデータ・ブロックは、一つのインジケータに対応する。仮想アドレスへの参照に回答して、その仮想アドレスは、ページ・テーブル・エントリーを用いて物理アドレスへと変換され、データが、その物理アドレスからキャッ

50

シユ内へフェッチされる。データをフェッチすると、インジケータがセットされていることに応答して、対応するデータ・ブロックがゼロ化される。

【図面の簡単な説明】

【0018】

【図1】本発明の一実施形態による方法を示すフローチャートである。

【発明を実施するための形態】

【0019】

以降では、図1を参照して、本発明による方法を例として説明する。

【0020】

図1のフローチャート100は、第1の処理ステップ101、第2の処理ステップ102、第3の処理ステップ103、および第4の処理ステップ104を含む。一連の矢印は、処理ステップ101から103を通過して、内部ストレージ110と合流し、最終的に第4の処理ステップ104へとわたるコントロールのフローを表している。このフロー内では、第1のシンボルから始まって第2のシンボルで終わる矢印は、コントロールが第1のシンボルから第2のシンボルへわたることを示している。

【0021】

目下の実施形態においては、本方法は、専用のメモリ・マネージメント・ユニット(MMU)によって適用される。マイクロプロセッサの設計においては、MMUとは、CPUによって要求されるメモリへのアクセスを処理することを担当する任意のコンピュータ・ハードウェア・コンポーネントを意味する。仮想メモリ・マネージメント、すなわち、仮想アドレスを物理アドレスへ変換することのほかに、MMUは、メモリ保護、バス・アービトレーション、および、8ビット・システムなどのよりシンプルなコンピュータ・アーキテクチャーにおけるバンク切り替えなどの技術をサポートすることができる。

【0022】

メモリは物理的に断片化されることがあり、二次ストレージ上へオーバーフローすることさえあるという事実の一方で、プロセスが、切れ目のないワーキング・メモリ(一般にはアドレス空間と呼ばれる)の概念に基づくことを可能にするために、MMUは、ページングをサポートし、ひいては、当技術分野においてページドMMU(PMMU)として知られている形式を取る。

【0023】

CPUがメモリにアクセスするのに必要とする平均時間を少なくするために、PMMUはさらに、CPUキャッシュ、すなわち、最も頻繁に使用されるメイン・メモリ・ロケーションからのデータのコピーを格納するように構成されている、より小さな、より高速なメモリをコントロールする。CPUは、メイン・メモリ内のあるロケーションからの読み取り、またはそのロケーションへの書き込みを行う必要がある場合には、まず、そのデータのコピーがキャッシュ内にあるかどうかをチェックする。そのデータのコピーがキャッシュ内にある場合には、CPUはすぐに、キャッシュからの読み取り、またはキャッシュへの書き込みを行い、これは、メイン・メモリからの読み取り、またはメイン・メモリへの書き込みよりも、本質的に高速である。キャッシングのコンテキストにおいては、メイン・メモリは、キャッシュのバッキング・ストアと呼ばれる場合がある。

【0024】

より具体的には、PMMUによってコントロールされるCPUキャッシュは、データのフェッチおよび格納をスピード・アップするためのデータ・キャッシュ、ならびに、実行可能な命令およびデータの両方に関する仮想メモリ・マネージメントをスピード・アップするための変換ルックアサイド・バッファ(TLB)を含む。TLBは、仮想アドレスを物理アドレスにマップするページ・テーブル・エントリ(PTE)を含む固定された数のスロットを含む。典型的な一実施態様においては、それぞれのPTEは、32から64の間のビットを含む。

【0025】

第1の処理ステップ101において、PMMUは、仮想アドレスを物理アドレスにマッ

10

20

30

40

50

プするための P T E を保持する。さらに P M M U は、複数のデータ・ブロックを含むキャッシュを保持し、それらのデータ・ブロックは一般に、C P U キャッシングのコンテキストにおいては、キャッシュ・ラインまたはキャッシュ・ブロックと呼ばれる。典型的に、それぞれのデータ・ブロックは、8 から 5 1 2 バイトのサイズにわたる。ハードウェア・アクセラレーションを最大にするために、目下の P M M U は、現況技術の 3 2 ビットおよび 6 4 ビットの C P U アーキテクチャーにおいては 1 バイトから 1 6 バイトにわたることが典型的である単一の C P U 命令によって要求される可能性があるデータの量よりも大きいデータ・ブロックを採用する。

【 0 0 2 6 】

第 1 の処理ステップ 1 0 1 において保持されている P T E は、複数のインジケータを含み、それぞれのデータ・ブロックは、1 つのインジケータに対応し、このインジケータは、その対応するデータ・ブロックを要求元のアプリケーションに提供する前にゼロ化する必要があるかどうかを示す。したがって、このインジケータは、フラグ、すなわち「真」または「偽」の値のいずれかを有するブール変数の形態を取る。P T E によって必要とされるストレージ容量を最小限に抑えるために、それぞれのインジケータは、ビットの形態を取り、そのバイナリー値のそれぞれは、割り当てられた意味を有する。C P U が複数のインジケータを単一のビット単位のオペレーション内にセットすることをさらに可能にするために、それらのビットは、コンピュータ計算においてはビットマスクとして知られている適切にサイズ設定されたベクトルへとグループ化される。代替実施形態においては、平均の書き込みスピードを犠牲にして P T E のサイズをさらに減らすために、それぞれのビットを、複数のデータ・ブロックに対応するように解釈して、ビットマスクを構成するビットの数を有効に減らすことができる。

【 0 0 2 7 】

中間のステップ (図示せず) において、O S は、アプリケーションによる要求に応じて、ネットワーク・インターフェース・カード、ネットワーク・アダプタ、またはローカル・エリア・ネットワーク (L A N) アダプタとしても知られている関連付けられたネットワーク・インターフェース・コントローラ (N I C) を用いて、コンピュータ・ネットワークを介してデータ・ブロックを受け取る。スループットを最大にするために、データ・ブロックは、ギガビット・イーサネット (G b E , 1 G i g E) 接続を用いて、すなわち、I E E E 8 0 2 . 3 - 2 0 0 8 標準によって定義されているように 1 ギガビット / 秒の速度で伝送される。C P U を用いて I / O データをコピーするというオーバーヘッドを避けることによってオペレーションをさらにスピード・アップするために、O S は、ページ・フリッピングを採用して、アプリケーションによって指定されたメモリ・ページ内での I / O データを格納する。メモリ・ページを構成するデータ・ブロックのうちの一部のみが、I / O オペレーションによって修正され、影響を受けない残りのデータ・ブロックは、すぐにはゼロ化されないと想定すると、O S は、ゼロ化されるデータ・ブロックに対応する、P T E のビットマスク内のビットをセットし、その一方で、修正されたデータ・ブロックに対応するビットをクリアする。

【 0 0 2 8 】

イーサネットはパケット交換プロトコルであるため、I / O データは、パケット、またはこのコンテキストにおいてはイーサネット・フレームと呼ばれるデータ伝送単位の形態を取る。このフレームは、単なる確認応答に必要とされる場合などの 6 4 バイトの下限と、コンピュータ・ネットワーキングの技術分野においては一般に最大伝送単位 (M T U) と呼ばれる 9 0 0 0 バイト以上の上限との間でさまざまなサイズとすることができる。4 から 6 4 キロバイトの典型的なメモリ・ページ・サイズを考慮に入れると、従来のメモリ・マネジメント・システムならば、このシナリオにおいて、それぞれのパケットを受け取った際に消去する必要があるメモリ・ページの多大な部分に起因して多大なゼロ化のオーバーヘッドを C P U およびメモリに課すことになる。したがって、本発明のさらなる利点は、小規模から中規模のパケット・データを受け取ることに對するその特有の適性にある。

10

20

30

40

50

【 0 0 2 9 】

ほとんどの場合において、NICを通じて受け取られるデータの量は、PTEによって定義されたデータ・ブロックの固定サイズの倍数と一致しない可能性がある。そのような場合においては、I/Oによって影響を受ける最初のおよび/または最後のデータ・ブロックの少なくとも一部が、定義されないままとなる。その中に含まれているデータを要求元のアプリケーションに開示することを避けるために、OSは、ページ・フリッピングの前に、データ・ブロックのうちのその残りの部分をすぐにゼロ化する必要がある。

【 0 0 3 0 】

代替実施形態においては、NICを通じてデータを受け取る代わりに、プログラムをCPUによって実行するためにメイン・メモリ内へロードするようOSに要求することができる。この場合には、実行ファイル全体の物理的な転送を避けるために、プログラム・ローダは、ページ・テーブルおよびTLBを更新し、実行ファイル用に指定された仮想アドレスを、関連付けられたオブジェクト・ファイルのコンテンツにマップすることを有効に宣言する。このようにしてマップされた仮想メモリのセグメントは、メモリマップされたファイルと一般に呼ばれ、使用されていないプログラム・コードをメイン・メモリ内へロードする必要性をいっさいなくすることができるという利点を有する。

10

【 0 0 3 1 】

場合によっては、ロードされるプログラムは、BSS (block started by symbol) として知られているセクションを含み、これは、BSSセグメントと呼ばれる場合が多く、はじめに、すなわち、プログラムの実行が開始するときに、ゼロの値のデータで満たされることを予期されている静的に割り当てられた変数を含む。この場合には、実行の前にBSSセクション全体をゼロ化することを避けるために、OSは、ゼロ化されるデータ・ブロックに対応する、PTEのビットマスク内のビットをセットし、その一方で、プログラムのその他のデータ・セグメントに対応するビットをクリアすることによって、同様に上述のメカニズムを採用する。

20

【 0 0 3 2 】

第2の処理ステップ102において、仮想アドレスがCPU命令によって参照されたことに応答して、PMMUは、仮想アドレスを物理アドレスへと変換することによって、仮想メモリ・マネージメントというその主要な機能を果たす。このために、PMMUは、第1の処理ステップ101のTLB内に含まれているPTEを採用する。より詳細には、PMMUは最初に、目下の仮想アドレスに対応するPTEを探してTLBを検索する。マッチするものが見つかり、すなわち、仮想メモリ・マネージメントの技術分野においてTLBヒットとして知られている状況になると、そのマッチするPTEから直接、物理アドレスを取り出すことができる。マッチするものが見つからないと、すなわち、一般にTLBミスと呼ばれる状況になると、PMMUには、OSによって保持されているページ・テーブルを調べることが要求される。このページ・テーブルが、マッチするPTEを含む場合には、PMMUは、そのPTEをTLB内にキャッシュし、失敗しているCPU命令を再実行する。

30

【 0 0 3 3 】

第3の処理ステップ103において、PMMUは、第2の処理ステップ102において得られた物理アドレスにアクセスし、その中に含まれているデータ・ブロックをCPUのデータ・キャッシュ内へフェッチし始める。データ・キャッシュがこの時点で既に満杯である場合には、占有されているキャッシュ・ラインを上書きする必要が生じることがあり、新たにフェッチされたデータ・ブロックを収容するために、その前のコンテンツを破棄する必要が生じることがある。どのキャッシュ・ラインの占有を解消するかを選択するためのさまざまなアルゴリズムが、当技術分野において知られている。

40

【 0 0 3 4 】

データ・キャッシュ上のいかなる書き込みオペレーションもスピード・アップするために、データ・キャッシュは、ライトバック・キャッシュ (ライトビハインド・キャッシュと呼ばれる場合もある) の形態を取る。このポリシーによれば、書き込みは、すぐにはバ

50

ッキング・ストアに反映されない。代わりに、PMMUは、どのキャッシュ・ラインがCPUによって上書きされているかを追跡把握し、それに応じて、それらのロケーションを「ダーティー」としてマークする。ダーティーなキャッシュ・ラインの占有を解消する前に、それぞれのデータ・ブロックがバックグ・ストアに書き戻されて、その更新されたコンテンツがその後の読み取りアクセスのために保存される。このアプローチは、コンピュータ・サイエンスにおいては「レイジー・ライト」として知られている。データ・ブロックを書き戻したら、PMMUは、そのデータ・ブロックがいつか取り出された場合にゼロ化されることを防止するためにPTEのビットマスク内の対応するビットをクリアする必要がある。

【0035】

10

それぞれのデータ・ブロックをフェッチするために、空いているキャッシュ・ラインが識別されると、または占有されているキャッシュ・ラインが解放されると、PMMUは、自分の内部ストレージ110内で見つかったPTE内に含まれているビットマスクを調べる。その対応するビットがセットされている場合には、PMMUは、メイン・メモリからデータ・ブロックを取り出さずに、第4の処理ステップ104においてそれぞれのキャッシュ・ラインをゼロ化してから、そのデータ・ブロックを要求元のアプリケーションが利用できるようにする。

【0036】

上述のさまざまな方法のステップは、プログラムされたコンピュータによって実行することができるということを当業者なら容易に認識するであろう。本明細書においては、いくつかの実施形態は、プログラム・ストレージ・デバイス、たとえばデジタル・データ・ストレージ・メディアをカバーするように意図されており、それらのプログラム・ストレージ・デバイスは、マシンまたはコンピュータによって読み取ることができ、マシンによって実行可能なまたはコンピュータによって実行可能な命令のプログラムをエンコードし、前記命令は、本明細書に記載の方法のステップのうちのいくつかまたはすべてを実行する。それらのプログラム・ストレージ・デバイスは、たとえば、デジタル・メモリ、磁気ディスクもしくは磁気テープなどの磁気ストレージ・メディア、ハード・ドライブ、または光学的に読み取り可能なデジタル・データ・ストレージ・メディアとすることができる。これらの実施形態はまた、本明細書に記載の方法の前記ステップを実行するようにプログラムされたコンピュータをカバーするように意図されている。

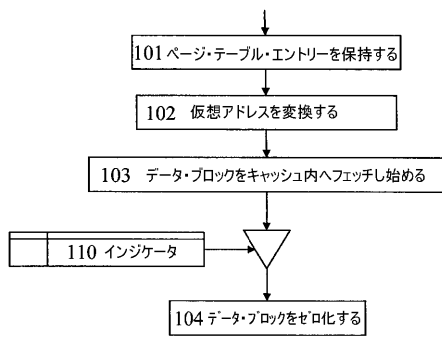
20

30

【0037】

本発明は、その他の特定の装置および/または方法で具体化することもできる。記載した実施形態は、すべての点で例示的なものにすぎず、限定的なものではないとみなすべきである。とりわけ、本発明の範囲は、本明細書に記載の説明および図ではなく、添付の特許請求の範囲によって示される。特許請求の範囲の均等物の意味および範囲内に収まるすべての変更は、特許請求の範囲内に包含される。

【図 1】
100



フロントページの続き

- (72)発明者 マレンダー, セイブ
ベルギー ベー - 2 6 0 0 アントワープ, アンジャベルシュトラート 1 5
- (72)発明者 マッキー, ジェームス バルマー
アメリカ合衆国 0 7 9 7 4 ニュージャーシィ, マレー ヒル, チェスナットヒル ドライヴ
1 5
- (72)発明者 ピアネーゼ, ファビオ
ベルギー ベー - 1 0 0 0 ブリュッセル, ル ドゥ ラ マドレーヌ 6 1
- (72)発明者 エヴァンス, ノア
ベルギー ベー - 2 0 0 0 アントワープ, スワールドシュトラート 1 4

審査官 野田 佳邦

- (56)参考文献 特開平06 - 2 0 8 5 1 1 (J P , A)
米国特許第0 5 9 2 0 8 9 5 (U S , A)
特開平0 2 - 1 7 6 9 5 1 (J P , A)
特開2 0 0 8 - 2 7 6 7 7 8 (J P , A)
Xiaolan Zhang, Suzanne McIntosh, Pankaj Rohatgi, and John Linwood Griffin , XenSocket:
A High-Throughput Interdomain Transport for Virtual Machines , Middleware '07 Proceedi
ngs of the ACM/IFIP/USENIX 2007 International Conference on Middleware , 2 0 0 7 年 , p.
184-203

- (58)調査した分野(Int.Cl. , D B 名)
G 0 6 F 1 2 / 0 8 - 1 2 / 1 2
G 0 6 F 1 2 / 1 4