

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4500014号
(P4500014)

(45) 発行日 平成22年7月14日 (2010. 7. 14)

(24) 登録日 平成22年4月23日 (2010. 4. 23)

(51) Int. Cl.	F I
G06F 12/08 (2006.01)	G06F 12/08 531E
	G06F 12/08 551C
	G06F 12/08 523C
	G06F 12/08 519Z

請求項の数 3 (全 13 頁)

(21) 出願番号	特願2003-181949 (P2003-181949)	(73) 特許権者	503003854
(22) 出願日	平成15年6月26日 (2003. 6. 26)		ヒューレット・パカード デベロップメント カンパニー エル. ビー.
(65) 公開番号	特開2004-54931 (P2004-54931A)		アメリカ合衆国 テキサス州 77070
(43) 公開日	平成16年2月19日 (2004. 2. 19)		ヒューストン コンパック センタ ド
審査請求日	平成18年6月15日 (2006. 6. 15)		ライブ ウェスト 11445
(31) 優先権主張番号	10/201, 180	(74) 代理人	110000246
(32) 優先日	平成14年7月23日 (2002. 7. 23)		特許業務法人オカダ・フシミ・ヒラノ
(33) 優先権主張国	米国 (US)	(74) 代理人	100081721
			弁理士 岡田 次生
		(74) 代理人	100105393
			弁理士 伏見 直哉
		(74) 代理人	100111969
			弁理士 平野 ゆかり

最終頁に続く

(54) 【発明の名称】 分散メモリマルチプロセッサシステムにおけるメモリ移行のためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

分散メモリマルチプロセッサシステムであって、

複数のセルが互いに通信可能に結合され、該複数のセルは全体として複数のプロセッサと複数のキャッシュと複数のメインメモリと複数のセルコントローラとを含み、

前記セルの各々は、前記プロセッサのうちの少なくとも1つと、前記キャッシュのうちの少なくとも1つと、前記メインメモリのうちの1つと、前記セルコントローラのうちの1つとを含み、

前記セルの各々は、当該マルチプロセッサシステムのオペレーティングシステムによる介入なしに、前記メインメモリのうちの第1のメインメモリから前記メインメモリのうちの第2のメインメモリにメモリを移行させるメモリ移行機能を実行するように構成され、

前記システムは、前記移行の完了にตอบสนองして、前記プロセッサのアドレスレジスタを更新することにより、新たな要求を前記第1のメインメモリではなく前記第2のメインメモリに向けるように構成される、分散メモリマルチプロセッサシステム。

【請求項 2】

各々が、少なくとも1つのプロセッサと、少なくとも1つのキャッシュと、メインメモリと、セルコントローラと、キャッシュ整合性ディレクトリとを含む、複数のセルを有する分散メモリマルチプロセッサシステムにおいてメモリを移行させる方法であって、

前記セルのうちの第1のセルと前記セルのうちの第2のセルとの間のメモリ移行トランザクションを開始するステップと、

該メモリ移行トランザクション中に、前記第1のセルの前記メインメモリにおける第1のメモリ部から前記第2のセルの前記メインメモリにおける第2のメモリ部に、データをコピーするステップと、

前記第1のセルと前記第2のセルとの前記キャッシュ整合性ディレクトリに、前記メモリ移行トランザクション中の前記第1のメモリ部と前記第2のメモリ部との移行状態を示す移行ステータス情報を格納するステップとを含む方法。

【請求項3】

ディレクトリベースのキャッシュ整合性を備えた分散メモリマルチプロセッサシステムであって、

複数のセルを含み、各セルは少なくとも1つのプロセッサ、少なくとも1つのキャッシュ、メインメモリ、セルコントローラおよびキャッシュ整合性ディレクトリを含み、

セルの各々は、一方のセルのメインメモリから他方のセルのメインメモリへメモリを移行するためのメモリ移行機能を実行するように構成され、

キャッシュ整合性ディレクトリの各々は、前記キャッシュ整合性ディレクトリを含む前記セルの前記メインメモリにおけるメモリの移行状態を示す移行状態情報を蓄積するように構成される、システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、概してコンピュータシステムに関する。より詳細には、分散メモリマルチプロセッサシステムにおけるメモリ移行(migration)に関する。

【0002】

【従来の技術】

従来、メインメモリは、物理的に中央バス上に位置していた。このタイプのシステム内では、完全な物理アドレスからなるメモリ要求がメモリサブシステムに転送され、データが返されていた。分散メモリシステムでは、メインメモリは、物理的に多くの異なるセルにわたって分散される。セルは、複数のプロセッサと、1つまたは複数の入出力(I/O)装置と、セルコントローラと、メモリとからなり得る。各セルは、メインメモリ空間の異なる部分を保持する。各プロセッサは、ローカルメモリだけでなく、1つまたは複数のクロスバススイッチ等のセル通信リンク回路を介して他のセルのメモリにもアクセスすることができる。

【0003】

キャッシングは、メモリアクセスに関連するパフォーマンス上の制限を改良することができる。キャッシングは、メインメモリより小さく高速なキャッシュメモリに、メインメモリの内容のサブセットを格納することを含む。キャッシュ内容がデータに対する要求を予測する確率を増大させるために、あらゆる戦略が使用される。たとえば、メモリアドレス空間における要求されたワードに近いデータは、その要求されたワードと時間的に近接して要求される可能性が比較的高いため、大抵のキャッシュはマルチワードラインをフェッチし格納する。単一のキャッシュラインに格納されるワードの数は、システムのラインサイズを画定する。たとえば、キャッシュラインは、8ワード長であってよい。

【0004】

キャッシュは、通常、メインメモリよりはるかに少ないライン記憶ロケーションを有する。通常、キャッシュデータを所持するメインメモリラインアドレスを一意に指示するために、各キャッシュロケーションにおいてデータとともに「タグ」が格納される。

【0005】

シングルプロセッサシステムとマルチプロセッサシステムとの両方において、キャッシュとメインメモリとの間の「整合性(coherency: コヒーレンシ)」を保証するという課題がある。たとえば、プロセッサがキャッシュに格納されたデータを変更すると、その変更がメインメモリにおいて反映されなければならない。通常、キャッシュにおいてデータが

10

20

30

40

50

変更される時刻とメインメモリにおいて変更が反映される時刻との間に幾分かのレイテンシがある。このレイテンシ中、メインメモリにおける未変更データは無効である。メインメモリデータが無効である間は読み出されないことを保証するための処置がとられなければならない。

【 0 0 0 6 】

各プロセッサまたは入出力モジュールがキャッシュメモリを有する、分散メモリマルチプロセッサシステムの場合、キャッシュメモリを有するシングルプロセッサシステムの場合より状況は幾分か複雑である。マルチプロセッサシステムでは、特定のメインメモリアドレスに対応する現データを、1つまたは複数のキャッシュメモリおよび/またはメインメモリに格納してよい。キャッシュメモリのデータは、プロセッサによって操作された可能性があり、その結果、値がメインメモリに格納された値と異なることになる。このため、いかなるアドレスの現データ値も、そのデータ値がどこに存在するかとは無関係に提供されることを保証するために、「キャッシュ整合性方式 (cache coherency scheme)」が実施される。

10

【 0 0 0 7 】

通常、キャッシュデータを変更するためには「許可」が必要である。通常、データが正確に1つのキャッシュに格納されている場合にのみ、その許可は与えられる。複数のキャッシュに格納されたデータは、しばしば読取専用として扱われる。各キャッシュラインは、そのラインに格納されているデータを変更する許可が与えられるか否かを示す1つまたは複数の状態ビットを含むことができる。状態の正確な特質はシステムによって決まるが、通常、変更する許可を示すために「プライバシー (privacy)」状態ビットが使用される。プライバシービットが「プライベート (private)」を示す場合、1つのキャッシュのみがそのデータを保持し、関連するプロセッサはそのデータを変更する許可を有する。プライバシービットが「パブリック (public)」を示す場合、いかなる数のキャッシュもデータを保持することができ、いかなるプロセッサもデータを変更することができない。

20

【 0 0 0 8 】

マルチプロセッサシステムでは、データを読み出すかまたは変更することを望むプロセッサに対し、通常、あるとすればいずれのキャッシュがそのデータのコピーを有するかと、そのデータの変更に対し許可が与えられるかに関する判断がなされる。「スヌーピング (Snooping)」は、その判断を行うために複数のキャッシュの内容を検査することを含む。要求されたデータがローカルキャッシュにおいて見つからない場合、リモートキャッシュを「スヌーピングする」ことができる。プライベートデータを、別のプロセッサが読み出すことができるようにパブリックにするよう要求するリコール (recall) を発行することができ、もしくは、いくつかのキャッシュのパブリックデータを、別のキャッシュが変更することができるように無効にするリコールを発行することができる。

30

【 0 0 0 9 】

多数のプロセッサおよびキャッシュに対し、網羅的なスヌーピングは、パフォーマンスを低下させる可能性がある。この理由により、分散メモリマルチプロセッサシステムによっては、セル内でスヌーピングし、セル間整合性についてはディレクトリベースのキャッシュ整合性に頼るものがある。ディレクトリベースのキャッシュ整合性を用いる分散メモリマルチプロセッサシステムは、1997年8月25日に出版され、2000年4月25日に発行され、「DISTRIBUTED MEMORY MULTIPROCESSOR COMPUTER SYSTEM WITH DIRECTORY BASED CACHE COHERENCY WITH AMBIGUOUS MAPPING OF CACHED DATA TO MAIN-MEMORY LOCATIONS」と題された米国特許第6,055,610号に記載されている。

40

【 0 0 1 0 】

ディレクトリベースのキャッシュ整合性を使用する分散メモリシステムでは、各セルのメインメモリは、通常、ディレクトリエントリをメモリの各ラインに関連付ける。各ディレクトリエントリは、通常、ラインをキャッシュするセルと、データのラインがパブリックであるかプライベートであるかとを特定する。また、ディレクトリエントリは、データをキャッシュするセル内の特定のキャッシュ (1つまたは複数) を特定してもよく、および

50

ノまたはスヌーピングを使用してセル内のいずれのキャッシュ（1つまたは複数）がデータを有しているかを判断してよい。このように、各セルは、そのメインメモリに格納されたデータのキャッシュされたコピーのロケーションを示すディレクトリを含む。

【0011】

例として、8セルシステムでは、各ディレクトリエントリは9ビット長であってよい。8つのセルの各々について、各々の「サイト(site)」ビットは、各々のセルがラインのキャッシュされたコピーを含むか否かを示す。9番目の「プライベート」ビットは、データがプライベートに保持されているかパブリックに保持されているかを示す。

【0012】

時に、セルからセルへまたは特定のセル内でメモリを移動または移行することが望ましい。たとえば、メモリを、欠陥のあるメモリデバイスから予備のメモリデバイスに移行させることができる。他の例として、1つまたは複数のメモリデバイスを含むボードを、恐らくはそのボードが欠陥のあるコンポーネントを含むため、そのボードがより新しい改訂版によって置き換えられているため、または他の何らかの理由により、システムから取り外す必要のある場合がある。ボードを取り外す前に、ボードからメモリを別のロケーションに移行することが望ましい場合がある。

【0013】

【発明が解決しようとする課題】

メモリ移行は、通常オペレーティングシステム介入によって発生し、メモリは、まず割付解除され、後に所望の宛先に再割付される。かかる従来技術によるメモリ移行技法では、移行されているメモリにアクセスしているプロセスが停止する可能性があり、もしくは、システムが、プロセスが終了するのを待たなければメモリを移行することができない可能性がある。このため、従来技術によるメモリ移行技法は、ソフトウェアの動作に影響を与え、時に、システムをある期間使用不可能にする。さらに、従来技術による移行技法を使用すると、オペレーティングシステムおよびファームウェアが必要とするいくつかのページを容易に移行することができない。

【0014】

また、メモリはインタリーブされる可能性もあり、従来技法を使用するメモリ移行がさらに困難となる。メモリをデインタリーブすることは単純なタスクではなく、時に、デインタリーブソリューションは存在しない。

【0015】

【課題を解決するための手段】

本発明の一形態では、互いに通信可能に結合され、全体として複数のプロセッサとキャッシュとメインメモリとセルコントローラとを含む、複数のセルを有する、分散メモリマルチプロセッサシステムを提供する。セルの各々は、プロセッサのうちの少なくとも1つと、キャッシュのうちの少なくとも1つと、メインメモリのうちの1つと、セルコントローラのうちの1つとを含む。セルの各々は、本システムのオペレーティングシステムに対して不可視である方法で、メインメモリのうちの第1のメインメモリからメインメモリのうちの第2のメインメモリにメモリを移行させるメモリ移行機能を実行するように構成される。

【0016】

【発明の実施の形態】

好ましい実施形態の以下の詳細な説明では、実施形態の一部を形成し、本発明を実施することができる特定の実施形態を例示として示す、添付図面を参照する。他の実施形態を利用してよく、本発明の範囲から逸脱することなしに構造的変更または論理の変更を行ってよい、ということを理解しなければならない。したがって、以下の詳細な説明は、限定する意味でとられるべきものではなく、本発明の範囲は、併記特許請求項によって規定される。

【0017】

図1は、本発明の一実施形態による、オペレーティングシステム介入なしにメモリを移行

10

20

30

40

50

するように構成された分散メモリマルチプロセッサシステム100を示すブロック図である。システム100は、8つのセル102、104、106、108、110、112、114および116を有し、それらはセル通信リンク118を介して通信可能に結合される。セル102は、4つのメモリアクセス装置120A~120Dと、4つのキャッシュ124A~124Dと、高速整合性ディレクトリまたはディレクトリキャッシュ126と、セルコントローラ128と、メインメモリ136とを有する。同様に、セル116は、4つのメモリアクセス装置164A~164Dと、4つのキャッシュ168A~168Dと、高速整合性ディレクトリ150と、セルコントローラ152と、メインメモリ160とを有する。

【0018】

一実施形態では、メモリアクセス装置120A~120Bおよび164A~164Bはプロセッサであり、メモリアクセス装置120C~120Dおよび164C~164Dは入出力(I/O)モジュールである。メモリアクセス装置120A~120Dは、それぞれアドレスレジスタ122A~122Dを有する。メモリアクセス装置164A~164Dは、それぞれアドレスレジスタ166A~166Dを有する。セルコントローラ128は、ファームウェア130と、コンフィギュレーション・ステータスレジスタ(CSR)132と、順序付きアクセスキュー(OAQ)134とを有する。セルコントローラ152は、ファームウェア154と、コンフィギュレーション・ステータスレジスタ156と、順序付きアクセスキュー158と、を有する。本発明の一形態では、セル104、106、108、110、112および114は、実質的にはセル102および116と同じである。

【0019】

本発明の一実施形態を、各セルが複数のプロセッサと複数のI/Oモジュールとを含むマルチセルシステムのコンテキストで説明するが、当業者には、本明細書で説明するメモリ移行技法は他のシステム構成にも適用可能である、ということが明らかとなる。たとえば、本発明の代替実施形態は、図1に示すもののようなセルを有するのではなく、単一プロセッサ(キャッシュ付き)ビルディングブロック、単一I/Oモジュール(キャッシュ付き)ビルディングブロックまたは他のビルディングブロック等、他のシステムビルディングブロックを組み込んでよい。

【0020】

一実施形態によれば、標準動作時、システム100は、ディレクトリベースのキャッシュ整合性を使用して従来からの方法でメモリにアクセスするように構成される。たとえば、プロセッサ120Aによりメインメモリ136からワードがフェッチされると、そのワードはキャッシュ124Aのキャッシュラインに格納される。さらに、要求されたワードに隣接するワードもまた、要求されたワードとともにフェッチされキャッシュラインに格納される。

【0021】

プロセッサ120Aによるデータに対する要求には、システム100のメインメモリローケーションのうちの1つを一意に特定するメインメモリワードアドレスが含まれる。キャッシュ124Aは、メインメモリアドレスの複数の最下位ビットを除去することにより、プロセッサ120Aからのメインメモリワードアドレスをラインアドレスに変換する。このラインアドレスがセルコントローラ128に転送されることにより、要求されたデータの位置が特定される。要求が所有者セルに転送されなければならない場合、セルコントローラ128は、ラインアドレスの複数の最上位ビットをセルIDに復号し、要求を満たすことができるように、アドレスの残りのビットを適当なセルに転送する。一実施形態では、セルコントローラは、セル内およびセル間のすべてのメモリトランザクションを処理する。

【0022】

キャッシュ(たとえば、キャッシュ124A~124Dおよび168A~168D)は、ラインアドレスの複数のビットを使用してキャッシュラインを特定する。そして、ライン

10

20

30

40

50

アドレスの残りのビットが、特定されたキャッシュラインに格納されたタグと比較される。「ヒット」の場合(すなわち、タグがラインアドレスの残りのビットと一致する場合)、ラインアドレスの複数の最下位ビットを使用して、プロセッサ(またはI/Oモジュール)に転送するために、特定されたキャッシュラインに格納されたワードのうちの1つが選択される。ミスの場合(すなわち、タグが一致しない場合)、最終的にメインメモリからフェッチされるラインは、特定されたキャッシュラインにおいてデータのラインに上書きし、特定されたキャッシュラインにおけるタグが更新される。最後に、要求されたワードは、キャッシュラインからプロセッサ(またはI/Oモジュール)に転送される。

【0023】

本発明の一形態では、キャッシュ(たとえば、キャッシュ124A~124Dおよび168A~168D)における各キャッシュラインについて状態ビットを含み、各セルのメインメモリの整合性ディレクトリ(たとえば、整合性ディレクトリ138および162)に整合性情報を格納し、各セルの高速整合性ディレクトリ(たとえば、高速整合性ディレクトリ126および150)に整合性情報を格納することにより、システム100において従来からの方法で整合性が実施される。

【0024】

一実施形態では、各キャッシュラインにタグビットおよびユーザデータビットを格納することに加えて、各キャッシュラインには「有効性」状態ビットおよび「プライベート」状態ビットもまた格納される。有効性状態ビットは、キャッシュラインが有効であるか無効であるかを示す。プライベート状態ビットは、キャッシュラインに格納されたデータがパブリックであるかプライベートであるかを示す。プロセッサは、そのキャッシュのいかなる有効データも読み出すことができる。しかしながら、一実施形態では、プロセッサは、そのキャッシュがプライベートに保持するデータしか変更することができない。プロセッサが、パブリックに保持するデータを変更する必要がある場合、そのデータはまずプライベートにされる。プロセッサがその関連するキャッシュにないデータを変更する必要がある場合、そのデータはそのキャッシュにプライベートとして入れられる。データは、他のキャッシュによって使用中である場合、プライベートにされる前にそのキャッシュからリコールされる。

【0025】

一実施形態では、スヌーピングを使用して、同じセルの他のプロセッサ(またはI/Oモジュール)に関連するキャッシュにおいて要求されたデータのコピーが突き止められる。このため、プロセッサ120Aが、パブリックに保持するデータを変更するよう要求する場合、セルコントローラ128は、スヌーピングを使用してローカルキャッシュ124A~124Dのすべてのコピーのリコールを行う。リコールは、いかなるプライベートに保持されたコピーも可能な限り迅速にパブリックに変換されるように、およびパブリックなコピーが無効化されるように要求する役割を果たす。一旦未解決のデータのコピーがなくなると、データのプライベートなコピーをプロセッサ120Aに提供することができ、あるいは、データのパブリックなコピーをプライベートにすることができる。そして、プロセッサ120Aは、データのそのプライベートなコピーを変更することができる。

【0026】

一実施形態では、セル間整合性は、システム100においてディレクトリベースである。要求をセル内で満足させることができない場合、要求は要求されたデータを所有しているセルのセルコントローラに転送される。たとえば、プロセッサ120Aがメインメモリ160内のアドレスをアサートする場合、セル116は要求されたデータを所有する。セルコントローラ152には、要求されたデータのコピーをシステム全体で探す責任が課される。この探索に必要な情報は、メインメモリ160においてユーザデータとともにラインベースで格納される整合性ディレクトリ162に維持される。本発明の一形態では、メインメモリの各ラインは、サイトビットと状態ビットとを格納する。一実施形態では、サイトビットは、各セルに対しそのセルがラインのコピーを保持するか否かを示し、スヌーピングを使用して、ラインのコピーを保持するセル内の特定のキャッシュが特定される。代

10

20

30

40

50

替実施形態では、サイトビットは、各セルの各キャッシュに対し、そのキャッシュがラインのコピーを保持しているか否かを示す。本発明の一形態では、メインディレクトリ状態ビットは、「プライバシー」状態ビットと「共有」状態ビットとを含む。プライバシーメインディレクトリ状態ビットは、データがパブリックに保持されるかプライベートに保持されるかを示す。共有メインディレクトリ状態ビットは、データが「アイドル」であるか複数のキャッシュによってキャッシュされるか否かを示す。共有状態ビットが、データがアイドルであることを示す場合、そのデータはキャッシュされないか、または単一キャッシュのみによってキャッシュされる。

【0027】

セルコントローラ152は、整合性ディレクトリ162の状態ビットから、システム100のいずれのセルが要求されたデータのコピーを保持しているかと、要求されたデータがプライベートに保持されているかパブリックに保持されているかを判断することができる。それにしたがって、リコールを、特定されたセルに向けることができる。

【0028】

高速ディレクトリ150は、予測リコールを起動することができる。一実施形態では、高速ディレクトリ150はユーザデータ情報を格納しないが、メインメモリ160のメインディレクトリ162に格納された整合性ディレクトリ情報のサブセットを格納する。当業者には、高速ディレクトリを実施し予測リコールを起動する技法は既知である。

【0029】

本発明の一形態では、整合性情報を格納することに加えて、システム100のメインメモリにおける整合性ディレクトリ（たとえば、ディレクトリ138および162）はまた、メモリ移行トランザクション中にラインの移行状態を特定するために使用される移行ステータス情報も格納する。一実施形態では、移行ステータス情報は、4つの移行状態、すなわち（1）Home_Cell_Ownership（ホームセル所有権）、（2）Waiting_For_Migration（移行待機）、（3）In_Migration（移行中）または（4）Migrated（移行済み）のうちの1つを特定する。これらの移行状態の各々については、図2および図3Aないし図3Cを参照して後により詳細に論考する。各セルコントローラ（たとえば、セル102のセルコントローラ128、およびセル116のセルコントローラ152）は、セルコントローラに対し本明細書で説明するようなメモリ移行機能を実行させるファームウェア（たとえば、セルコントローラ128のファームウェア130およびセルコントローラ152のファームウェア154）を含む。

【0030】

一実施形態では、メモリ移行は、キャッシュラインベースで行われる。代替実施形態では、単一ライン以外のメモリサイズを使用してよい。ラインが移行される先のセルを特定するために、「新セル」という用語を使用し、ラインが移行される元のセルを特定するために、「旧セル」という用語を使用する。一実施形態では、移行は同じセル内で発生する可能性があり、そこでは、ラインは、セル内のメモリの1つの物理的口ケーションから同じセル内のメモリの別の物理的口ケーションに移動される。

【0031】

図2は、本発明の一実施形態によるメモリ移行プロセス200を示すフローチャートである。メモリ移行プロセス200では、メモリがセル102からセル116に移行されており、そのためセル102を「旧セル」と呼び、セル116を「新セル」と呼ぶと仮定する。ステップ202において、セルコントローラ128のファームウェア130は、旧セル102のコンフィギュレーション・ステータスレジスタ132にビットが書き込まれるようにし、その後新セル116のコンフィギュレーション・ステータスレジスタ156にビットを書き込むことにより、移行を開始する。旧セル102のコンフィギュレーション・ステータスレジスタ132への書き込みは、セルに対し、メモリがそのセルから移行されていることを通知し、新セル116のコンフィギュレーション・ステータスレジスタ156への書き込みは、セルに対し、メモリがそのセルに移行されていることを通知する。この時点で、プロセッサおよびI/Oモジュールアドレス範囲レジスタ（たとえば、旧セル

10

20

30

40

50

102のレジスタ122A~122Dおよび新セル116のレジスタ166A~166D)は、まだ、旧セル102を移行されるラインの所有者として指している。

【0032】

ステップ204において、セルコントローラ152は、メインメモリ160のメモリラインのうち選択された1つ(すなわち、「新ライン」)に対しディレクトリ162の移行状態を「Waiting_For_Migration」に設定する。ステップ206において、高速ディレクトリ150の新ラインへのいかなる参照もフラッシュされる。一実施形態では、「Waiting_For_Migration」状態の新ラインに対するいかなる要求も不当であり、適当な誤り回復/ロギングステップが起動される。

【0033】

ステップ208において、新セル116のセルコントローラ152は、所望のライン(すなわち「旧ライン」)に対する所有者の変更の意図とともに「フェッチ要求」を旧セル102に送出する。フェッチ要求は、旧セル102に対し、新セル116がそのラインのホームセルとなるよう要求していることを指示する。代替実施形態では、フェッチ要求を、新セル116以外のエンティティによって起動することができる。

【0034】

ステップ210において、セルコントローラ128は、フェッチ要求をその順序付きアクセスキュー134に入れる。ステップ212において、旧ラインに対するいかなる先の要求も、セルコントローラ128によって通常の方法で処理される。フェッチ要求の順番に達すると、ステップ214において、セルコントローラ128はフェッチ要求を処理する。ステップ216において、セルコントローラ128は、旧ラインを、それを所有しているいかなるエンティティ(たとえば、プロセッサまたはI/Oモジュール)からもリコールする。

【0035】

ステップ218において、セルコントローラ128は、フェッチ要求に対する応答を返し、「Migrate_Idle_Data(アイドルデータ移行)」トランザクションを通してホームセル所有権とともに旧ラインデータを転送する。ステップ220において、セルコントローラ128は、旧ラインに対しディレクトリ138の移行状態を「In_Migration」に設定する。この時点で旧ラインに対するいかなる要求も、順序付きアクセスキュー134に入れられる(空きがある場合)か、または「否定応答(nack)される」(すなわち、肯定応答されない)。一実施形態では、旧セル102は、「In_Migration」状態にある間、旧ラインに対するいかなる要求も処理しない。

【0036】

ステップ222において、セルコントローラ152は、「Migrate_Idle_Data」トランザクションを受け取り、「Ack」(すなわち、肯定応答)トランザクションを送出する。ステップ224において、セルコントローラ152は、受け取ったラインデータを新ラインにコピーし、ラインのホーム所有権を想定し、ディレクトリ162においてラインを「アイドル」としてマークする。他のいずれのエンティティもこのラインを有していないため、ラインは「アイドル」としてマークされる。上述したように、ステップ216において、ラインは先のすべての保持者からリコールされた。

【0037】

ステップ226において、セルコントローラ128は、セルコントローラ152から「Ack」トランザクションを受け取り、旧ラインに対しディレクトリ138の移行状態を「Migrated」状態に遷移させる。ここで、ステップ228において、旧ラインに対する順序付きアクセスキューに保留中の任意のアクセスまたは旧ラインに対する任意の新たなアクセスは、セルコントローラ128により旧ラインのための新ホームセル116に向けられる。代替実施形態では、セルコントローラ128は、旧ラインに対する要求者に応答し、それらに対し、それらの要求を新セル116に送るよう要求する。ステップ230において、旧セル102が旧ラインに対する順序付きアクセスキュー134にそれ以上保留中のエントリを有していない場合、セルコントローラ128は、旧ラインに対する保留中要求

10

20

30

40

50

のすべてが処理されたことを示すために、そのコンフィギュレーション・ステータスレジスタ 132 にステータスピットを設定する。一実施形態では、この時点で、旧ラインに対するいかなる新たな要求も順序付きアクセスキュー 134 に入れられていない。旧ラインに対する新たな要求はすべて、セルコントローラ 128 により新セル 116 に転送される。

【0038】

ステップ 232 において、ファームウェア 130 は、コンフィギュレーション・ステータスレジスタ 132 を読み出し、ステータスピットが設定された（ステップ 230）と判断する。そして、ファームウェア 130 は、ステータスピットをリセットする。ステップ 234 において、ファームウェア 130 は、旧セル 102 からさらなるラインが移行されるか否かを判断する。さらなるラインが移行される場合、かかるラインの各々に対しステップ 202 ~ 232 が繰り返される。さらなるラインが移行されるか否かに関らず、ステップ 236 ~ 240 が実行されることにより、第 1 のラインの移行が完了する。

10

【0039】

ステップ 236 において、ファームウェア 130 は、すべてのメモリアクセス装置（たとえば、プロセッサまたは I/O モジュール）においてアドレス範囲レジスタ（たとえば、アドレス範囲レジスタ 122A ~ 122D および 166A ~ 166D）を変更することにより、旧ラインに対する新しい要求がすべて新セル 116 に向かうことを保証する。アドレスレジスタの変更は、即時には発生しない。変更中、旧ラインに対するいずれかの要求が旧セル 102 に送出されると、これらの要求は、旧セル 102 によって新セル 116 に転送される。代替実施形態では、旧セル 102 は、かかる要求を要求者に戻し、要求者に対して要求を新セル 116 に送るように通知する。ステップ 236 における変更後、メモリアクセス装置のすべてに対するアドレス範囲レジスタは、移行されたラインに対するホームセルとして新セル 116 を指す。一実施形態では、アドレス範囲レジスタを、セルマップの形態で実施してよい。

20

【0040】

ステップ 238 において、旧ラインに対する未解決の要求（たとえば、まだ処理されていない、順序付きアクセスキューにおける保留中の要求）がまだいくつかある可能性があるため、ファームウェア 130 は、すべてのあり得る要求者から旧セル 102 に対し「プランジ (plunge)」を起動することにより、旧セル 102 に対するいかなる先の要求も旧セル 102 に達した（および新セル 116 に向けられた）ことを保証する。プランジトランザクションは、あり得るすべてのメモリアクセス装置から旧セル 102 に対して送出される。プランジトランザクションが他のメモリ要求トランザクションと同様に待ち行列に入れられるため、旧セル 102 がプランジトランザクションのすべてを受け取るまでに、ファームウェア 130 は、旧ラインに対するすべての要求が受け取られた（および新セル 116 に転送された）ことを知る。

30

【0041】

ステップ 240 において、ファームウェア 130 は、コンフィギュレーション・ステータスレジスタ 132 が、プランジが完了したことを示すまで待機する。コンフィギュレーション・ステータスレジスタ 132 は、ファームウェア 130 に対し、旧ラインに対しそれ以上未解決の要求がないことを示す。代替実施形態では、ファームウェア 130 は、プランジトランザクションを使用するのではなく、旧ラインに対するすべての未解決の要求が旧セル 102 に達し新セル 116 に転送されるために十分長い所定期間待機する。

40

【0042】

ステップ 242 によって示すように、プランジトランザクションが完了した後（または所定期間が経過した後）、ラインの移行が完了する。

【0043】

図 3A は、本発明の一実施形態による、メモリ移行プロセス 200 中に旧セル 102 から移行されるメモリラインに対するディレクトリ 138 のディレクトリ状態遷移を示す状態図である。図 3A において状態 S1 によって示すように、旧セル 102 における旧ライン

50

に対するディレクトリ 1 3 8 の開始移行状態は、「Home_Cell_Ownership」であり、それは、旧セル 1 0 2 が、移行されるラインに対する現ホームセルであることを示す。旧ラインに対する第 2 の移行状態 S 2 は、「In_Migration」である。上述したように、旧ラインの移行状態は、プロセス 2 0 0 のステップ 2 2 0 中に「Home_Cell_Ownership」から「In_Migration」に遷移する。旧ラインに対する第 3 の移行状態 S 3 は、「Migrated」である。旧ラインの移行状態は、プロセス 2 0 0 のステップ 2 2 6 中に「In_Migration」から「Migrated」に遷移する。

【 0 0 4 4 】

図 3 B は、本発明の一実施形態による、メモリ移行プロセス 2 0 0 中に新セル 1 1 6 におけるメモリラインに対するディレクトリ 1 6 2 のディレクトリ状態遷移を示す状態図である。図 3 B において状態 S 4 によって示すように、新セル 1 1 6 における新ラインに対するディレクトリ 1 6 2 の移行状態は、「Waiting_For_Migration」であり、それはプロセス 2 0 0 のステップ 2 0 4 中に設定される。新ラインに対する次の移行状態 S 5 は、「Home_Cell_Ownership」であり、それはセル 1 1 6 が移行されたラインに対する新ホームセルであることを示す。新ラインの移行状態は、プロセス 2 0 0 のステップ 2 2 4 中に「Waiting_For_Migration」から「Home_Cell_Ownership」に遷移する。

【 0 0 4 5 】

図 3 C は、図 3 A および図 3 B に示す状態の時間順序を示す図である。図 3 C に示すように、メモリ移行プロセス 2 0 0 は、状態 S 1 (Home_Cell_Ownership) の旧セル 1 0 2 の旧ラインで開始する。次に、メモリ移行プロセス 2 0 0 中、新セル 1 1 6 の新ラインの移行状態は、状態 S 4 (Waiting_For_Migration) に設定される。後に、移行プロセス 2 0 0 中、旧セル 1 0 2 における旧ラインの移行状態は、状態 S 2 (In_Migration) に遷移する。次に、新セル 1 1 6 における新ラインの移行状態は、状態 S 5 (Home_Cell_Ownership) に遷移する。最後に、旧セル 1 0 2 における旧ラインの移行状態は、状態 S 3 (Migrated) に遷移する。

【 0 0 4 6 】

一実施形態では、移行シーケンス中、メモリ 1 つまたは複数のアドレス範囲が移行されるが、移行の粒度は 1 メモリラインである。所望のアドレス範囲における複数のメモリラインは、常に種々の異なる移行状態にあってよい。本発明の一形態では、指定された範囲内のすべてのメモリラインが移行されるまで、アドレス範囲レジスタは移行を反映するように更新されない(たとえば、図 2 のステップ 2 3 6)。

【 0 0 4 7 】

本発明の実施形態は、従来技術によるメモリ移行技法に対し多くの利点を提供する。本発明の一形態は、ディレクトリベースのキャッシュ整合性を用いる分散メモリマルチプロセッサシステムにおいて、オペレーティングシステムの関与を必要とすることなしに、インタリーブされるか否かに関らずメモリをロケーションからロケーションに動的に移行させるシステムおよび方法を提供する。一実施形態では、システムファームウェアによるかまたはユーティリティファームウェアによって提供される何らかのソフトウェアの助けを借りて、移行の際にハードウェアを使用する。本発明の一形態では、移行されたメモリにアクセスするソフトウェアは、移行が発生している間にシームレスにメモリにアクセスし続けることができ、そのため、ユーザに対しサービスが中断されない。本発明の一形態では、オペレーティングシステムおよびアプリケーションソフトウェアの関与なしに、またはそれらに悪影響を及ぼすことなく、移行が発生する。一実施形態では、移行はオペレーティングシステムに対して「不可視」であり、オペレーティングシステムとアプリケーションソフトウェアとは、移行が発生しているかまたは発生したことが通知される必要がない。

【 0 0 4 8 】

本発明の一形態は、プロセッサインタフェースから独立したメモリ移行プロセスを提供し、移行機能を実施するためにプロセッサインタフェースプロトコルがいかなる新たなトランザクションをサポートすることも必要としない。メモリ移行プロセスの一実施形態と

10

20

30

40

50

もに、いかなるプロセッサまたはI/Oコントローラ設計をも使用することができる。

【0049】

一実施形態では、本明細書で説明した技法を使用して、さらなるメモリがシステムに追加された時にメモリを移行させることができ、従来からの誤り検出および訂正方式とともに使用して、欠陥のあるメモリロケーションを予備のメモリロケーションに移行させることができる。欠陥のあるメモリが置換されるかまたは他の方法で修復されると、予備のメモリロケーションを新たなまたは修復されたメモリに戻るよう移行させることができる。

【0050】

本発明の好ましい実施形態の説明の目的のために、本明細書では、特定の実施形態を例示し説明したが、当業者には、本発明の範囲から逸脱することなく、示し説明した特定の実施形態を、多種多様な代替および/または等価実施態様に置き換えてよいことは認められよう。化学、機械、電子機械、電気およびコンピュータ技術における当業者は、本発明を多種多様な実施形態で実施してよい、ということ容易に認めるであろう。この出願は、本明細書で論考した好ましい実施形態のいかなる適用形態または変形形態をも包含するように意図される。したがって、この発明は、特許請求項の範囲とその等価物とによってのみ限定されることが明示的に意図されている。

10

【図面の簡単な説明】

【図1】本発明の一実施形態による、メモリを移行するように構成された分散メモリマルチプロセッサシステムを示すブロック図。

【図2】本発明の一実施形態によるメモリ移行プロセスを示すフローチャート。

20

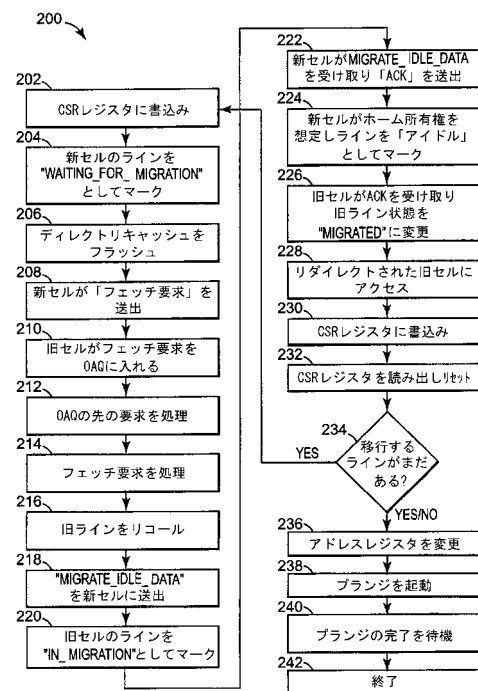
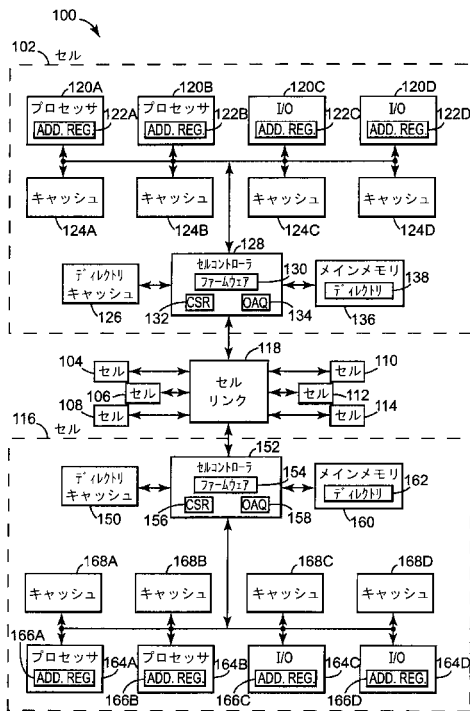
【図3A】本発明の一実施形態によるメモリ移行シーケンス中の「旧セル」に対するディレクトリ状態遷移を示す状態図。

【図3B】本発明の一実施形態によるメモリ移行シーケンス中の「新セル」に対するディレクトリ状態遷移を示す状態図。

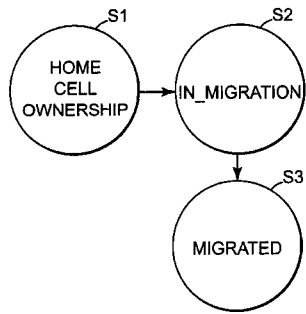
【図3C】図3Aおよび図3Bに示す状態の時間順序を示す図。

【図1】

【図2】



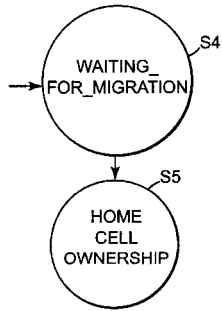
【 3 A 】



【 3 C 】



【 3 B 】



フロントページの続き

- (72)発明者 デベンドラ・ダス・シャルマ
アメリカ合衆国 9 5 0 5 0 カリフォルニア州サンタ・クララ、アカシア・コート 2 0 4 3
- (72)発明者 アシシュ・グプタ
アメリカ合衆国 9 5 1 2 9 カリフォルニア州サン・ノゼ、オラ・ストリート 5 6 3 7
- (72)発明者 ウィリアム・アール・ブライグ
アメリカ合衆国 9 5 0 7 0 カリフォルニア州サラトガ、ペレゴ・ウェイ 1 8 6 3 0

審査官 清木 泰

- (56)参考文献 特開平 1 1 - 2 3 8 0 4 7 (J P , A)
特開平 1 1 - 1 5 4 1 1 5 (J P , A)
特表平 1 0 - 5 0 3 3 1 0 (J P , A)
特開平 0 4 - 2 1 1 8 4 8 (J P , A)
福田宗弘, 村田浩樹, 清水茂則, キャッシュ・オンリ・メモリ・アーキテクチャ向きディレクトリ機構の提案, 情報処理学会研究報告, 日本, 社団法人情報処理学会, 1 9 9 3 年 1 0 月 2 1 日, Vol:93, No:91, (93-ARC-102), Pages:41-48
- (58)調査した分野(Int.Cl., D B 名)
G06F12/08-12/12
G06F15/16-15/177