

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2005-519391

(P2005-519391A)

(43) 公表日 平成17年6月30日(2005.6.30)

(51) Int. Cl.⁷

G06F 12/08

G06F 12/06

F I

G06F 12/08 531B

G06F 12/08 531E

G06F 12/08 551C

G06F 12/06 530A

テーマコード(参考)

5B005

5B060

審査請求有 予備審査請求有 (全21頁)

(21) 出願番号 特願2003-573550 (P2003-573550)
 (86) (22) 出願日 平成14年2月28日(2002.2.28)
 (85) 翻訳文提出日 平成16年10月25日(2004.10.25)
 (86) 国際出願番号 PCT/US2002/005779
 (87) 国際公開番号 W02003/075162
 (87) 国際公開日 平成15年9月12日(2003.9.12)
 (81) 指定国 EP(AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), JP

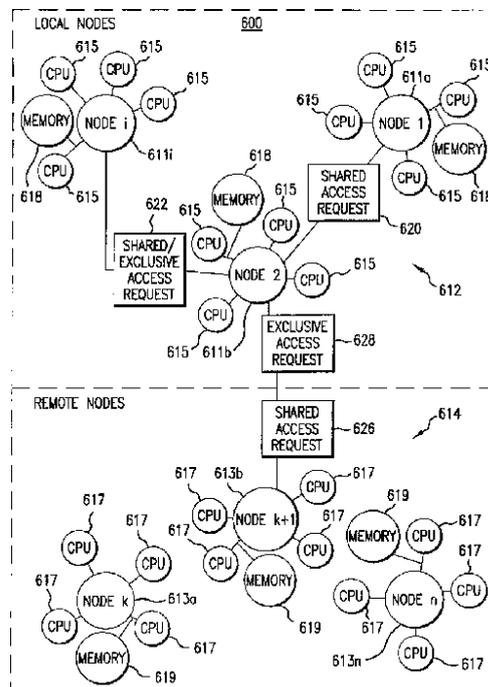
(71) 出願人 597132333
 シリコン、グラフィクス、インコーポレイテッド
 SILICON GRAPHICS, INC.
 アメリカ合衆国カリフォルニア州、マウンテンビュー、アンフィシアター、パークウェイ、1600
 (74) 代理人 100078282
 弁理士 山本 秀策
 (74) 代理人 100062409
 弁理士 安村 高明
 (74) 代理人 100113413
 弁理士 森下 夏樹

最終頁に続く

(54) 【発明の名称】 共有ベクトルの増加を伴わないDSMマルチプロセッサシステムにおけるキャッシュコヒーレンスのための方法およびシステム

(57) 【要約】

本発明は、分散共有記憶(DSM)マルチプロセッサシステム(600)において、キャッシュコヒーレンスを保持するための方法とシステムを目的としている。この方法は、受信ノードによって、共有アクセス要求(620、626)が受信されることによって開始される。ここで、受信ノードとは、アクセスしたい情報を持つ少なくとも1つの主記憶装置を有する任意のノードである。次に、その共有アクセス要求が、ローカルノード(611)から発生したのかりモートノード(613)から発生したのかを判断する。共有アクセス要求がローカルノードから発生した場合、共有アクセス要求は、共有アクセス要求として処理される。共有アクセス要求が認可され、共有ベクトルが生成または更新されて共有ローカルノードに反映される。



【特許請求の範囲】**【請求項 1】**

分散共有記憶 (DSM) マルチプロセッサシステムにおけるキャッシュコヒーレンスのための方法であって、

(1) 情報の共有アクセス要求を受信するステップと、

(2) 前記共有アクセス要求は、ローカルノードから発生したのか、リモートノードから発生したのかを判断するステップと、

(3) 前記共有アクセス要求が前記リモートノードから発生した場合、前記共有アクセス要求を排他的アクセス要求として処理するステップと、

(4) 前記共有アクセス要求が前記ローカルノードから発生した場合、前記共有アクセス要求を処理するステップとを有する方法。 10

【請求項 2】

前記共有アクセス要求が前記ローカルノードから発生したのか、前記リモートノードから発生したのかを判断する前記ステップは、さらに前記情報が格納されているメモリ位置を判断するステップを有することを特徴とする請求項 1 に記載の方法。

【請求項 3】

排他的アクセス要求として前記共有アクセス要求を処理するステップは、さらに、他のローカルノードまたはリモートノードへ介入要求を送信し、他のローカルノードまたはリモートノードから前記要求情報のコピーを削除するステップを有することを特徴とする請求項 1 に記載の方法。 20

【請求項 4】

前記共有アクセス要求が前記ローカルノードから発生したのか、前記リモートノードから発生したのかを判断する前記ステップは、さらに共有アクセスリクエストのアドレス上位ビットと受信ノードのアドレス上位ビットを比較するステップを有することを特徴とする請求項 1 に記載の方法。

【請求項 5】

前記排他的アクセス要求として、前記共有アクセス要求を処理するステップは、さらに前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットとが一致する場合には、前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項 4 に記載の方法。 30

【請求項 6】

前記共有アクセス要求が前記ローカルノードから発生したのか、前記リモートノードから発生したのかを判断する前記ステップは、さらに、共有アクセス要求のアドレスと受信ノードに格納されているアドレス表とを比較するステップを有することを特徴とする請求項 1 に記載の方法。

【請求項 7】

前記排他的アクセス要求として、前記共有アクセス要求を処理するステップは、さらに、前記共有アクセスリクエストの前記アドレスと、前記受信ノードに格納されている前記アドレス表中の少なくとも 1 つのアドレスとが一致する場合には、前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項 6 に記載の方法。 40

【請求項 8】

さらに、

(5) 前記共有アクセス要求が前記排他的アクセス要求に変換される場合は、他のローカルノードまたはリモートノードから前記要求情報のコピーを削除するステップと、

(6) 前記共有アクセス要求がリモートから発生する場合には、共有アクセス要求へのポインタ、前記共有アクセス要求の関連位置または前記情報のキャッシュ場所を格納するステップと、

(7) 前記共有アクセス要求は有効かどうかを判断するステップと、

(8) 前記排他的アクセス要求が最早、有効でない場合には、排他的アクセスを終了す 50

るステップと、

(9) 前記共有アクセス要求によって要求された情報をキャッシュへ戻すステップと、
を有することを特徴とする請求項1に記載の方法。

【請求項9】

リモートノードからの共有アクセス要求を排他的アクセス要求として処理した場合において、さらに、

(5) 前記ローカルノードから後続の共有アクセス要求を受信ノードで受信するステップと、

(6) 前記ローカルノードが前記共有アクセス要求を送信する場合には、前記受信ノードからの介入要求を、前記リモートノードへ送信するステップと、

(7) 前記リモートノードの排他的アクセスを送信するステップと、

(8) 前記ローカルノードの前記共有アクセス要求を前記排他的アクセス要求として処理するステップとを有することを特徴とする請求項1に記載の方法。

【請求項10】

ステップ(6)は、さらに前記ローカルノードの共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットとを比較するステップを有することを特徴とする請求項9に記載の方法。

【請求項11】

ステップ(8)は、さらに前記ローカルノードの前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットとが一致する場合には、前記ローカルノードの前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項10に記載の方法。

【請求項12】

別の介入要求を送信するステップ(6)は、さらに、前記ローカルノードの前記共有アクセス要求のアドレス上位ビットと前記受信ノードに格納されているアドレス表とを比較するステップを有することを特徴とする請求項9に記載の方法。

【請求項13】

ステップ(8)は、さらに前記ローカルノードの前記共有アクセス要求の前記アドレスと、前記受信ノードに格納されている前記アドレス表中の少なくとも1つのアドレスとが一致する場合には、前記ローカルノードの前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項12に記載の方法。

【請求項14】

さらに、

(9) 前記ローカルノードの前記共有アクセス要求へのポインタを格納するステップと、

(10) 前記ローカルノードの前記共有アクセス要求が有効かどうかを判断するステップと、

(11) 前記排他的アクセス要求が最早、有効でない場合には、前記ローカルノードの前記排他的アクセス要求を終了するステップと、

(12) 前記ローカルノードの前記共有アクセス要求によって要求された情報を戻すステップとを有することを特徴とする請求項9に記載の方法。

【請求項15】

ディレクトリに基づくプロトコルを有する分散共有記憶(DSM)マルチプロセッサシステムにおけるキャッシュコヒーレンスのための方法であって、

(1) 受信ノードによって、情報の共有アクセス要求を受信するステップと、

(2) 前記共有アクセス要求は、ローカルノードから発生したのか、リモートノードから発生したのかを判断するステップと、

(3) 前記共有アクセス要求が前記リモートノードから発生した場合、前記共有アクセス要求を、排他的アクセス要求として処理するステップと、

(4) 前記共有アクセス要求が前記ローカルノードから発生した場合、前記共有アクセ

10

20

30

40

50

ス要求を処理するステップと、

(5) 他のローカルノードまたはリモートノードへ介入要求を送信し、他のローカルノードまたはリモートノードから前記要求情報のコピーを削除するステップと、

(6) 前記共有アクセス要求が前記排他的アクセス要求に変換される場合は、他のローカルノードまたはリモートノードから前記要求情報のコピーを削除するステップと、

(7) 前記共有アクセス要求がリモートから発生する場合には、共有アクセス要求へのポインタを格納するステップと、

(8) 前記共有アクセス要求は有効かどうかを判断するステップと、

(9) 前記排他的アクセス要求が最早、有効でない場合には、排他的アクセスを終了するステップと、

(10) 前記共有アクセス要求によって要求された情報をキャッシュへ戻すステップとを有することを特徴とする方法。

10

【請求項16】

ステップ(2)は、さらに前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットを比較するステップを有することを特徴とする請求項15に記載の方法。

【請求項17】

ステップ(3)は、さらに前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットとが一致する場合には、前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項16に記載の方法。

20

【請求項18】

ステップ(2)は、さらに前記共有アクセス要求のアドレスと前記受信ノードに格納されているアドレス表とを比較するステップを有することを特徴とする請求項15に記載の方法。

【請求項19】

ステップ(3)は、さらに前記共有アクセスリクエストの前記アドレスと、前記受信ノードに格納されている前記アドレス表中の少なくとも1つのアドレスとが一致する場合には、前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項18に記載の方法。

【請求項20】

さらに、

(11) 前記ローカルノードから共有アクセス要求を受信するステップと、

(12) 前記ローカルノードが前記共有アクセス要求を送信する場合には、前記受信ノードからの介入要求を、前記リモートノードへ送信するステップと、

(13) 前記リモートノードの前記排他的アクセスを送信するステップと、

(14) 前記ローカルノードの前記共有アクセス要求を前記排他的アクセス要求として処理するステップとを有することを特徴とする請求項15に記載の方法。

30

【請求項21】

ステップ(12)は、さらに前記ローカルノードの前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットとを比較するステップを有することを特徴とする請求項20に記載の方法。

40

【請求項22】

ステップ(14)は、さらに前記ローカルノードの前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードのアドレス上位ビットとが一致する場合には、前記ローカルノードの前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項21に記載の方法。

【請求項23】

ステップ(12)は、さらに前記ローカルノードの前記共有アクセスリクエストのアドレス上位ビットと前記受信ノードに格納されているアドレス表とを比較するステップを有することを特徴とする請求項20に記載の方法。

50

【請求項 24】

ステップ(14)は、さらに前記ローカルノードの前記共有アクセスリクエストの前記アドレスと、前記受信ノードに格納されている前記アドレス表中の少なくとも1つのアドレスとが一致する場合には、前記ローカルノードの前記共有アクセス要求を前記排他的アクセス要求に変換するステップを有することを特徴とする請求項23に記載の方法。

【請求項 25】

さらに、

(15)前記ローカルノードの前記共有アクセス要求へのポインタを格納するステップと、

(16)前記ローカルノードの前記共有アクセス要求が有効かどうかを判断するステップと、

(17)前記排他的アクセス要求が最早、有効でない場合には、前記ローカルノードの前記排他的アクセス要求を終了するステップと、

(18)前記ローカルノードの前記共有アクセス要求によって要求された情報をキャッシュへ戻すステップとを有することを特徴とする請求項20に記載の方法。

【請求項 26】

分散共有記憶(DSM)マルチプロセッサシステムにおけるキャッシュコヒーレンスのための方法であって、

(1)共有情報を有する受信ノードによって、第1共有アクセス要求を受信するステップと、

(2)前記受信ノードによって、第2共有アクセス要求を受信するステップと、

(3)第1共有アクセスリクエストおよび第2共有アクセスリクエストの一方を、排他的アクセス要求として処理するステップと、

(4)前記第1共有アクセスリクエストおよび前記第2共有アクセスリクエストの他方へ前記介入要求を送信するステップとを有することを特徴とする方法。

【請求項 27】

ステップ(3)は、さらに、

(5)前記第1共有アクセスリクエストおよび前記第2共有アクセスリクエストの一方は、ローカルノードから発生したのか、リモートノードから発生したのかを判断するステップと、

(6)前記第1共有アクセス要求が前記リモートノードから発生した場合、前記第1共有アクセス要求を、前記排他的アクセス要求として処理するステップと、

(7)第2共有アクセスリクエストへ介入要求を送信するステップとを有することを特徴とする請求項26に記載の方法。

【請求項 28】

ステップ(3)は、さらに、

(5)前記第1共有アクセスリクエストおよび前記第2共有アクセスリクエストの一方は、ローカルノードから発生したのか、リモートノードから発生したのかを判断するステップと、

(6)前記第1共有アクセス要求が前記リモートノードから発生した場合、前記第2共有アクセス要求を、前記排他的アクセス要求として処理するステップと、

(7)第1共有アクセスリクエストへ介入要求を送信するステップとを有することを特徴とする請求項26に記載の方法。

【請求項 29】

さらに、

前記第1共有アクセス要求および前記第2共有アクセス要求の他方が、ローカルノードから発生した場合には、前記第1共有アクセスリクエストおよび前記第2共有アクセスリクエストの一方を共有アクセス要求として処理するステップを有することを特徴とする請求項26に記載の方法。

【請求項 30】

10

20

30

40

50

さらに、

前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の他方が、ローカルノードから発生した場合には、リモートノードから発生した前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の一方の排他的アクセスを無効とするステップを有することを特徴とする請求項 26 に記載の方法。

【請求項 31】

前記無効にするステップは、さらに、

(5) 前記第 1 共有アクセスリクエストおよび前記第 2 共有アクセスリクエストの一方のアドレス上位ビットと受信ノードのアドレス上位ビットを比較するステップと、

(6) 前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の一方のアドレス上位ビットと受信ノードのアドレス上位ビットとが一致する場合には、前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の一方を前記排他的アクセス要求として処理するステップとを有することを特徴とする請求項 30 に記載の方法。

10

【請求項 32】

前記無効にするステップは、さらに、

(5) 前記第 1 共有アクセスリクエストおよび前記第 2 共有アクセスリクエストの一方のアドレス上位ビットと受信ノードに格納されたアドレス表とを比較するステップと、

(6) 前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の一方のアドレス上位ビットと受信ノードに格納されたアドレス表とが一致する場合には、前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の一方を前記排他的アクセス要求として処理するステップとを有することを特徴とする請求項 30 に記載の方法。

20

【請求項 33】

さらに、

(7) 前記第 1 共有アクセスリクエストおよび前記第 2 共有アクセスリクエストの一方へのポインタを格納するステップと、

(8) 前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の前記一方が有効かどうかを判断するステップと、

(9) 前記第 1 共有アクセス要求および前記第 2 共有アクセス要求の一方が最早、有効でない場合には、前記第 1 共有アクセスリクエストおよび前記第 2 共有アクセスリクエストの一方の排他的アクセスを終了するステップと、

30

(10) 前記第 1 共有アクセスリクエストおよび前記第 2 共有アクセスリクエストの一方によって要求された情報をキャッシュへ戻すステップとを有することを特徴とする請求項 32 に記載の方法。

【請求項 34】

分散共有記憶 (DSM) マルチプロセッサシステムにおけるキャッシュコヒーレンスを保持するためのシステムであって、

複数の処理ノードであって、各々の前記処理ノードは少なくとも 1 つの中央処理装置 (CPU) および記憶装置に接続されている処理ノードと、

共有アクセス要求を受信するための少なくとも 1 つの受信ノードと、

前記共有アクセス要求が前記ローカルノードから発生したか、または前記リモートノードから発生したかを判断するための手段と、

40

複数の共有メモリベクトルであって、各々の前記共有メモリベクトルは、前記受信ノードから情報共有を要求する前記ローカルノードまたは前記受信ノードから情報共有を要求する前記リモートノードのいずれかを定義する共有メモリベクトルと、

排他的メモリアクセスポインタであって、前記排他的メモリアクセスのポインタは、前記受信ノードの情報に対して排他的アクセスを認可された前記 DSM マルチプロセッサシステムにおける前記リモートノードを指示すメモリアクセスポインタと、

前記排他的メモリアクセスポインタを格納するための記憶装置と、

前記排他的メモリアクセスポインタが無効になった場合には、前記排他的メモリアクセスポインタを削除するための削除手段とを有することを特徴とするシステム。

50

【請求項 35】

前記システムは、ディレクトリに基づくプロトコル環境を使用して動作する請求項 34 に記載のシステム。

【請求項 36】

前記削除手段はさらに、一度、排他的メモリアクセスポインタが前記 DSM マルチプロセッサシステムにおける前記リモートノードを指し示すと、前記他の処理ノードへ介入要求を送信する手段を有することを特徴とする請求項 34 に記載のシステム。

【請求項 37】

前記システムはさらに、前記受信ノードが前記ローカルノードから前記共有アクセス要求を受信した場合には、前記リモートノードへ別の介入要求を送信する手段を有することを特徴とする請求項 34 に記載のシステム。 10

【請求項 38】

ディレクトリに基づくプロトコル環境を使用して、分散共有記憶 (DSM) マルチプロセッサシステムにおけるキャッシュコヒーレンスを保持するためのシステムであって、

複数の処理ノードであって、各々の前記処理ノードは少なくとも 1 つの中央処理装置 (CPU) および記憶装置に接続されている処理ノードと、

共有アクセス要求を受信するための少なくとも 1 つの受信ノードと、

前記共有アクセス要求が前記ローカルノードから発生したか、または前記リモートノードから発生したかを判断するための手段と、

複数の共有メモリベクトルであって、各々の前記共有メモリベクトルは、前記受信ノードから情報共有を要求する前記ローカルノードまたは前記受信ノードから情報共有を要求する前記リモートノードのいずれかを定義する共有メモリベクトルと、 20

排他的メモリアクセスポインタであって、前記排他的メモリアクセスのポインタは、前記受信ノードの情報に対して排他的アクセスを認可された前記 DSM マルチプロセッサシステムにおける前記リモートノードを指示すメモリアクセスポインタと、

前記排他的メモリアクセスポインタが、一度、前記 DSM マルチプロセッサシステムにおける前記リモートノードを指示すと、前記他の処理ノードへ介入要求を送信する手段と

、前記排他的メモリアクセスポインタを格納するための記憶装置と、

前記排他的メモリアクセスポインタが無効になった場合には、前記排他的メモリアクセスポインタを削除するための削除手段とを有することを特徴とするシステム。 30

【請求項 39】

ディレクトリに基づくプロトコル環境を使用して、分散共有記憶 (DSM) マルチプロセッサシステムにおけるキャッシュコヒーレンスを保持するためのシステムであって、

情報の共有アクセス要求を受信する手段と、

前記共有アクセス要求は、ローカルノードから発生したのか、リモートノードから発生したのかを判断する手段と、

前記共有アクセス要求が前記リモートノードから発生した場合、前記共有アクセス要求を、排他的アクセス要求として処理する手段と、

前記共有アクセス要求が前記ローカルノードから発生した場合、前記共有アクセス要求を処理する手段とを有することを特徴とするシステム。 40

【請求項 40】

ディレクトリに基づく環境を使用して、分散共有記憶 (DSM) システムにおけるキャッシュコヒーレンスを保持するためのシステムであって、

情報のアクセス要求を受信する手段と、

前記アクセス要求は、ローカルノードから発生したのか、リモートノードから発生したのかを判断する手段と、

前記アクセス要求が前記ローカルノードから発生したか前記リモートノードから発生したかによって、前記アクセス要求を、異なった仕方で処理する手段と、

を処理する手段とを有することを特徴とするシステム。 50

【請求項 4 1】

コンピュータプログラムロジックをその中に備えるコンピュータ使用可能媒体を有するコンピュータプログラムプロダクトであって、前記コンピュータプログラムロジックは、分散共有記憶(DSM)システムにおけるキャッシュコピーを保持することを、コンピュータシステムによって可能にし、前記コンピュータプログラムロジックは、

情報の共有アクセス要求を受信するのをコンピュータで可能にする第1機能と、

前記共有アクセス要求は、ローカルノードから発生したのか、リモートノードから発生したのかを判断するのをコンピュータで可能にする第2機能と、

前記共有アクセス要求が前記リモートノードから発生した場合、前記共有アクセス要求を、排他的アクセス要求として処理するのをコンピュータで可能にする第3機能と、

前記共有アクセス要求が前記ローカルノードから発生した場合、前記共有アクセス要求を処理するのをコンピュータで可能にする第4機能とを有することを特徴とするコンピュータプログラムプロダクト。

10

【発明の詳細な説明】

【技術分野】

【0001】

(発明の分野)

本発明は、分散共有記憶(DSM)マルチプロセッサシステムにおけるキャッシュコピーのための方法およびシステムに関する。

【背景技術】

20

【0002】

(背景技術)

単一プロセッサコンピュータシステムおよびマルチプロセッサコンピュータシステムを含むコンピュータシステムにおいては、一般的に一時点で、多数のプロセスまたはスレッドが実行される。各プロセスは、いくらかの物理メモリを必要とする。多くの場合、物理メモリは限られていて、異なるプロセス間で割り当てられることが必要である。

【0003】

コンピュータシステムは、一般にメモリアクセスタイムを低減するために、主記憶と各プロセッサとの間に、一段以上のキャッシュメモリを使用している。キャッシュメモリは、主記憶から取り出された情報を格納する。プロセッサによって取り出された情報は、プロセッサに達するために、一段以上のキャッシュを経る必要がある。キャッシュは、小さくてプロセッサに物理的に接近していることが多く、プロセッサと同一チップ上に位置することがある。このため、キャッシュされた情報は、概して主記憶に格納された情報より、はるかに高速でアクセス可能である。よって、キャッシュは一般にプロセッサによって、繰り返しアクセスが必要な情報を格納する。

30

【0004】

キャッシュの整合性を保持するためのシステムおよび方法には、ディレクトリプロトコルを含み、このディレクトリプロトコルにおいて、集中ディレクトリにメモリの状態が保持されている。情報は、その情報を「共有」している異なるプロセスによって、多数の場所にキャッシュされている。あるいは、1つのプロセスは、ある期間中、その情報に対して「排他的」権限をもつことができる。主記憶またはキャッシュ場所のいずれかにある情報をプロセスが変更した場合は、その情報の他のインスタンスを無効とするかまたは更新する必要がある。これを、キャッシュ整合性の保持と呼ぶ。分散共有記憶(DSM)システムにおいては、ディレクトリを分散することができる。集中制御装置が、共有情報の整合性の保持を担当する。あるメモリ位置に格納されている情報が変更された場合はいつでも、集中ディレクトリのチェックが行われ、当該情報のコピーがキャッシュに格納されているかどうか判断される。格納されている場合には、各コピーは、更新されるかまたは無効にされる。

40

【0005】

大規模なDSMマルチプロセスシステムにおいては、キャッシュコピーの保持は

50

、困難な作業となることがある。D S Mシステムでは、一般的に、当該情報をどこにキャッシュしているかを識別するのに、共有ベクトルを使用している。しかし、D S Mシステムの規模が増加するにつれて（例えば、プロセスまたは処理ノードの数）、共有ベクトルが、それにつれて増大し、これにより処理速度が低下し、キャッシュにある情報が他のプロセッサによって使用できない期間が増える。

【発明の開示】

【発明が解決しようとする課題】

【0006】

したがって、D S Mマルチプロセッサシステムにおけるキャッシュコヒーレンスに関する、より優れたシステムが必要となっている。

10

【課題を解決するための手段】

【0007】

（発明の要旨）

本発明は、分散共有記憶（D S M）マルチプロセッサシステムにおけるキャッシュコヒーレンスのためのシステムおよび方法に関する。D S Mマルチプロセッサシステムは、複数のノードを有することがある。各ノードは、さらに少なくとも1つの中央処理装置（C P U）、キャッシュ記憶装置およびオプションの入出力（I / O）装置を有する。D S Mマルチプロセッサシステムは、一般に、プロセッサに対して制御（例えば、スレッドスケジューリング、メモリアクセス等）を維持する1つ以上のオペレーティングシステムを備えている。

20

【0008】

一般に、リクエスト（例えば、プロセッサ上で実行されるプロセス）が、情報へのアクセスを要求すると、その情報のコピーが、その情報を格納しているメモリアドレスからリクエストに関連付けられたキャッシュアドレスへ送信される。次に、リクエストは、キャッシュアドレスから当該情報をアクセスする。キャッシュコヒーレンスの目的のために、状態表示は一般に当該メモリアドレスに関連付けられていて、当該情報をキャッシュしている場所を示している。

【0009】

アクセス要求は、排他的アクセスまたは共有アクセスとすることができる。上記の状態表示は、一般にそのリクエストに認可されたアクセスの種類を表示を含む。例えば、排他的アクセスが認可された場合、その状態表示には当該リクエストに対して、排他的アクセスが認可されたことが識別される。共有アクセスが認可された場合には、当該状態表示は、一般に当該情報のコピーが存在する1つ以上のキャッシュ位置を示す共有ベクトルを含んでいる。付加的な規則および手順を、競合する規則または/およびその他のシナリオを処理するために実装することができる。

30

【0010】

当今のD S Mマルチプロセッサシステムは、数十、数百、さらに数千ものC P Uを備えることができる。より多くのリクエストが当該情報の共有アクセスを取得すると、その共有ベクトルのサイズが増加し、場合により劇的に増加する。これにより、メモリ空間が占有され、処理速度が減少する。

40

【0011】

本発明では、ローカルリクエストおよびリモートリクエストの概念を導入することによって、発生し得る共有ベクトルのサイズを減少する。一般的に、「ローカル」および「リモート」という用語は、リクエストの要求されたメモリ位置への物理的な近接度に対して定義される。リモートおよびローカルに対するその他の定義も使用可能である。

【0012】

運用中、リクエストが情報の共有アクセスを要求すると、当該情報が格納されているメモリ位置に関してそのリクエストがローカルかまたはリモートかを判断する。リクエストがローカルの場合、その共有アクセス要求は、通常共有アクセス要求手順に従って、処理される。リクエストがリモートの場合には、その共有アクセス要求は、排他的要求に変

50

換され排他的要求の手順に従って処理される。したがって、リモートリクエストのキャッシュ整合性を維持するのに共有アクセスベクトルは必要ではない。このことは、メモリ要求および処理要求をかなり低減し、特に大規模なDSMマルチプロセッサシステムにおいてはこの低減が著しい。

【0013】

ある実施の形態においては、追加の規則を実行できる。例えば、前述のように、リモートリクエストからの共有アクセス要求が排他的アクセス要求に変換され、ローカルリクエストから後続の共有アクセス要求が受信された状況に対して、新しい規則を実行することができる。新しい規則の例としては、リモートリクエストに認可された排他的アクセスを終了し、当該ローカルリクエストへの共有アクセスを認可することである。当該状態表示が、排他的アクセスの表示から共有ベクトルの表示に変更され、当該ローカルリクエストへの共有アクセスが表示される。リモートリクエストが、その情報を再度必要とする場合は、再要求をする必要がある。

10

【0014】

(好ましい実施形態の詳細な説明)

(1. 概要)

本発明は、分散共有記憶(DSM)マルチプロセッサシステムにおいて、キャッシュコヒーレンスを向上するための方法およびシステムに関する。DSMシステムは一般に、複数のノードを含み、これらの複数のノードは、これらのノードに接続された複数の中央処理装置(CPU)および記憶装置ならびにノードから情報を送受信するための入出力(I/O)モジュールを含む。ノード内の記憶装置は、種々の機能や手順を実行するために、他のノードまたはリクエストで必要な情報を含むことがある。1つ以上のリクエストが、要求する情報への共有アクセスを認可された場合には、共有ベクトルが共有リクエストを識別する。この共有ベクトルは、キャッシュ整合性を維持するために使用される。当該リクエストが、要求する情報への排他的アクセスを認可されると、ポインタ(要求元ノードの2進数)が、共有ベクトルの代わりに格納される。例えば、当該共有ベクトルの大きさがNビットならば、ポインタの大きさが $\log(N)$ ビットである。ここでNはシステム内のノード数である。

20

【0015】

共有リクエストの数が増加すると、共有ベクトルの大きさが増大する。共有ベクトルの大きさが増大するにつれて、より多くのメモリを消費し、システム全体および各特定のノードの機能と手順の処理速度を低下させる。したがって、共有ベクトルの大きさを低減するより良好なシステムが必要である。

30

【0016】

本発明によれば、要求情報が存在するノードから離れたノードから、共有アクセスの要求が発生すると、当該共有アクセス要求は排他的アクセス要求に変換される。これによって、当該要求情報の共有ベクトルを保持する必要がなくなる。ある実施例では、リクエストが要求された情報から離れているかどうかの判定は、当該リクエストおよび当該要求情報の物理アドレスに基づく。

【0017】

1つの実施例では、アドレス比較は、要求元ノードのアドレス上位ビットと要求された情報のアドレス上位ビットとの間で行われる。別の実施例では、あるノードは共有アクセスを認可されたり、されていない他のノードのアドレス表を格納している。換言すれば、1つのノードが別のノードへの共有情報アクセス要求を送信すると、後者のノードは要求元ノードのアドレスをアドレス表と比較する。当該アドレスが表中のいずれかのアドレス一致すると、共有情報アクセス要求は認可されるか、または排他的要求に変換される。

40

【0018】

ある実施例において、あるローカルノードが、リモートノードによって排他的に保持されている情報へのアクセス要求を送信すると、当該リモートノードのアクセスは除去され、そのローカルノードに対するアクセスが認可される。

50

【 0 0 1 9 】

(2 . 事例環境)

本発明は、分散共有記憶 (D S M) マルチプロセッサシステムにおいて実現できる。しかし、上記の実現は、本発明の使用を D S M システムに限定するものとして考えられるべきではない。

以下は、本発明が実現される環境を理解するのに有用となるシステムの検討である。

【 0 0 2 0 】

図 1 において、単一プロセッサシステム 1 1 0 は、主記憶 1 1 2 へ接続された単一のプロセッサおよびキャッシュノード 1 1 4 を備える。主記憶 1 1 2 は、プロセッサおよびキャッシュノード 1 1 4 によって使用される情報を格納する。入出力 I / O システム 1 1 6 は、例えばコンピュータ端末および記憶ディスクを含むユーザインタフェースおよび記憶装置などの周辺装置へのインタフェースを提供する。

10

【 0 0 2 1 】

図 4 において、プロセッサおよびキャッシュノード 1 1 4 は、プロセッサおよびキャッシュノード 4 1 0 として実現できる。プロセッサおよびキャッシュノード 4 1 0 は、主記憶に接続されたプロセッサ 4 1 2 を備え、これはキャッシュ 4 1 6 を介した主記憶 1 1 2 であってもよい。明確さのために、単一のプロセッサ 4 1 2 およびキャッシュメモリ 4 1 6 を表示してある。当業者は、マルチプロセッサおよび多数段のキャッシュが使用できることを理解するであろう。

【 0 0 2 2 】

キャッシュ 4 1 6 は、主記憶 1 1 2 のような主記憶から取り出した情報を、キャッシュするために設けられている。情報がキャッシュ 4 1 6 にキャッシュされるとすぐに、プロセッサ 4 1 2 は、キャッシュ 4 1 6 から情報を取り出すことができる。プロセッサ 4 1 2 は、一般に、主記憶 1 1 2 内の情報のアクセスより高速で、キャッシュ 4 1 6 から情報を取り出すことができる。この理由は、キャッシュ 4 1 6 がプロセッサ 4 1 2 に近接しており、またキャッシュ 4 1 6 を構成しているメモリ部品が、主記憶を構成しているメモリ部品より高速であるからである。キャッシュ 4 1 6 は、ユーザの必要性によって指定された 1 段以上のキャッシュを備えることができる。

20

【 0 0 2 3 】

運用中、プロセッサ 4 1 2 は、1 つ以上のプロセスのスレッドを処理する。プロセッサ 4 1 2 が、主記憶 1 1 2 に格納されている情報のアクセスを必要とするときは、アクセス要求が送信される。プロセッサ 4 1 2 が、要求した情報へのアクセスを許可されると、主記憶 1 1 2 は、要求された情報をキャッシュ 4 1 6 へ返却する。一度、要求された情報がキャッシュ 4 1 6 に格納されると、プロセッサ 4 1 2 は、その情報を必要に応じてアクセスできる。以後、プロセッサ 4 1 2 は主記憶 1 1 2 をアクセスせずに、キャッシュ 4 1 6 内の情報をアクセスできる。

30

【 0 0 2 4 】

図 2 において、集中化共有メモリ対称型マルチ処理 (S M P) システム 2 1 0 の例は、複数のプロセッサおよびキャッシュノード 2 1 2 から 2 1 8 を備えている。S M P 2 1 0 は、任意の数のノード 2 1 2 から 2 1 8 を含むことができる。プロセッサおよびキャッシュノード 2 1 2 から 2 1 8 は、バス 2 2 2 を介して、集中化共有主記憶 2 2 0 に接続されている。I / O システム 2 2 4 を、S M P 2 1 0 と、コンピュータ端末および記憶ディスクのような種々の外部装置および周辺装置とのインターフェースのために設けることができる。

40

【 0 0 2 5 】

プロセッサおよびキャッシュノード 2 1 2 から 2 1 8 は、例えば、上記の図 4 のプロセッサおよびキャッシュノード 4 1 0 として実装できる。あるいは、1 つ以上のプロセッサおよびキャッシュノード 2 1 2 および 2 1 8 は、複数のプロセッサ 4 1 2 およびキャッシュ 4 1 6 を使用できる。どちらの実装においても、S M P 2 1 0 は、多数のプロセッサ 4 1 2 に複数タスクの並列処理を可能にする。集中化共有メモリ 2 2 0 は、マルチプロセッ

50

サ 4 1 2 がタスク間の情報を共有するのを許可する。

【 0 0 2 6 】

図 3 において、分散共有記憶 (D S M) システム 3 1 0 は、相互接続ネットワーク 3 4 4 を介して相互接続されたいくつかの処理ノード 3 5 0 から 3 6 0 を含む。 D S M 3 1 0 は、任意の数の処理ノード 3 5 0 から 3 6 0 を含むことができる。各処理ノード 3 5 0 から 3 6 0 は、プロセッサおよびキャッシュノード 3 1 2 から 3 2 2 および分散共有記憶 3 2 8 から 3 3 8 の一部と共に描かれている。当業者には明らかなように、1 つ以上の処理ノード 3 5 0 から 3 6 0 では、プロセッサおよびキャッシュノードを使用する必要がない。

【 0 0 2 7 】

プロセッサおよびキャッシュノード 3 1 2 から 3 2 2 は、例えば、上記の図 4 のプロセッサおよびキャッシュノード 4 1 0 として実装できる。ここでは、各プロセッサ 4 1 2 は、1 段以上のキャッシュ 4 1 6 を経て、共有メモリ 3 2 8 から 3 3 8 の一部をアクセスする。あるいは、1 つ以上のプロセッサおよびキャッシュノード 3 1 2 から 3 2 2 は、複数のプロセッサ 4 1 2 およびキャッシュ 4 1 6 をもつことができる。

【 0 0 2 8 】

分散共有メモリ部分 3 2 8 から 3 3 8 は、単一の連続ブロックの物理メモリから形成されている場合は、処理ノード 3 5 0 から 3 6 0 内のプロセッサによってアクセスされる。当業者には明らかなように、1 つ以上の処理ノード 3 5 0 から 3 6 0 は、共有メモリの一部を使用する必要がない。

【 0 0 2 9 】

図 3 の例において、各処理ノード 3 5 0 から 3 6 0 は、オプションの入出力 (I / O) 装置と共に表示されている。当業者には明らかなように、1 つ以上の処理ノード 3 5 0 から 3 6 0 は、I / O 装置を有する必要がない。さらに、異なる種類の I / O 装置および外部周辺機器および資源の組合せが D S M システムで使用できる。よって、1 つ以上の処理ノード 3 5 0 から 3 6 0 は、プロセッサまたはプロセッサなし、共有メモリまたは共有メモリなし、および I / O ありまたは I / O なしのどのような組合せをも含むことができる。

【 0 0 3 0 】

図 3 の例において、各処理ノード 3 5 0 から 3 6 0 は、キャッシュコヒーレンシディレクトリと共に表示されている。ディレクトリ情報は、メモリ制御装置および/またはキャッシュタグと関連付けることができる。

【 0 0 3 1 】

D S M 3 1 0 全体で、物理記憶または主記憶 3 2 8 から 3 3 8 を区別することによって、各処理ノード 3 5 0 から 3 6 0 は、主記憶の一部を含むことができる。このプロセッサとメモリとの間の物理的な近接によって、処理ノード内のプロセッサおよびメモリに対する記憶待ち時間が減少される。

【 0 0 3 2 】

図 1 から図 3 と関連して記載したシステム 1 1 0、2 1 0 および 3 1 0 のような単一プロセッサシステム、S M P S および D S M は周知のものである。このようなシステムのさらに詳細については、例えば、パターンソンとヘネシー、コンピュータアーキテクチャの数量的なアプローチ、第 2 版 (H e n n e s s y a n d P a t t e r s o n , C o m p u t e r A r c h i t e c t u r e A Q u a n t i t a t i v e A p p r o a c h , 2 d E d . (M o r g a n a n d K a u f m a n n P u b l . : U S A 1 9 9 6)) に見られ、ここに参照のため、その全体が組込まれている。

【 0 0 3 3 】

(3 . D S M マルチプロセッサシステムにおけるキャッシュコヒーレンスのためのシステムおよび方法)

本発明は、分散共有記憶マルチプロセッサシステムにおいて、いくつかの共有アクセス要求を排他的アクセス要求へ変換することによって、発生し得る共有ベクトルの大きさを

10

20

30

40

50

減少する。これによって、これらの環境において共有ベクトルの必要性を除去する。本発明は、ソフトウェア、ハードウェア、ファームウェアまたはこれらのいずれかの組合せによって実現できる。

【0034】

大きなベクトルの問題は、図3に関連して記述する。処理ノード352が情報を必要とする場合、一般に、要求は相互接続ネットワーク344に渡って送信される。要求された情報が、処理ノード350の主記憶328に存在する場合、処理ノード352からその要求は、処理ノード350へと送られる。その要求が共有アクセス要求で、その要求が認可されたものならば、その主記憶328に関連付けられた共有ベクトルが生成または更新され、処理ノード352に対して当該アクセスがなされたことが反映される。

10

【0035】

その要求が、排他的アクセス要求であれば、共有ベクトルの必要はない。その代わりに、処理ノード352へ排他的アクセスが認可された表示が格納される。共有ベクトルおよび排他的アクセスの表示が、キャッシュコヒーレンシー機構によって使用される。

【0036】

さらに多くの処理ノードが、同一情報に対して、共有アクセスを要求すると、その共有ベクトルは、大きさが増大し、より多くのメモリ空間を必要として、システム全体の処理速度を低下させる。

【0037】

本発明によれば、発生し得る共有エントリの数の減少によって、発生し得る共有ベクトルの大きさは減少する。これは、ローカルノードとリモートノードとを区別することによって行われる。ローカルノードとリモートノードは、要求情報が存在するところの「ホーム」ノードに対比して定義される。上記の例において、処理ノード352が処理ノード350から情報を要求する場合には、処理ノード350はホームノードと呼ばれる。この要求に関して、ホームノード350に対して、リクエストがローカルまたはリモートと定義される。本発明によれば、共有アクセス要求がリモートノードから発生する場合、その共有要求は排他的要求に変換される。その共有アクセス要求がローカルノードから発生する場合、その要求は共有アクセス要求として処理される。

20

【0038】

したがって、ノード350から360は、そのようなノードが情報アクセスを必要とする場合、他のノードのキャッシュ装置へ向かう共有ベクトルを有することがある。図3のシステム310には、6個の処理ノードが表示されている。しかし、当業者にとっては、より小さいシステムおよびより大きいシステムが可能なが理解される。一般に、多数の処理ノードは、潜在的に大きな共有ベクトルとなる。大きな共有ベクトルは、情報処理の速度を減少する。これは、一般に大きな共有ベクトルを扱うには、時間を要するからである。

30

【0039】

本発明の例を、図6を参照して記述する。図6において本発明の実施形態にかかわるDSMシステム600を示す。システム600は、複数の処理ノードつまりリクエスト611および613を有する。処理ノード611はローカルノードまたはローカルリクエストを表す。処理ノード613は、リモートノードまたはリモートリクエストを表す。ノードは、システム内の任意の処理装置、キャッシュメモリ位置、または情報および/または処理機能または手順のいずれかを格納可能な他の任意の媒体とすることができる。

40

【0040】

「ホームノード」という用語は、要求された情報が存在するノードである。「ホームノード」という用語は、処理装置またはメモリ位置または情報を格納可能な、およびDSMマルチプロセッサシステムにおいて、他の処理装置が情報を要求する他の任意の媒体を表現できる。

【0041】

「ホームノード」という用語は、相対的なものであり本発明の範囲を限定するものでは

50

なく、説明のために使用している。したがって、本検討のためだけに、ノード2をホームノード611bに指定する。当業者は、一般にDSMマルチプロセッサシステム内の任意のノードは、同上システムの他の任意のノードおよびホームノードに指定したノードからも情報を要求できることを理解している。

【0042】

ローカルノードは、ホームノード611bに比較的に近接したノードグループ612に位置している。リモートノードは、ノードグループ614に位置していて、概してホームノード611bに比較的に近接していない。

【0043】

図6の例において、ノードグループ612内に*i*個のローカル処理ノード611があり、ノードグループ614内に*n*個のリモート処理ノードがある。システム600の要求に従って、*i*および*n*の数は可変であることを当業者は理解している。

【0044】

各ローカルノード611は、複数の中央処理装置(CPU)615および複数の記憶装置618を有することがある(図6に示すように、各ノード611毎に、4個のCPU615および単一記憶装置618がある)。同様に、各リモートノード613は複数のCPU617および複数の記憶装置619を有することがある(図6に示すように、各ノード613毎に、4個のCPU617および単一の記憶装置619がある)。当業者は、他の実施形態が可能であることを理解している。

【0045】

さて、本発明の使用例を説明する。ローカルノード611aは、共有アクセス要求620をホームノード611bへ送信する。ノード611aは、ホームノード611bに対してローカルノードであるので、その要求620はホームノード611bによって共有アクセス要求として処理される。

【0046】

共有アクセス要求がリモートノードから到着すると、排他的アクセス要求に変換される。例えば、リモートノード613bが共有アクセス要求626をホームノード611bへ送信すると、その共有アクセス要求626は、排他的アクセス要求628に変換される。

【0047】

一度、排他的アクセス要求628が認可されると、ホームノード611bは、リモートノード613bのアドレスへのポインタを格納する。そのアクセスは、排他的であるので、共有ベクトルを有する必要はない。換言すれば、当該情報への共有アクセスを要求する各ノード毎のビットを潜在的に格納しなければならぬ代わりに、当該リモートノード613bの位置を参照するホームノード611bによって、2進数番号(またはノード番号またはCPU番号)が格納される。

【0048】

ある実施例において、ローカルノード611iからの後続の共有アクセス要求622が、リモートノード613bの排他的アクセスを終了する。次に、ローカルノード611iは、当該情報への共有アクセスまたは排他的アクセスを認可される。

【0049】

以下は、本発明の実施例にかかわる方法510の記載である。下記に説明のために、DSM600に関する方法510を示す。しかし、方法510は、DSM600には限定されない。ここにある記述に基づいて、当業者は、方法510が、その他のDSMシステムにおいても実現できることを理解されるであろう。

【0050】

図5aは、本発明の実施例にかかわるキャッシュコヒーレンス維持の方法510を示すフローチャートである。方法510はステップ512から始まり、ここで前述のように共有アクセス要求が、情報を保持するノード、つまりホームノードによって受信される。ある実施例では、ホームノードは、主記憶の一部を含むDSMマルチプロセッサシステムの任意のノードである。ホームノードの主記憶装置は、マルチプロセッサシステム内の種々

10

20

30

40

50

のノードによって要求される情報を含む。このような情報は、当該システム中の他のノードが実行している種々の機能および手順の実行を要求することがある。要求は、当該主記憶の位置から読まれたり書かれたりする。あるノードが、ホームノードの主記憶装置に格納されている特定の情報を使用する必要がある場合、このようなノードは、そのような情報を求めてアクセス要求をホームノードへ送信する。このアクセス要求は、共有アクセス要求または排他的アクセス要求のどちらでもよい。

【0051】

ステップ514において、共有アクセスがローカルノード（ローカルリクエスト）が発生したかまたはリモートノード（リモートリクエスト）発生したかが判断される。ある実施例では、ステップ514は、リクエストのアドレスと要求情報のアドレスとの比較に基づいて行われる。あるいは、ホームノードがシステム内のノードのアドレス表を格納している、このアドレス表ではアドレスをノードの近接パラメータによってソートしていることができる。近接パラメータによって、共有アクセスを要求するノードが、ホームノードに対して、リモートノードであるかローカルノードであるかが判断される。リクエストノードが共有アクセスの要求を送信すると、ホームノードは、要求元ノードのアドレスを読み、要求元ノードのアドレスとアドレス表とを比較する。要求元ノードのアドレスがアドレス表のアドレスと一致すると、要求元ノードは、ローカルノードと判断されて、共有アクセス要求が、共有アクセス要求として処理される。要求元ノードのアドレスが、アドレス表に格納されているアドレスと一致しない場合は、その要求元ノードは、リモートノードと判断されて、その共有アクセス要求は、排他的アクセス要求に変換される。

10

20

【0052】

他の実施例では、要求元ノードのアドレス上位ビットとホームノードのアドレス上位ビットとの比較が行われる。このビットが一致した場合は、要求元ノードローカルノードと判断されて、共有アクセス要求が共有アクセス要求として処理される。一致しない場合は、要求元ノードはリモートノードと判断されて、その共有アクセス要求は排他的アクセス要求に変換される。

【0053】

ステップ514において、共有アクセス要求がローカルノードから発生したと判断されると、処理はステップ516へ進み、ここではその共有アクセス要求は共有アクセス要求として処理される。共有アクセス要求が認可された場合は、大体、共有ベクトルが生成され、および/または更新されてその要求元ローカルノードは要求情報への共有アクセスを有することが識別される。

30

【0054】

ステップ514では、共有アクセス要求がリモートノードから発生したと判断されると、処理はステップ518へ進み、ここではその共有アクセス要求は排他的アクセス要求に変換される。これを図6に示す。ここでは、共有アクセス要求626が排他的アクセス要求628に変換されている。次に、排他的アクセス要求は、排他的アクセス要求手順、例えばここで記述したようにステップ520および/またはステップ524などに従って処理される。

【0055】

例えば、その要求情報が他のリクエストの間で共有されている場合、および当該リモートノードに排他的アクセスが与えられている場合には、ステップ520において、当該要求情報のコピーを以前にキャッシュした他のノードへ、ホームノードは無効命令を送信する。これによって、そのリモートノードがその情報への排他的アクセスを有している間は、他のプロセスによって、その情報のキャッシュされたコピーが使用されるのが防止される。

40

【0056】

そのリモートノードは、排他的アクセスを認可されているので共有ベクトルは必要ない。その代わりにステップ524において、ポインタが格納されて要求情報への排他的アクセスが認可されたリモートノードを識別する。このポインタは、共有ベクトルの代わりに

50

キャッシュコヒーレンス機構で使用される。このようにして、本発明はローカルノードの数に基づいた共有ベクトルの発生しうる大きさを制限する。換言すれば、ローカルノードと共にリモートノードを潜在的に識別する大きな共有ベクトルを有するという問題が、大幅に解消される。

【0057】

図5bにおいて、本発明の別の実施例の方法510を表示している。この実施例では、リモートノードが、ホームノードによって排他的アクセス要求を認可された場合でも、図5aのステップ512から524に記載されているように、ローカルノードはリモートノードの排他的アクセス権を終了することができる。これを図6に示す。ここでは、リモートノード613bは排他的アクセスが認可され、続いてローカルノード611iが共有アクセスを要求する。

10

【0058】

図5bにおいて、ステップ530で方法510は、アクセス要求を受信する。次に、方法510はステップ532に進む。ここで、受信したアクセス要求が排他的要求かどうか判定処理される。ステップ532において、アクセス要求が排他的アクセス要求でなければ（すなわち、共有アクセス要求であれば）、ステップ534において、その共有アクセス要求がリモートノードから発生したかどうか判定処理される。ステップ534は、前述したように実行することができる。例えば、アドレス比較を行うことができる。

【0059】

要求が排他的要求（ステップ532で判断される）またはリモートノードから発生する場合は（ステップ534で判断される）、本方法では、処理がステップ536に進む。ここでは、要求情報が他のリクエストによって排他的に保持されているかどうか判定処理される。ステップ536において、要求情報が現時点で他のリクエストによって排他的に保持されている場合は、ステップ548に示すように、介入要求が当該情報の排他的所有者に送信される。

20

【0060】

ステップ548の機能は、現在の排他的所有者が情報を修正した場合には、新しいリクエストが排他的アクセスを認可される前に、修正情報がホームノードへ返却されることを保証する。一度、要求情報の排他的所有者によって介入要求が送信されると、処理はステップ550に進む。

30

【0061】

ステップ550において、新しいリクエスト（例えば、リモートノード）は、要求情報への排他的アクセスを認可される。例えば図6において、ノード613bは、ノード611bに位置している要求情報への排他的アクセスを有している新しいリクエストであってもよい。

【0062】

新しいリクエストが、現在、要求情報への排他的アクセスを有することを表示するために、ステップ552において、ホームノードによって、当該リモートリクエストに排他的アクセスが与えられたことを示すポインタが格納される。

【0063】

ステップ536に戻って、当該要求情報は他のリクエストによって現在、排他的に保持されていないとシステムが判断すると、ステップ554に示すように、本方法では同一情報への共有アクセスを現在保持している全ノードへ無効化命令を送信する。次に、処理は前述のようにステップ550および552進む。

40

【0064】

ステップ532に戻って、受信したアクセス要求が排他的要求ではない場合は、処理はステップ534へ進む。また、ステップ534に戻って、受信したアクセス要求がリモートノードから発生したのではないと判定処理された場合は、処理はステップ538へ進む。

【0065】

50

ステップ538において(ステップ536と同様)、当該要求情報が、現在、他のリクエストによって排他的に保持されているかどうかを判定処理する。ステップ538において、要求情報が現在、他のリクエストによって排他的に保持されていない場合には、処理はステップ544へ進む。

【0066】

ステップ544はステップ550と同様であり、ここではリクエストは、要求情報へのアクセスを認可される。しかし、ステップ544では、リクエストはローカルノード(ステップ534で判断される)であり、したがって、要求情報へ共有アクセスが与えられる。

【0067】

ステップ544に続くステップ546において、新しい共有ベクトルは、要求情報が位置しているノードによって格納される。例えば図6において、ローカルノード611iが、ノード611bに位置している要求情報への共有アクセスを要求すると、ノード611bは新しい共有ベクトルを格納する。

【0068】

ステップ538に戻って、要求情報が現在、他のリクエストによって排他的に保持されていると判定処理されると、処理はステップ540に進む。ステップ540において、リクエストの位置に対して、ホームノードがリモートに位置付けられているかが判定処理される。要求情報の現在の排他的所有者が、ホームノードに対して、リモートであると判断されると、処理はステップ548へ進む。ステップ548およびこれに続くステップは上記のとおりである。

【0069】

ステップ540において、ホームノードに対してリクエストがローカルであると判定処理されると、処理はステップ542に進む。ステップ542(ステップ548と同様)において実行される機能によって、現在の排他的所有者が情報を修正した場合は新しいリクエストが排他的アクセスを認可される前に当該修正情報がホームノードへ返却されることが保証される。ステップ542において、ローカルリクエストは当初のリクエストの排他的情報アクセス終了する。当初のリクエストの情報要求は、共有アクセス要求となる。

【0070】

一度、介入要求が、要求情報の排他的所有者へ送信されると、処理はステップ544に進む。ステップ544および続くステップ546については前述されている。ステップ544において、情報アクセスはローカルリクエストに与えられ、ステップ546において共有ベクトルがシステムによって格納される。

【0071】

リモートノードが、さらに当該情報へアクセスする必要がある場合は、再度、ステップ530を実行するなどによって、その情報を再要求せねばならない。あるいは、ローカルリクエストと共に当該情報へ共有アクセスできるように、リモートノードの排他的アクセスを共有アクセスに置き換える。

【0072】

上述の方法510は、ディレクトリに基づくプロトコルに実装することができるが、他の実装も可能である。

【0073】

(4. 結論)

本発明の方法、システム、構成要素の実施事例を、ここに記載した。他の所で言及したように、これらの実施事例は説明だけの目的で記載されたもので、限定するものではない。その他の実施例も可能であり、本発明によって包含される。このような実施例は、ここに含まれている教示に基づいて、関連技術の当業者にとっては明白である。よって、本発明の広さと範囲は上記の実施事例によって限定されるべきものではなく、以下の請求項およびその等価物に従ってのみ定義されるべきである。

【図面の簡単な説明】

10

20

30

40

50

【 0 0 7 4 】

本発明は、添付の図面を参照して記載する。各図面において、同一の参照番号は、同一の要素または類似機能の要素を示す。さらに、参照番号の左端の数字は、その参照番号が最初に現れる図面を示す。

【 図 1 】 図 1 は、単一プロセッサシステムのブロック図である。

【 図 2 】 図 2 は、集中対称型共有記憶マルチプロセッサ (S M P) システムのブロック図である。

【 図 3 】 図 3 は、分散共有記憶マルチプロセッサシステム (D S M) のブロック図である。

【 図 4 】 図 4 は、図 1 から図 3 に示すプロセッサシステムのいずれでも使用可能なプロセッサとキャッシュのブロック図である。

【 図 5 a 】 図 5 a は、D S M マルチプロセッサシステムにおけるキャッシュコヒーレンスを保持するための方法の実施例のフローチャートである。

【 図 5 b 】 図 5 b は、D S M マルチプロセッサシステムにおけるキャッシュコヒーレンスを保持するための方法の別の実施例のフローチャートである。

【 図 6 】 図 6 は、本発明を適用できる D S M マルチプロセッサシステムの実施例のブロック図である。

10

【 図 1 】

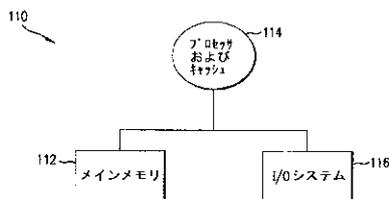


FIG. 1

【 図 2 】

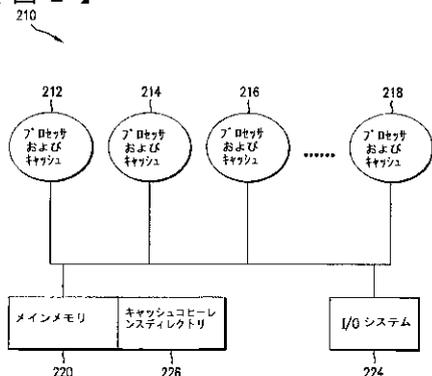


FIG. 2

【 図 3 】

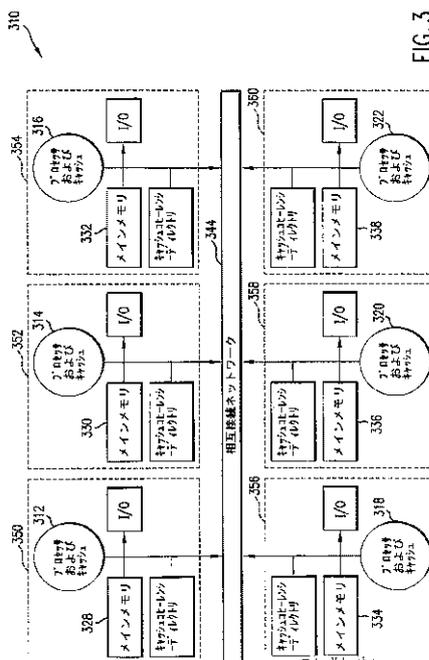


FIG. 3

【図4】

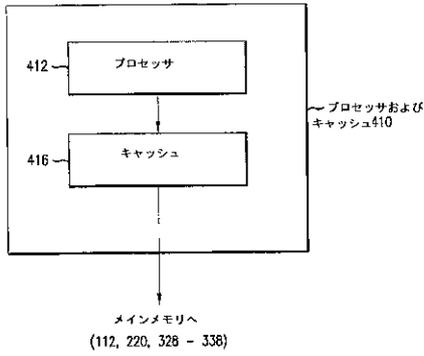


FIG. 4

【図5a】

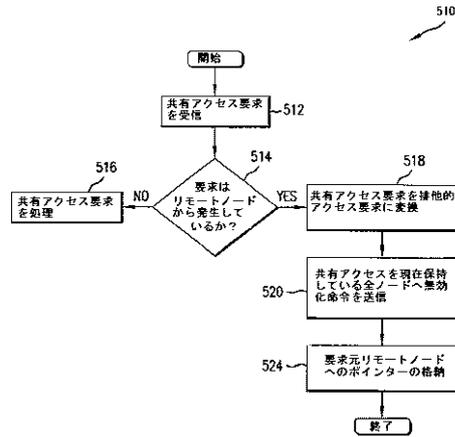


FIG. 5a

【図5b】

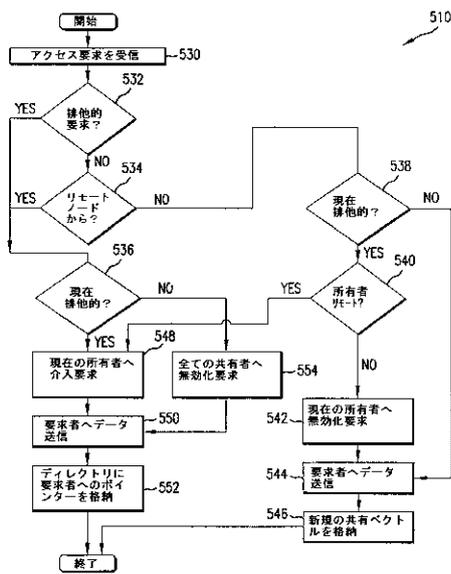


FIG. 5b

【図6】

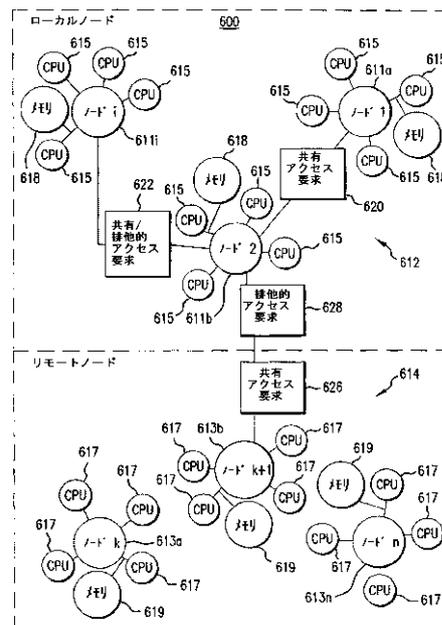
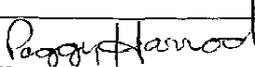


FIG. 6

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT		International application No. PCT/US02/05779																								
A. CLASSIFICATION OF SUBJECT MATTER IPC(7) : G06F 12/00 US CL : 711/147, 141, 148, 152, 153, 154 According to International Patent Classification (IPC) or to both national classification and IPC																										
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) U.S. : 711/147, 141, 148, 152, 153, 154 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) WEST																										
C. DOCUMENTS CONSIDERED TO BE RELEVANT																										
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.																								
Y,P A A A A	US 6,370,622 B1 (CHIOU et al.) 9 APRIL 2002, entire US 6,026,474 A (CARTER et al.) 15 February 2000, entire US 5,940,860 A (HAGERSTEN et al.) 17 AUGUST 1999, entire US 5,617,577 A (YAMADA et al.) 1, APRIL 1997, entire US 5,592,625 A (SANDBERG) 7 JANUARY 1997, entire	1-41 1-41 1-41 1-41 1-41																								
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.																										
<table border="0"> <tr> <td colspan="2">* Special categories of cited documents:</td> <td>"T"</td> <td>later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</td> </tr> <tr> <td>"A"</td> <td>document defining the general state of the art which is not considered to be of particular relevance</td> <td>"X"</td> <td>document of particular relevance, the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone</td> </tr> <tr> <td>"E"</td> <td>earlier document published on or after the international filing date</td> <td>"Y"</td> <td>document of particular relevance, the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art</td> </tr> <tr> <td>"L"</td> <td>document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)</td> <td>"&"</td> <td>document member of the same patent family</td> </tr> <tr> <td>"O"</td> <td>document referring to an oral disclosure, use, exhibition or other means</td> <td></td> <td></td> </tr> <tr> <td>"P"</td> <td>document published prior to the international filing date but later than the priority date claimed</td> <td></td> <td></td> </tr> </table>			* Special categories of cited documents:		"T"	later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention	"A"	document defining the general state of the art which is not considered to be of particular relevance	"X"	document of particular relevance, the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone	"E"	earlier document published on or after the international filing date	"Y"	document of particular relevance, the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art	"L"	document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&"	document member of the same patent family	"O"	document referring to an oral disclosure, use, exhibition or other means			"P"	document published prior to the international filing date but later than the priority date claimed		
* Special categories of cited documents:		"T"	later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention																							
"A"	document defining the general state of the art which is not considered to be of particular relevance	"X"	document of particular relevance, the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone																							
"E"	earlier document published on or after the international filing date	"Y"	document of particular relevance, the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art																							
"L"	document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"&"	document member of the same patent family																							
"O"	document referring to an oral disclosure, use, exhibition or other means																									
"P"	document published prior to the international filing date but later than the priority date claimed																									
Date of the actual completion of the international search 03 JUNE 2002		Date of mailing of the international search report 25 JUN 2002																								
Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box FCT Washington, D.C. 20231 Facsimile No. (703) 305-3230		Authorized officer  KIMBERLY MCLEAN Telephone No. (703) 308-8900																								

フロントページの続き

(72)発明者 デネロフ, マーティン エム.

アメリカ合衆国 ニュージャージー 07755, オークハースト, オールド ファーム ロ
ード 24

Fターム(参考) 5B005 KK02 MM01 PP11

5B060 KA01 KA02 KA06