

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3849951号
(P3849951)

(45) 発行日 平成18年11月22日(2006.11.22)

(24) 登録日 平成18年9月8日(2006.9.8)

(51) Int. Cl. F I
G06F 12/08 (2006.01) G06F 12/08 531B
 G06F 12/08 531E

請求項の数 3 (全 24 頁)

<p>(21) 出願番号 特願平9-59914 (22) 出願日 平成9年2月27日(1997.2.27) (65) 公開番号 特開平10-240707 (43) 公開日 平成10年9月11日(1998.9.11) 審査請求日 平成16年2月18日(2004.2.18)</p>	<p>(73) 特許権者 000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号 (74) 代理人 100099298 弁理士 伊藤 修 (72) 発明者 垂井 俊明 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内 (72) 発明者 岡澤 宏一 神奈川県川崎市幸区鹿島田890番地 株 会社日立製作所 情報・通信開発本部内 (72) 発明者 岡田 康行 神奈川県海老名市下今泉810番地 株 会社日立製作所 オフィスシステム事業部 内 最終頁に続く</p>
---	--

(54) 【発明の名称】 主記憶共有型マルチプロセッサ

(57) 【特許請求の範囲】

【請求項1】

1つ以上のCPU、キャッシュと、主記憶を備える複数のノードと、ノード間を結ぶネットワークとからなり、ネットワークを用いてノード間でキャッシュコヒーレントチェックのためのコマンドとその応答のやり取りを行ってキャッシュコヒーレント制御を行う主記憶共有型マルチプロセッサにおいて、

前記各ノードは、自ノードの主記憶の各ページ対応に、該当するページが他のノードからアクセスされたかどうかを記憶する第1のビットが1ビット割り当てられ、該第1のビットはシステムの初期化時にリセットされ、主記憶の該当するページが他のノードからアクセスされた場合に、ハードウェアによりセットされるテーブルと、

自ノードのCPUが自ノードの主記憶をアクセスする際に、アクセスするページに該当する前記テーブルの前記第1のビットを検査し、該第1のビットがセットされていた場合には他のノードへのキャッシュコヒーレントチェックのためのコマンドの送出を行い、該第1のビットがセットされていなかった場合には他のノードへのキャッシュコヒーレントチェックのためのコマンドの送出を行わない手段を備えることを特徴とする主記憶共有型マルチプロセッサ。

【請求項2】

請求項1記載の主記憶共有型マルチプロセッサにおいて、

システムソフトウェアが主記憶のページをアロケートする際に、アロケートするページに対応する前記テーブルのビットを該システムソフトウェアがリセットすることを特徴と

する主記憶共有型マルチプロセッサ。

【請求項3】

請求項1記載の主記憶共有型マルチプロセッサにおいて、

前記テーブルに、前記主記憶の各ページ対応に該ページに対してキャッシュコヒーレント制御が必要で無いことを記憶する第2のビットを1ビット割り当て、

前記手段は、自ノードのCPUが自ノードの主記憶をアクセスする際に該第2のビットを検査し、該ビットがセットされていない場合には、前記第1のビットの値にしたがって他のノードへのキャッシュコヒーレント制御の要否を判断し、前記第2のビットがセットされていた場合には他のノードへのキャッシュコヒーレントチェックのためのコマンドの送出不行を特徴とする主記憶共有型マルチプロセッサ。

10

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は情報処理装置、特に、パーソナルコンピュータ(PC)、ワークステーション(WS)、サーバ機等に用いられる主記憶共有型の並列計算機システムに係り、特に、主記憶の制御方式に関する。

【0002】

【従来の技術】

近年PC、WSの上位機種及びサーバ機などでは、主記憶共有型のマルチプロセッサ(SMP)構成が広まっており、性能を向上させるために、20~30台以上の多数台のプロセッサの間で主記憶を共有することが重要な課題になってきている。

20

主記憶共有型のマルチプロセッサの構成方法として広く使われている方式として共有バスがあるが、バスではバスのスループットがネックになるため、接続可能なプロセッサの数は高々8台程度である。従って、多数台のプロセッサを接続する方式としては適さない。

【0003】

現在行われている、多数台のプロセッサを接続する主記憶共有マルチプロセッサの構成方法には、大きく2つの方式がある。

その一つに、クロスバススイッチによる構成があり、例えば、「進化したシステム・アーキテクチャ」(SunWorld誌1996年1月号第29頁~第32頁)に開示されている。

30

この方式では、プロセッサと主記憶を持つ各ボードを、高速なクロスバススイッチで接続し、プロセッサ間のキャッシュ一貫性を保持する。

この方式では、キャッシュ一貫性の保持が高速にできるという利点がある。

しかし、キャッシュの一貫性を保持するためのトランザクションが全プロセッサにブロードキャストされるため、クロスバススイッチにかかるトラフィックが非常に高く、性能的に隘路になるとともに、高速なスイッチが必要になるためコスト高を招くという欠点がある。

さらに、キャッシュ一貫性保持のためのトランザクションをブロードキャストしなければならないため、台数が非常に多いシステムを実現することは困難であり、数十台程度が限度である。

40

以下ではこの方式をスイッチ型SMP(Symmetrical Multi Processor)と呼ぶ。

【0004】

それに対して、ディレクトリ方式によるマルチプロセッサの構成があり、たとえば「The Stanford FLASH Multiprocessor」(第21回ISCA Proceedings)に開示されている。

この方式では、主記憶のデータライン毎に、そのデータラインがどこのプロセッサのキャッシュに接続されているかを示すビットマップであるディレクトリを設けることにより、必要なプロセッサにのみキャッシュ一貫性を保持するためのトランザクションを送る。

それにより、スイッチにかかるトラフィックを大幅に削減することができ、スイッチの八

50

ードウェアコストを削減することができる。

しかし、キャッシュ一貫性を保持するためのトランザクションを出す際には、必ず、主記憶に置かれたディレクトリの内容をチェックしなければならないため、アクセスレーテンシが大幅に増えるという欠点がある。

さらに、ディレクトリを置くためのメモリのコストが高くなるという欠点を持つ。

【0005】

上記のようにスイッチ型SMPとディレクトリ方式は一長一短である。

一般にスイッチ型SMPの方がハードウェア規模は大きくなり、台数が多くなった場合のスケラビリティは良くないが、高性能を達成できる。

したがって、PCやサーバ機等の、台数のそれほど多くない(30台程度までの)システムでは、スイッチ型SMPで実現する方が得策である。

10

【0006】

ここで、主記憶共有型マルチプロセッサを構成する上でのもう一つの問題点として、信頼性の問題がある。

従来の主記憶共有型マルチプロセッサは、システム全体で一つのOSを持つ。この方式は、システムの全てのプロセッサを一つのOSで管理できるため、柔軟なシステム運用(負荷分散等)をできるという利点を持つ。

しかし、多数台のプロセッサを主記憶共有のマルチプロセッサ構成で接続した場合、システムの信頼性が低下するという欠点を持つ。

複数のプロセッサをネットワークで接続したクラスタ構成のサーバや、MPP(Massively Parallel Processors)では、ノード毎にOSは別なので、OSなどのバグ等のためにシステムダウンしても、システムダウンするのは当該のノードのみである。

20

それに対して主記憶共有型のマルチプロセッサで、システム全体を1つのOSで制御する場合、あるプロセッサがシステムバグ等のためにダウンすると、OSがダウンしてしまうため、全てのプロセッサが影響を受けてしまう。

【0007】

上記の問題を避けるために、主記憶共有型のマルチプロセッサにおいて、複数のOSを走らせる方式が「Hive: Fault Containment for Shared-Memory Multiprocessors」(第15回 ACM Symposium on Operating Systems Principles)において開示されている。

30

この方式では、ディレクトリ方式の主記憶共有のマルチプロセッサに、以下の2つの機構を持たせる。

(1)システム全体を複数のセル(パーティション)にわけ、各パーティション毎に、独立したOSを走らせる。全体のアドレス空間は一つであり、OSごとに担当するアドレス範囲が異なる。

(2)主記憶のページ毎に書き込みアクセス可能なプロセッサを表すビットマップを設け、書き込みアクセスをビットマップが1であるプロセッサのみに許す。

各プロセッサの主記憶に書き込みが行われる場合(Fetch & Invalidate命令によりデータがキャッシングされる場合、もしくはWrite Back要求が到来した場合)、ビットマップの内容がチェックされ、ビットマップが1であるプロセッサからのアクセスのみが許される。

40

上記の(1)の機構により、たとえ、1つのパーティションのOSがダウンしても、他のパーティションがダウンする事を避けることができる。

さらに、(2)の機構を設けることにより、バグにより暴走したパーティションのプロセッサが、他のパーティションが使用するデータを破壊することを防止することができる。以上のように、主記憶共有型のマルチプロセッサ内を複数のパーティションに分けることにより、システムの信頼性を大幅に向上させることができる。

【0008】

50

【発明が解決しようとする課題】

上記従来技術で述べた、スイッチ型SMPを構成し、さらにSMP内をパーティションに分けようとする場合、以下に示す3つの問題点がある。

【0009】

(A) ローカル主記憶のアクセスが遅い

プロセッサが同一ボード内の主記憶をアクセスする場合、理想的にはクロスバススイッチを経由せずに高速にアクセスを行うことができるはずである。

しかし、実際は、必ず他のプロセッサへのキャッシュ一貫性保持のためのトランザクションを出し、他のプロセッサのキャッシュのチェック(以下ではこの処理をCCC: Cache Coherent Checkと呼ぶ)を行わなければならない。

10

なぜなら、他のプロセッサのキャッシュの上に、アクセスされたデータのコピーがキャッシングされている可能性があるからである。

実際に他のプロセッサのキャッシュにデータがキャッシングされていた場合は、上記のCCCは無駄にはならない。

しかし、アクセスされたデータが、他のプロセッサから一回もアクセスされていないローカルなデータの場合は、該当するデータが他のプロセッサのキャッシュ上にキャッシングされている可能性は0であり、本来は、CCCの処理は全く不要である。

そのため、無駄なCCCのために、アクセスレーテンシを増加させてしまうのみならず、スイッチ上のトラフィックを増大させてしまうという問題点がある。

【0010】

20

ディレクトリ方式では、キャッシュライン単位でどのプロセッサにキャッシングされているかが管理されているため、無駄なCCCは生じない。

しかし、先に述べたように、ディレクトリのためのハードウェア量が多いのみならず、ディレクトリを管理するためのオーバヘッドが非常に大きいという欠点を持つ。

たとえば、16プロセッサ、4GB主記憶、64B/ラインのシステムのディレクトリとしては、

$$4GB / 64B * 16bit = 128MB$$

もの主記憶が必要になる。

従って、大幅なハードウェア量の削減が必要になる。

【0011】

30

(B) パーティションのアドレスが0番地から始まらない

上記の従来のパーティション管理機構では、システム全体でアドレス空間が一つである。

従って、各パーティションの分担するアドレスが0番地から始まらない。

例えば、パーティション数が2、各パーティションの主記憶容量が1MBとすると、パーティション0は0番地から1M番地までのアドレス空間を持つのに対し、パーティション1は1M番地から2M番地までのアドレス空間を持たなければならないことになる。

既存のOSは、アドレスが0番地から主記憶が実装されていることを前提としているため、上記の制限は従来OSを使用する場合、大きな障害になる。

【0012】

(C) パーティション管理のためのハードウェア量が多い

40

上記従来例のパーティション管理機構を用いる場合、4KBのページ毎に、各プロセッサが該当するページへのアクセスを許されるかどうかビットマップで記憶されている。

したがって、該当するビットマップのハードウェア量が非常に大きいという問題点がある。

例えば、プロセッサの台数を16台、システムの主記憶容量を4GBとすると、

$$4GB / 4KB * 16 = 16MB$$

ものメモリがパーティション管理のために必要になり、コストの増大を招く。

【0013】

従って、

本発明の第一の目的は、他のプロセッサから全くアクセスされていないローカルなデータ

50

を、他のノードへのCCCを行わずに高速にアクセスすることが可能な、主記憶共有型のマルチプロセッサを、少ないハードウェアオーバーヘッドで実現することである。

本発明のもう一つの目的は、主記憶共有型のマルチプロセッサをパーティションに分けた際に、各パーティションの持つローカルな主記憶が独立したアドレス空間を持つことによりローカルな主記憶を0番地から始めることができ、かつ、必要な領域を共有することができる主記憶共有型のマルチプロセッサを構成することである。

本発明のさらなる目的は、上記のパーティション管理を少ないハードウェア量で実現することである。

【0014】

【課題を解決するための手段】

上記目的を達成するため、本発明は、

1つ以上のCPU、キャッシュと、主記憶を備える複数のノードと、ノード間を結ぶネットワークとからなり、ネットワークを用いてノード間でキャッシュコヒーレント制御を行う主記憶共有型マルチプロセッサにおいて、

各ノードは、自ノードの主記憶の各ページ対応に、該当するページが他のノードからアクセスされたかどうかを記憶する第1のビットが1ビット割り当てられ、該第1のビットはシステムの初期化時にリセットされ、主記憶の該当するページが他のノードからアクセスされた場合に、ハードウェアによりセットされるテーブルと、自ノードのCPUが自ノードの主記憶をアクセスする際に、アクセスするページに該当する前記テーブルの前記第1のビットを検査し、該第1のビットがセットされていた場合には他のノードへのキャッシュコヒーレント制御を行い、該第1のビットがセットされていなかった場合には他のノードへのキャッシュコヒーレント制御を行わない手段を備えるようにしている。

【0015】

さらに、システムソフトウェアが主記憶のページをアロケートする際に、アロケートするページに対応する前記テーブルのビットを該システムソフトウェアがリセットするようにしている。

【0016】

また、前記テーブルに、前記主記憶の各ページ対応に該ページに対してキャッシュコヒーレント制御が必要で無いことを記憶する第2のビットを1ビット割り当て、

前記手段は、自ノードのCPUが自ノードの主記憶をアクセスする際に該第2のビットを検査し、該ビットがセットされていない場合には、前記第1のビットの値にしたがって他のノードへのキャッシュコヒーレント制御の要否を判断し、前記第2のビットがセットされていた場合には他のノードへのキャッシュコヒーレント制御を行わないようにしている。

【0017】

1つ以上のCPU、キャッシュと、主記憶を備える複数のノードと、ノード間を結ぶネットワークとからなり、ネットワークを使ってノード間でキャッシュコヒーレント制御を行い、主記憶を共有している各ノードを1つ以上のノードからなる複数のパーティションに分けることが可能な主記憶共有型マルチプロセッサにおいて、

各ノードの主記憶を、全ノードからアクセス可能な共有領域と、パーティション内からのみアクセス可能なローカル領域に分割し、各々の領域について別個の開始アドレスを指定するようにしている。

【0018】

さらに、前記各ノードは、アクセスされたアドレスがローカル領域であるか共有領域であるかを判定する手段と、パーティション内にどのノードが含まれているかを判定する手段を備え、

他のノードへキャッシュコヒーレンス制御のためのコマンドを出すとき、共有領域へのアクセスコマンドに関しては、システム内の全ノードにコマンドをブロードキャストし、ローカル領域へのアクセスコマンドに関しては、パーティション内のノードにのみコマンドをマルチキャストするようにしている。

10

20

30

40

50

【 0 0 1 9 】

また、前記各パーティションのローカル領域のアドレスが0番地から始まるようにしている。

また、他のノードからキャッシュコヒーレントコマンドが到来した際に、アクセスアドレスがローカル領域か共有領域かを判定する手段と、アクセス元のノードがパーティション内かパーティション外かを判定する手段を備え、

ローカル領域にパーティション外のノードからコマンドが到来したと判定された場合には、アクセスを抑止し、エラーを報告するようにしている。

また、前記各ノードは前記共有領域の構成情報を記憶するレジスタを備えるようにしている。

10

また、前記共有領域の構成情報は、共有領域の開始アドレスと、1プロセッサの担当する共有領域の大きさからなるようにしている。

また、前記共有領域の構成情報は、共有領域の開始アドレスと終了アドレスの組からなるようにしている。

また、前記各ノードは、パーティションの構成情報を記憶する手段として、パーティション内のノードの分布をビットマップで記憶する手段を備えるようにしている。

【 0 0 2 0 】

【 発明の実施の形態 】

以下、本発明に係る主記憶共有型マルチプロセッサを、図面に示した実施の形態を参照してさらに詳細に説明する。

20

【 0 0 2 1 】

(1) 装置の概要

図1は本発明に係る主記憶共有型マルチプロセッサのブロック図である。

本システムは64ノードのシステムであり、複数のノード、例えば100、200（これらはノード0、ノード63と呼ぶことがある）が、ネットワーク900により接続される。各ノードは同じ構造を有する。

すなわち、各ノードは、CPU110～112、（部分）主記憶160、主記憶アクセス回路130、ネットワークコマンド送信/受信回路180/190を持つ。180、190、900については、公知の技術であるので内部の説明は省略する。

主記憶160は、このシステムに共通の主記憶の一部を構成し、各ノードで実行されるプログラムおよびデータの一部を保持するもので、このシステムはいわゆる分散共有メモリ型の並列計算機システムである。

30

主記憶は、ローカル主記憶161、共有主記憶162に分けられ、ベースアドレスレジスタ1610、1620により、それぞれ別個のアドレスを指定できる。

ベースアドレスレジスタ1610、1620はシステムの立ち上げ時に、システムのアドレスマップに応じて、後に述べるパーティション/主記憶構成情報レジスタ150とともにセットされる。

図ではノード内のCPUはバスにより接続されているが、バス以外の結合方式、例えば一対一結合、スイッチによる結合でもよい。

これらノード内の接続方法については公知の技術であるので内部構造の詳細な説明は行わない。

40

【 0 0 2 2 】

主記憶アクセス回路130は、CPUから主記憶アクセス命令が発行されたときに、他のノードへのCCCのためのネットワークコマンドの発行、他ノードへの主記憶アクセスコマンドの発行、自ノード内の主記憶のアクセス、および、他のノードから送られてきたCCCコマンド、主記憶アクセスコマンドを実行するための回路である。

まず、送信側の回路を説明すると、131はCPUから送られてきたバスコマンドの受信/バスコマンドの分類を行うための回路である。

132はCPUからアクセスされたアドレスが内部（ノード内の主記憶のアドレス）かリモート（他のノードが持つ主記憶のアドレス）かを判断するための回路である。

50

132はパーティション構成情報150の内容を用いて内部ノリモートを判断する。

132、150は本実施の形態に特有の回路である。

【0023】

138は自ノード内の主記憶の各ページの属性(他のノードからアクセスされたか否か、及び他のノードへのCCCが不要であるかどうか)を、記憶するためのテーブルRAT(Remote Access Table)、133はCPUからアクセスされたアドレスのRATの値をチェックし、必要な動作を起動するための回路、148は他のノードからアクセスされたページのRATの値を変更するための回路である。

これらの138、133、148は本実施の形態に特有の回路である。

【0024】

137はRATを初期化等のためにCPUからアクセスするための回路である。

134はCCCコマンド、他ノード主記憶アクセスコマンド、他ノードへの返答コマンド等のネットワークコマンドを生成するための回路である。

139はアクセスアドレス、アクセスコマンドより、ネットワークコマンドをどのノードに出すかを判断する回路である。

139は本実施の形態に特有の動作を行う。

135はCPUが自ノードの主記憶上のデータをアクセスする際に、他ノードへ出したCCCの内容を記憶し、他のノードから帰ってきたCCCに対する返答コマンドを集計するためのCCC待ち合わせ回路A、136はCPUがアクセスした自ノードの主記憶160の内容をアクセスするための主記憶アクセス回路Aである。

【0025】

次に受信側の回路を説明すると、141はネットワークコマンドのアドレス、送信元ノード番号と、パーティション構成情報150の内容の整合性をチェックするための回路であり、本実施の形態に特有の回路である。

142は他のノードから送られてきたネットワークコマンドの分類を行う回路、147は他のノードから送られてきたコマンドのアクセスアドレスが内部リモートかを判断する回路である。147は、本実施の形態に特有の回路である。143は他のノードから送られてきたCCCコマンドなどを、ノード内のバスに出力するための回路である。

144は、他のノードへのデータフェッチ要求の回答を集計し、CPUに返送するデータを選択するための回路である。

145は、他のノードが自ノード内の主記憶をアクセスした際に、自ノードのバスに出したCCCコマンドの内容を記憶し、自ノードのバスからの返答を待ち合わせるためのCCC待ち回路B、146は、他のノードからの要求に応じて自ノード内の主記憶160をアクセスするための回路である。

149は他のノードからのアクセスコマンドが来た際に、アクセス元のノード番号を一時的に記憶するためのラッチである。返答先のノード番号を知るために用いられる。

【0026】

(2)バス、ネットワークコマンドの説明

ノード内のバス上では、以下の6つのコマンドが使われる。括弧内はこの実施例で使われる略号である。

・Fetch(F)

データのライン転送を要求する。

CPUの読み出しコマンドがミスした場合に出される。CCCコマンドの一つである。

・Fetch&Invalidate(FI)

データのライン転送と同時に、他のキャッシュ上のデータの無効化を要求する。CPUの書き込みコマンドがミスした場合に出される。CCCコマンドの一つである。

・Invalidate(I)

他のキャッシュ上のデータの無効化を要求する。

CPUが、他のキャッシュと共有されているキャッシュラインに対して書き込み要求を出した場合に出される。CCCコマンドの一つである。

10

20

30

40

50

・ Write Back (WB)

キャッシュラインの書き戻しを要求する。リプレースにより、データが追い出されたときに生じる。

・ Data (D)

データ転送を要求する。

F, FI コマンドへの返答。

・ No Data (ND)

F, FI コマンドに対し、どのCPUもDコマンドを出さない状態。

これは、バス上では明示的なコマンドとしては存在しないが、便宜的にコマンドとして扱う。

バス上では、コマンドに付随して、アドレスが転送され、さらにWB, Dコマンドではデータ(キャッシュライン)が転送される。

【0027】

ネットワークコマンドは以下の7種類が存在する。

バスコマンドと同じものについては意味の説明を略する。

・ Fetch (F)

・ Fetch & Invalidate (FI)

・ Invalidate (I)

・ Write Back (WB)

・ Data (D)

F, FI コマンドに対し、キャッシュ上のデータを返送するためのコマンド。

・ Data Mem (DM)

F, FI コマンドに対し、主記憶上のデータを返送するためのコマンド。

ここで、他のいずれかのノードのキャッシュよりDコマンドが来た場合は、DMコマンドで返送されたデータは無視されなければならない。

・ No Data (ND)

F, FI コマンドに対し、該当するノードではどのCPUもDコマンドを出さないことを示す返答。

【0028】

図6、図7にバスコマンドのフォーマットを示す。

図6はF, FI, I, NDコマンドであり、コマンドの他に、宛先ノード番号を表すビットマップ、アクセス元ノード番号、アクセスアドレスを含む。

ここで、宛先ノード番号はビットマップで表されているため、複数のビットを立てることにより、特定の複数のノード(例えばパーティション内のノード)に向けたマルチキャスト、さらには、全ビットを立てることにより、システムの全プロセッサに向けたブロードキャストを容易に実現することができる。

図7はWB, D, DMコマンドのフォーマットであり、上記の情報の他にキャッシュラインのデータを含む。

【0029】

(3) CPUのからのアクセスに対する動作

以下では、CPUからのアクセスに対する主記憶をアクセス回路130の動作を場合に分けて順に説明する。

共有バス上にコマンドが出されると、バスコマンド受信/バスコマンド分類回路131は、アクセスされたアドレスをリモート判定回路132に送り、アクセスされたコマンドが自ノードの主記憶のアドレスをアクセスしているか(内部)、他ノードの主記憶のアドレスをアクセスしているか(リモート)を判断する。

リモート判定回路132では、パーティション/主記憶構成情報150の内容を使用し、内部/リモートの判定を行う。

【0030】

図5にパーティション/主記憶構成情報150の内容を示す。

10

20

30

40

50

パーティションノ主記憶構成情報150は、パーティション内にどのノードが属しているかを示すビットマップ1500、パーティション内のノードの数を表すレジスタ1506を持ち、さらに、共有領域に関しては、共有領域の先頭アドレス(ベースアドレス)1501、1ノードあたりの共有領域の大きさ1502を記憶する。

ここでは、各ノードの共有領域は、全て同じ大きさである。

さらに、ローカル領域に関しては、パーティション内の各ノードに対して、ローカル領域の開始アドレス1503と終了アドレス1504を記憶する。

パーティションレジスタ1500の該当するビットが1であるノードに対してのみこの情報は有効である。

各ノード毎に別個のレジスタを持つことにより、ローカル領域に関しては各ノードの主記憶容量は異なっても良い。 10

【0031】

リモート判定回路132では、アクセスアドレスをAとすると、Aが自ノードのローカル領域の開始アドレス1503と終了アドレス1504の間にあるかを調べる。

さらに、自ノードの番号をN、共有領域ベースアドレスレジスタ1501の内容をB、共有領域サイズノードレジスタ1502の内容をSとすると、

$$B + S \times N \leq A < B + S \times (N + 1)$$

であるかどうかを調べる。

いずれかが成り立つ場合はバスコマンド受信ノバスコマンド分類回路131に、アクセスされたアドレスが内部であることを、そうでない場合はリモートの主記憶であることを知らせる。 20

その後、バスコマンド受信ノバスコマンド分類回路131は、アクセスコマンド、アクセスアドレスが内部であるかリモートであるかによって、異なった動作をする。

【0032】

(A)内部アドレスへのF、FIコマンドの場合

まず、131aを通してRATチェック回路133にアクセスアドレス、コマンドが送られる。

RATチェック回路133では、アクセスされたページに該当するRAT138の内容をチェックし、他のノードへのCCCが要するかどうかを判断する。

【0033】

図3にRATの詳細を示す。

RATでは、ページ毎にA、Nの2つのビットを記憶するメモリである。

Aビットは他のノードへのCCCが必要であるかどうかを示す。

Aが0の場合、該当するページが他のノードからアクセスされていない(従って他のノードへのCCCは不要である)ことを示し、Aが1の場合該当するページが他のノードからアクセスされたこと(従って、他のノードへのCCCが必要であること)を示す。

Aビットは立ち上げ時に0に初期化されるほか、ソフトウェアがあるページを割り当てる際にソフトウェアにより0にリセットされる。

【0034】

図9に本実施例におけるページ割り当てのアルゴリズムを示す。 40

これにより、ページがシステムソフトウェアにより回収され、再利用された場合にも、RATが正しく動作することを保証することができる。

Aビットはハードウェア(RAT変更回路148)で1にセットされる。これについては後ほど述べる。

【0035】

Nビットは通常のキャッシュコヒーレント管理が不要であるページを示すのに用いる。

Nビットが0の場合該当するページは通常のキャッシュコヒーレント管理の対象になることを示し、Nが1の場合は、該当するページは(ページ、フラッシュ等の命令を用いて)ソフトウェアによりコヒーレント管理が行われるため通常のキャッシュコヒーレント管理は不要であることを示す。 50

Nビットが1のページに対しては、他のノードへのCCCのためのコマンドは出されない。

I/O領域や、数値計算の配列がおかれる領域などで有効である。

Nビットはアプリケーションソフトウェアが定義したアドレス情報に基づき、システムソフトウェアが管理する。

【0036】

図2にRATチェック回路133の詳細を示す。

まず、1333、1334でアクセスされたページに該当するRATの値が読み出しされた後、ゲート1335により、 $A = 1$ かつ $N = 0$ であるかどうか(すなわち、他のノードへのCCCが必要であるかどうか)を判断する。

他のノードへのCCCが必要である場合、スイッチ1330、信号133aを通してネットワークコマンド生成回路134に他ノードへのコマンドを出すことを依頼するとともに、他ノードへのCCCが必要で、かつコマンドがFかFIの場合(コマンドデコード回路1332およびゲート1336により判断する)スイッチ1331、信号133bを通してCCC待ち回路AにCCCの結果を待ちあわせることを依頼する。

それに対して、他ノードへのCCCが必要でない場合($A = 0$ もしくは $N = 1$ の場合)は、ゲート1337でコマンドがFかFIかどうか判断され、FかFIで外部へのCCCが不要な場合には、スイッチ1338、信号133cを通じて主記憶アクセス回路A136に主記憶の内容を読み出すことを依頼する。

すなわち、外部へのCCCを省略し、即座に主記憶をアクセスできる。

【0037】

(A1)他のノードへのCCCが不要の場合

この場合、信号133cを通じて主記憶アクセス回路A136にアクセスアドレスが伝えられる。

主記憶アクセス回路A136はアクセスされたデータラインを自ノード内の主記憶160から読み出し(ローカル領域161、共有領域162の何れの場合も同じである)、読み出されたデータを136a、143を通じて自ノード内のCPUに返す。

【0038】

(A2)他のノードへのCCCが必要な場合

まず、信号133aを通じてネットワークコマンド生成回路134に他のノードにF、FIコマンドを送出することを依頼する。

それと同時に信号133bを通じてCCC待ち回路A135に対して他のノードからの返答を待ち合わせることを依頼する。

ネットワークコマンド生成回路134は、宛先生成回路139により判断されたノードにコマンドを送付する。

ここで、注意しなければならないのは、宛先ノードは図6、図7に示すように全ての(64個の)ノードに対応するビットマップで表されるため、複数のノードを指定すればブロードキャスト、マルチキャストを指定することができる。

【0039】

図11に宛先生成回路139の構成を示す。

宛先生成回路では、まず、アクセスアドレスが内部であるかリモートであるか(リモート判定回路1391)が判断され、さらに、アクセスアドレスが共有メモリであるかどうか(共有メモリ判定回路1392)が判断される。

1391の動作は132と同一である。

共有メモリ判定回路1392では、アクセスアドレスをA、共有領域ベースアドレスレジスタ1501の内容をB、共有領域サイズ/ノードレジスタ1502の内容をSとすると、

$$B - A < B + S \times 64$$

で有るかどうか調べられる(アクセスアドレスが上記の範囲内に入っていれば共有メモリである)。

10

20

30

40

50

さらに、回路1394により、アクセスアドレスのhomeノードが求められる。homeノードとは、アクセスアドレスを主記憶上に持つノードである。

具体的には以下の手順で求められる（以下ではアクセスアドレスをAとする）。

- ・Aと全てのノードに対応するローカル領域レジスタ1503、1504の値が比較され、Aが何れかのノードHLの開始アドレスと終了アドレスの間に有る場合、HLがhomeノードである（Aはローカル領域である）。

- ・さらに、共有領域ベースアドレスレジスタ1501の内容をB、共有領域サイズノードレジスタ1502の内容をSとすると、以下の計算を行う。

$HS = (A - B) / S$ （小数点以下切り捨て）

HSが0以上64未満の場合、HSがhomeノード番号である（Aは共有領域である） 10

。

【0040】

セレクタ1390は、上記の判断結果、及び、アクセスコマンドにより、図12に示す動作をする。

つまり、コマンドがF、FI、Iの場合は、共有領域の場合は全ノードへブロードキャストするためオール1のビットマップ1393aが選ばれ（ここでは説明を簡単にするため64個全ノードが実装されているとするが、実装されているノードが64未満でも、ビットマップを変えることにより対応できる）、ローカル領域の場合は、パーティション内のノードにのみマルチキャストするため、パーティションレジスタ1500の内容が選ばれる。 20

D、NDコマンドの場合、及びリモート主記憶からのDMコマンドの場合、アクセス元ノード番号レジスタ149の値をデコードした（ビットマップにした）結果が選ばれ、アクセス元のノードに結果が返送される。

リモートの主記憶へのWBコマンドの場合は、回路1394で求められたhomeノード番号をデコードした結果が選ばれ、データをhomeノードに書き戻す。

内部のアドレスへのWB、NDコマンドはエラーである（あり得ないはずである）。

この場合、F、FIコマンドであるので、ローカル領域の場合はパーティション内の全ノードにマルチキャストされ、共有領域の場合はシステム内の全ノードにブロードキャストされる。

他ノードへF、FIコマンドを出した結果は、DもしくはNDコマンドにより、返送される。 30

他のノードより返送されたコマンドはネットワークコマンドチェック回路141でパーティション/主記憶構成情報150との整合性をチェックされた後、ネットワークコマンド分類回路142に送付される。

【0041】

ネットワークコマンド分類回路142は、コマンドの種類、アクセスされたアドレスが内部（自ノードの主記憶）かリモート（他ノードの主記憶）か（リモート判定回路147を用いてを判断される、リモート判定回路147の動作は132と全く同一である）に応じ、該当する出力にコマンドを出す。

内部アドレスへのN、NDコマンドの場合は、信号142aを通じてCCC待ち回路A（135）にコマンドを送る。 40

【0042】

図14にネットワークコマンドチェック回路141の詳細を示す。

ネットワークコマンドチェック回路141では、ローカル領域へのパーティション外のノードへのアクセスをエラーとして検出するための回路である（ローカル領域を他のパーティションからの不正なアクセスから守るために、ローカル領域へはパーティション内のノードからのみしかアクセスが許されない）。

まず、アクセスアドレスが共有メモリか否かが共有メモリ判定回路1410により判断される（共有メモリ判定回路1410の動作は、1392と同一である）。

さらに、パーティション内外判定回路1411によりアクセス元のノードがパーティシ 50

ンの中であるかどうか判断される。

具体的には、ネットワークコマンドのアクセス元ノード番号をデコーダ1412でデコードした結果と、パーティションレジスタ1500（パーティション内のノードを表すビットマップ）をAND-ORゲートでゲートした結果により、アクセス元のノード番号に相当するビットがパーティションレジスタ1500の中で立っているかどうかチェックされる。

その結果、正しいアクセスの場合、つまりアクセスアドレスが共有メモリであるか、アクセス元のノードがパーティション内である場合は、スイッチ1415を通じてアクセスコマンドがネットワークコマンド分類回路142に送られる。

不正なアクセスの場合は、ゲート1416により、エラーが報告される（スイッチ1415によりアクセスは抑止される）。

これにより、ローカル領域への他のパーティションからの不正なアクセスを防止することができる。

【0043】

内部のアドレスに対して他のノードへF、FIコマンドを出した結果としては、全てのノードからNDが返る場合と、1つ以上のノードからDコマンドが返る（それ以外のノードからはNDが返る）場合に分けられる。

他のノードからの返答はCCC待ち回路A（135）で集計され、該当する動作がとられる。

【0044】

図10にCCC待ち回路A（135）の詳細を示す。

信号133bを通じて与えられたCCCアドレスはレジスタ13506に蓄えられる。

後に他のノードからCCCの返答（DもしくはNDコマンド）が来ると、まず、コンパレータ13507でCCCアドレス13506と比較される。

比較した結果が一致した場合は、コマンドがDであるかNDであるかが、デコーダ13500とゲート13501、13502で判断される。

Dコマンドが（1つでも）到来した場合（つまり他のキャッシュから最新のデータが送られて来た場合は、フリップフロップ（FF）13504がセットされ、Dコマンドが到来したことが記憶されると同時に、ラッチ13514にDコマンドの内容が記憶される（Dコマンドにより送られてきた値は、後にCPUに返される）。

それとは別に、D、NDコマンドが1個到来する毎に、ゲート13503を通じてカウンタ13505がカウントアップされ、今までに何個の返答が来たかが数えられる。

カウンタの値が信号13511aで示される、期待される返答数と一致すると、全ての返答が返ったとして、CCCの待ち合わせが終了する。

【0045】

その際にフリップフロップ13504が1の場合、スイッチ13515がONになり、ラッチ13514に記憶されていたDコマンドにより返送されたデータが信号135b、バスコマンド生成回路143を通じてCPUに返される。

それに対して、フリップフロップ13504が0の場合（全てのノードからNoDataコマンドが返された場合は、スイッチ13516がONになり、信号135aを通じて、CCCアドレスレジスタ13506に記憶されていたアクセスアドレスが主記憶アクセス回路A（136）に伝えられ、主記憶上のデータをアクセスすることが依頼される。

【0046】

信号13511aで示される期待される返答数は、次の手順で求められる。

まず、CCCアドレスレジスタ13506で記憶されているアクセスアドレスが共有主記憶であるかどうか、共有メモリ判定回路13508で判断され（共有メモリ判定回路13508の動作は1392と同一である）、セクタ13511に伝えられる。

セクタ13511は、アクセスアドレスが共有メモリの場合は64（システムの全ノード数、ただし、実装されているノード数が64未満の場合は実装されているノード数を指定する）を出力し、アクセスアドレスが共有メモリでない（ローカルメモリの場合は）は

10

20

30

40

50

パーティション内ノード数レジスタ1506の値を出力する。

これにより、該当するデータをアクセスするノードの数を求めることができる。

【0047】

(A2a)他のノードからの回答が全てNDの場合

CCC待ち回路A(135)より信号135aを通じて主記憶アクセス回路A(136)に、自ノード上の主記憶160をアクセスすることが依頼され、アクセスされたデータがCPUに返される。

【0048】

(A2b)何れかのノードよりDコマンドが返ってきた場合

CCC待ち回路A(135)より、信号135b、バスコマンド生成回路143を通じて、Dコマンドにより返送された最新のデータが、CPUに返される。 10

【0049】

(B)内部アドレスへのIコマンドの場合

この場合、RATチェック回路135により、他のノードへのCCCの要不要がチェックされるところまでは(A)と同様である。

【0050】

(B1)他のノードへのCCCが必要な場合

信号133aを通じてネットワークコマンド生成回路134に他のノードへのIコマンドの生成を依頼する(Iコマンドに返答はない)。

ネットワークコマンド生成回路134は宛先生成回路139の指定するノードにコマンドを送る。 20

つまり、ローカル領域へのIコマンドの場合、パーティション内の全ノードにマルチキャストされ、共有領域へのIコマンドの場合、システムの全ノードにブロードキャストされる。

【0051】

(B2)他のノードへのCCCが不要な場合

Iコマンドの場合、主記憶のデータアクセスは不要なので何も起こらない。

すなわち、コマンドデコード回路1332は0を出力するため、ゲート1337の出力は0であり、信号133cには何も出力されない。

【0052】

(C)内部アドレスへのWBコマンドの場合

この場合、バスコマンド受信/バスコマンド分類回路131は信号131cを通じて主記憶アクセス回路136に、書き戻されたデータを主記憶160に書き込むことを依頼する。

【0053】

(D)リモートアドレスへのF,FIコマンドの場合

この場合、信号131bを通じてネットワークコマンド生成回路に、他ノードへF,FIコマンドを生成する事を要求する。

ネットワークコマンド生成回路134は宛先生成回路139の指定するノードにコマンドを送る。 40

つまり、ローカル領域へのF,FIコマンドの場合、パーティション内の全ノードにマルチキャストされ、共有領域へのF,FIコマンドの場合、システムの全ノードにブロードキャストされる。

リモートアドレスへのF,FIに対しては、D,DM,NDの3種のコマンドが到来する。

ネットワークコマンド分類回路142は、リモートアドレスに対するD,DM,NDコマンドが到来した場合は、信号142eを通じて返答集計回路144に送出する。

【0054】

この場合、返答されるコマンドの組み合わせとしては、

(あ)homeノードからDMコマンドが、その他のノードからは全てNDが返る場合 50

(い) home ノードから DM コマンドが、他の何れか一つ以上のノードから D コマンドが返る (その他のノードからは ND コマンドが返る) 場合

(う) home ノードを含む何れか一つ以上のノードから D コマンドが返る (その他のノードからは ND コマンドが返る) 場合

に分けられる。

(あ) では DM コマンドより送られた home ノードの主記憶から読み出された値が使われるのに対し、(い) では、他のノードのキャッシュ上のデータを優先しなければならず (なぜならキャッシュ上のデータは変更されている可能性があるため)、DM コマンドで home ノードの主記憶から送られたデータは捨てられ、D コマンドによりキャッシュから送られたデータが活用される。

10

(う) では、D コマンドにより何れかのノードのキャッシュから送られてきたデータが用いられる。

各ノードからの返答は返答集計回路 144 で集計され、(あ)(い)(う)の何れかの場合であるかが判断され、アクセスされた結果 (データライン) がアクセス元の CPU に返される。

【0055】

図 13 に返答集計回路 144 の詳細を示す。

先ず、他ノードから送られてきたコマンドが 14400 によりデコードされる。

D コマンドの場合はラッチ 14401 により D コマンドの内容 (アドレス、データ) が記憶されるとともに、フリップフロップ 14403 により D コマンドが一つ以上到来したことが記憶される。

20

DM コマンドの場合はラッチ 14402 により DM コマンドの内容が記憶される。それと並行して、D、DM、ND コマンドが一つ到来する毎に、ゲート 14404 を通じてカウンタ 14405 がカウントアップされ、到着した返答の数を数える。

カウンタの値が信号 14413a で示される期待される返答数と一致した場合 (コンパレータ 14406 で判断される)、全ての返答が到着したと判断し、D コマンドが一つでも来ていた場合 (フリップフロップ 14403 が 1 の時) はゲート 14408、スイッチ 14410 により D コマンドの内容が、信号 144a を通じてバスコマンド生成回路 143 に供給される。

それに対し D コマンドが一つも来ていなかった場合は、ゲート 14407、スイッチ 14409 により、DM コマンドの内容がバスコマンド生成回路に供給される。

30

これにより、D コマンドが一つ以上来ていた場合は、D コマンドにより返送された他のノードのキャッシュ上のデータがアクセスを行った CPU に返され、D コマンドが一つも来ていなかった場合は、DM コマンドにより返送された、home ノードの主記憶上の値がアクセスを行った CPU に返される。

1411、1413 等の期待される返答数を求める回路は、アドレスラッチ 14414 にラッチされたアクセスアドレスに対して、CCC の返答の数を求める。

詳細な動作は、CCC 待ち回路 A (135) の回路 (13508、13511 等) と全く同一である。

【0056】

40

(E) リモートアドレスへの I コマンドの場合

この場合は RAT のチェックは行われず、信号 131b により直接ネットワークコマンド生成回路 134 に他のノードへのコマンドの送りの依頼が行われる。

その後の動作は (B1) と同様である。

【0057】

(F) リモートアドレスへの WB コマンドの場合

この場合、バスコマンド受信 / バスコマンド分類回路 131 は信号 131b を通じてネットワークコマンド生成回路 134 に、WB コマンドを送付することを依頼する。

宛先生成回路が出す宛先は、home ノードであり、WB コマンドは home ノードに送られる。

50

【 0 0 5 8 】

(4) 他のノードからのアクセスに対する動作

ここでは、主記憶アクセス回路 1 3 0 が、ネットワーク 9 0 0 を通じて送られてきた、他のノードからのコマンドに対してどのように動作するかを述べる。

ここで、D, DM, ND コマンドに対する動作は既に (3) で述べてあるので、その他の F, FI, I, WB に対する動作を述べる。

他のノードから送られてきたコマンドは、ネットワークコマンドチェック回路 1 4 1 によりチェックされた後、ネットワークコマンド分類回路 1 4 2 に送られる。

ネットワークコマンド分類回路 1 4 2 では、アクセスコマンドの種類、アクセスアドレスが内部かリモートか (リモート判定回路 1 4 7 により判断される) により、該当する出力にコマンドを送る。

10

また、F, FI, I コマンドの場合は、アクセス元ノード番号レジスタ 1 4 9 にアクセス元ノード番号をセットする。

【 0 0 5 9 】

(A) 内部アドレスへの F, FI コマンドの場合

ネットワークコマンド分類回路 1 4 2 は、信号 1 4 2 d を通じてバスコマンド生成回路 1 4 3 にバスコマンドの生成を依頼する。

バスコマンド生成回路 1 4 3 は、ノード内の共有バス 1 2 0 を通じてノード内の CPU に F, FI コマンドを出す。

それと同時に、信号 1 4 2 c を通じて CCC 待ち回路 B (1 4 5) に、ノード内の共有バスに出されたコマンドを待ち合わせることを依頼する。

20

【 0 0 6 0 】

図 1 5 に CCC 待ち回路 B (1 4 5) の詳細を示す。

信号 1 4 2 c を通じて送られてきた、CCC を行うアドレスは、CCC アドレスレジスタ 1 4 5 1 に記憶される。

後にノード内の共有バス 1 2 0 からバスコマンド受信 / バスコマンド分類回路 1 3 1 を通じて、CCC の返答 (D もしくは ND コマンド) が来ると、まず、コンパレータ 1 4 5 2 で CCC アドレス 1 4 5 1 と比較される。

比較した結果が一致した場合は、コマンドが ND であるかどうか、デコーダ 1 4 5 0 とゲート 1 4 5 3 で判断される。

30

ND コマンドが返送された場合は、スイッチ 1 4 5 5 を通じて、CCC アドレスレジスタ 1 4 5 1 に記憶されていたアクセスアドレスが、信号 1 4 5 a を通じて主記憶アクセス回路 B 1 4 6 に伝えられ、主記憶上のデータをアクセスすることが依頼される。

【 0 0 6 1 】

ノード内の共有バスに出した F, FI コマンドの結果は、D (ノード内のいずれかの CPU のキャッシュに該当するデータがある場合)、もしくは ND (ノード内の何れの CPU のキャッシュにも該当するデータが無い場合) により返送される。

D コマンドの場合は、データも同時に返送される。

バスコマンド受信 / バスコマンド分類回路 1 3 1 は、コマンド、および、アクセスアドレスが内部であるか、リモートであるか (リモート判定回路 1 3 2 により判定される) に応じて、該当する出力にコマンドを送る。

40

【 0 0 6 2 】

(A a) 内部アドレスへのアクセスに対し、D コマンドが返送された場合

自ノードのキャッシュ上にアクセスデータがあった場合、バスコマンド受信 / バスコマンド分類回路 1 3 1 は、D コマンド (コマンド、アドレス、データ) を、信号 1 3 1 b を通じて、ネットワークコマンド生成回路 1 3 4 に送る。

ネットワークコマンド生成回路 1 3 4 は、アクセス元のノードにアクセスデータを D コマンドを用いて返送する (宛先生成回路 1 3 9 は、アクセス元ノード番号 1 4 9 を選択する) 。

これにより、自ノードの CPU のキャッシュ上のデータがアクセス元のノードに返送され

50

る。

【 0 0 6 3 】

(A b) 内部アドレスへのアクセスに対し、 N D コマンドが返送された場合。

【 0 0 6 4 】

自ノードのキャッシュ上にはアクセスデータが無かった場合、バスコマンド受信/バスコマンド分類回路 1 3 1 は、 N D コマンドを信号 1 3 1 d を通じて C C C 待ち回路 B (1 4 5) に送付する。

C C C 待ち回路 B は信号 1 4 5 a を通じて、主記憶アクセス回路 B (1 4 6) に主記憶 1 6 0 上のアクセスアドレスの内容 (キャッシュライン) を読み出すことを依頼する。

主記憶アクセス回路 B (1 4 6) は、信号 1 6 0 a を通じてアクセスラインを読み出すと、信号 1 4 6 a を通じてネットワークコマンド生成回路 1 3 4 に送る。 10

ネットワークコマンド生成回路 1 3 4 は、アクセス元のノードにアクセスデータを D M コマンドを用いて返送する (宛先生成回路 1 3 9 は、アクセス元ノード番号 1 4 9 を選択する) 。

これにより、自ノードの主記憶上のデータがアクセス元のノードに返送される。

それと同時に、主記憶アクセス回路 B (1 4 6) は、信号 1 4 6 b を通じて R A T 変更回路 1 4 8 にアクセスされたページに該当する、 R A T 1 3 8 の A ビットに 1 を設定することを依頼する。

【 0 0 6 5 】

図 4 に R A T 変更回路 1 4 8 の詳細を示す。 20

アクセスされたアドレスのページ番号に対応する R A T の A ビットに 1 が書き込まれる。

【 0 0 6 6 】

(B) 内部アドレスへの I コマンドの場合

ネットワークコマンド分類回路 1 4 2 は、信号 1 4 2 d を通じてバスコマンド生成回路 1 4 3 にバスコマンドの生成を依頼する。

バスコマンド生成回路 1 4 3 は、ノード内の共有バス 1 2 0 を通じてノード内の C P U に I コマンドを出す (I コマンドには返答は無い) 。

【 0 0 6 7 】

(C) 内部アドレスへの W B コマンドの場合

ネットワークコマンド分類回路 1 4 2 は、信号 1 4 2 b を通じて、主記憶アクセス回路 B (1 4 6) に、 W B されたデータを主記憶に書き込むことを依頼する。 30

主記憶アクセス回路 B (1 4 6) は、信号 1 6 0 a を通じて主記憶 1 6 0 へ W B されたデータを書き込む。

それと同時に、主記憶アクセス回路 B (1 4 6) は、信号 1 4 6 b を通じて R A T 変更回路 1 4 8 にアクセスされたページに該当する、 R A T 1 3 8 の A ビットに 1 を設定することを依頼する。

【 0 0 6 8 】

(D) リモートアドレスへの F , F I コマンドの場合

ネットワークコマンド分類回路 1 4 2 は、信号 1 4 2 d を通じてバスコマンド生成回路 1 4 3 にバスコマンドの生成を依頼する。 40

バスコマンド生成回路 1 4 3 は、ノード内の共有バス 1 2 0 を通じてノード内の C P U に F , F I コマンドを出す。

ノード内の共有バスに出した F , F I コマンドの結果は、 D (ノード内のいずれかの C P U のキャッシュに該当するデータがある場合) 、もしくは N D (ノード内の何れの C P U のキャッシュにも該当するデータが無い場合) により返送される。

D コマンドの場合は、データも同時に返送される。

【 0 0 6 9 】

リモートアドレスに対する D , N D コマンドの場合、何れの場合も、バスコマンド受信/バスコマンド分類回路 1 3 1 は、 D コマンド (コマンド、アドレス、データ) もしくは N D コマンド (コマンド、アドレス) を、信号 1 3 1 b を通じて、ネットワークコマンド生 50

成回路134に送る。

ネットワークコマンド生成回路134は、アクセス元のノードにアクセス結果をD又はNDコマンドを用いて返送する(宛先生成回路139は、アクセス元ノード番号149を選択する)。

これにより、Dコマンドの場合、自ノードのCPUのキャッシュ上のデータがアクセス元のノードに返送され、NDコマンドの場合、自ノードのCPUのキャッシュ上には該当するデータが無いことが伝えられる。

【0070】

(E)リモートアドレスへのIコマンドの場合

この場合、内部アドレスへのIコマンドの場合(B)と全く同じ動作を行う。

【0071】

(F)リモートアドレスへのWBコマンドの場合

他のノードからリモートアドレスへのWBコマンドが来ることはあり得ない。

【0072】

ネットワークコマンド分類回路142はエラーを報告する。

【0073】

以上述べた手順によりクロスバネットワークにより接続されたノード間でキャッシュコピーレンスをとることができる。

その際に、リモートアクセステーブル(RAT)138を用いて他のノードへのCCCを削減することができる。

さらに、主記憶共有マルチプロセッサ内をパーティションに分ける際に、主記憶を、各パーティションのローカル領域/パーティション間共有領域の2つに分け、ローカル領域に対しては各パーティション間で独立したアドレス空間とすることにより、各パーティションのアドレスを0番地より始めるようにすることができる。

さらに、パーティション/主記憶構成情報150をレジスタにより記憶することにより、少ないハードウェア量でパーティションの管理を実現することができる。

【0074】

<変形例>

本発明は以上の実施の形態に限定されるのではなくいろいろの変形例にも適用可能である。

例えば、

(1)以上においては、ノード内のCPUはバス120により接続されているが、その他の接続形態(スイッチによる接続、主記憶アクセス回路130への一対一接続)も可能である。

(2)また、以上においては、RAT148は、主記憶アクセス回路130に内蔵される、専用メモリにより、構成されているが、外付けのメモリにすることも可能である。

さらに、RATを物理的にはローカル主記憶161上に置くことも可能である。

さらに、この場合、RATの内容を主記憶アクセス回路130内にキャッシングする事も可能である。

【0075】

(3)上記実施の形態において、パーティション/主記憶構成情報150では、共有領域の構成情報は、共有領域ベースアドレス1501及び、共有領域サイズ/ノードレジスタ1502の2つのレジスタにより記憶されている。

しかし、ローカル領域レジスタ1503、1504のように、ノード毎に開始アドレス、終了アドレスの組で覚えることも可能である(これにより、各ノードの共有領域の大きさを変えることができる)。

この場合、上記実施の形態においては共有領域ベースアドレス1501及び、共有領域サイズ/ノードレジスタ1502を用いて計算で求めていた、共有メモリ判定回路、homeノード判定回路は、コンパレータにより構成することができる。

【0076】

10

20

30

40

50

(4) 以上においては、ノード内の各CPU(110~112)は独立したキャッシュを持っているが、複数のCPUに共有される外付けの2次キャッシュを持たせることも可能である。

各ノードの主記憶アクセス回路130に、CPUのキャッシュTAGのコピーを持ち、他のノードから到来するキャッシュコヒーレントトランザクションをフィルタリングする事も可能である。

【0077】

(5) 以上においては、ノード間はクロスバネットワークにより接続されているが、他の形式のネットワーク(多段網等)により接続することも可能である。

(6) 上記実施の形態においては、他のノードへのネットワークコマンドを出す際に、宛先のノードをビットマップにより指定することにより、ブロードキャスト、マルチキャストを指示しているが、ネットワークコマンド生成回路134が、宛先のノード毎に複数のコマンドを出すことにより、ブロードキャスト、マルチキャストを実現することも可能である。

【0078】

【発明の効果】

本発明によれば、主記憶共有型のマルチプロセッサにおいて、リモートアクセステーブル(RAT)を置くことにより、少ないハードウェア量の追加により、自ノード内の主記憶をアクセスする際に、他のノードからアクセスされていないページに関しては、ノード間のCCCを省くことができる。

従って、アクセスレテンシを削減するとともに、ノード間のネットワークにかかるトラフィックを削減することができる。

さらに、本発明によれば、主記憶共有マルチプロセッサ内を複数のパーティションに分ける際に、主記憶をパーティション毎のローカルな領域とパーティション間で共有される領域に分け、ローカル領域に対しては各ノードのアドレス空間を独立することにより、各パーティションの開始アドレスを0番地から始めることを可能にすることができる。

さらに、各パーティションのローカル領域、共有領域の範囲をレジスタにより記憶することにより、パーティション管理のためのハードウェア量を従来のページ毎の管理と比較して大幅に削減することができる。

【図面の簡単な説明】

【図1】本発明のキャッシュコヒーレンス保持機構を持つ主記憶共有型マルチプロセッサである。

【図2】各ノードのRATチェック回路のブロック図である。

【図3】各ノードのRATのブロック図である。

【図4】各ノードのRAT変更回路のブロック図である。

【図5】各ノードにおいて、パーティション/主記憶構成情報を表すレジスタ群の詳細図である。

【図6】ネットワーク上の、F, FI, I, NDコマンドのパケットフォーマットである。

【図7】ネットワーク上の、WB, D, DMコマンドのパケットフォーマットである。

【図8】各ノードの主記憶のブロック図である。

【図9】本発明のマルチプロセッサシステムにおいて、主記憶のページをアロケートする際のフロー図である。

【図10】各ノードのCCC待ち回路Aのブロック図である。

【図11】各ノードの宛先生成回路のブロック図である。

【図12】各ノードの宛先生成回路内の、宛先セクタの入力と出力の関係の表を示す図である。

【図13】各ノードの返答集計回路のブロック図である。

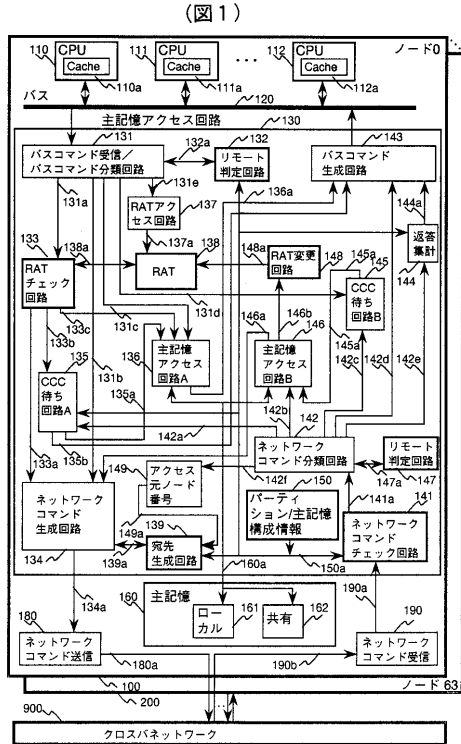
【図14】各ノードのネットワークコマンドチェック回路のブロック図である。

【図15】各ノードのCCC待ち回路Bのブロック図である。

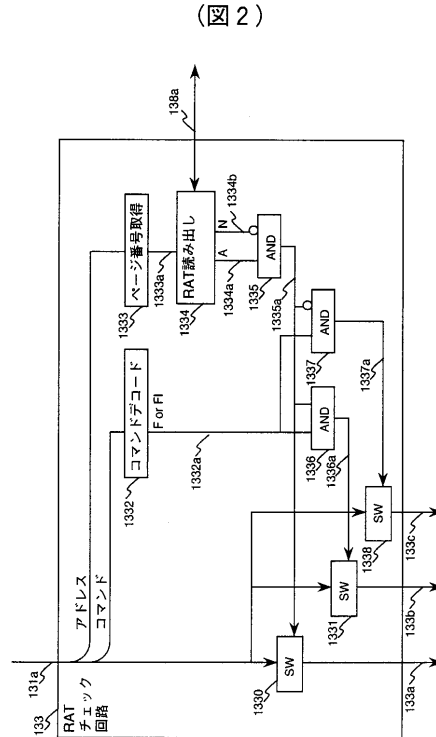
【符号の説明】

1 0 0、2 0 0	ノード	
1 1 0 ~ 1 1 2	C P U	
1 1 0 a ~ 1 1 2 a	キャッシュ	
1 2 0	バス	
1 3 1	バスコマンド受信ノバスコマンド分類回路	
1 3 2、1 4 7	リモート判定回路	
1 3 3	R A Tチェック回路	
1 3 4	ネットワークコマンド生成回路	
1 3 5	C C C待ち回路 A	10
1 3 6	主記憶アクセス回路 A	
1 3 7	R A Tアクセス回路	
1 3 8	R A T (R e m o t e A c c e s s T a b l e)	
1 3 9	宛先生成回路	
1 4 1	ネットワークコマンドチェック回路	
1 4 2	ネットワークコマンド分類回路	
1 4 3	バスコマンド生成回路	
1 4 4	返答集計回路	
1 4 5	C C C待ち回路 B	
1 4 6	主記憶アクセス回路 B	20
1 4 8	R A T変更回路	
1 4 9	アクセス元ノード番号レジスタ	
1 5 0	パーティションノ主記憶構成情報	
1 6 0	主記憶	
1 6 1	ローカル主記憶	
1 6 2	共有主記憶	
1 8 0	ネットワークコマンド送信回路	
1 9 0	ネットワークコマンド受信回路	
9 0 0	ノード間クロスバネットワーク	

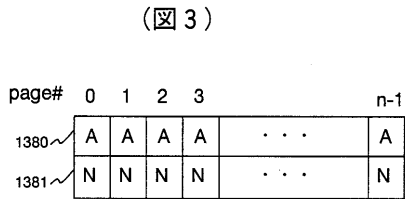
【 図 1 】



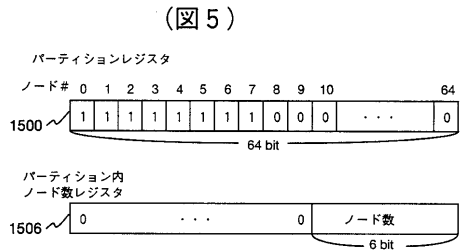
【 図 2 】



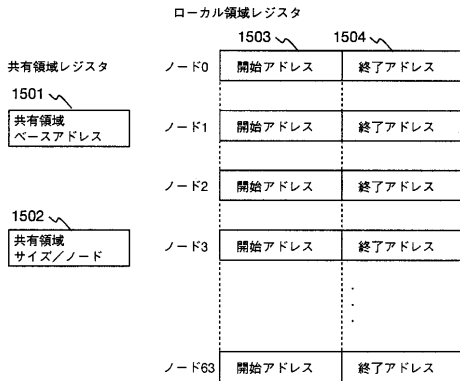
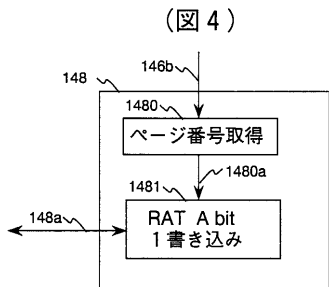
【 図 3 】



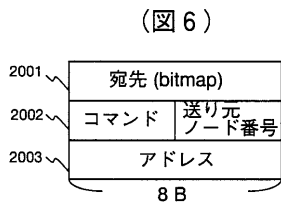
【 図 5 】



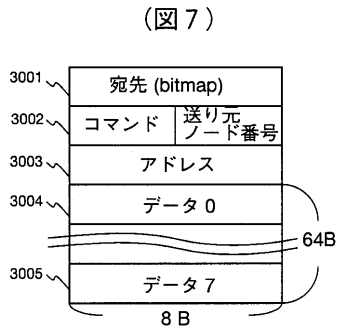
【 図 4 】



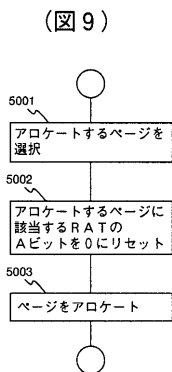
【 図 6 】



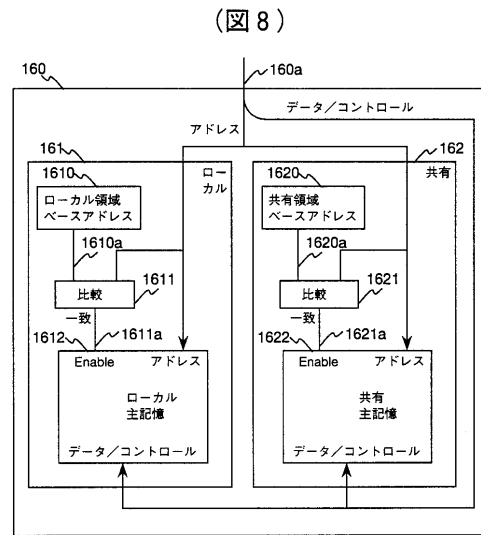
【 図 7 】



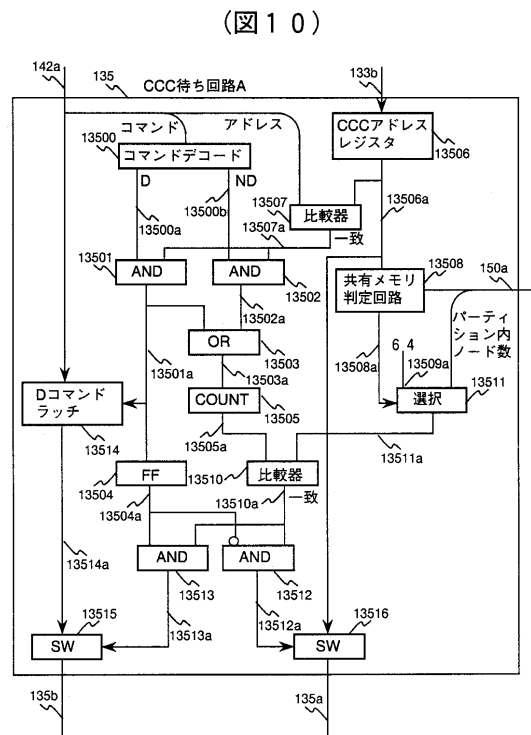
【 図 9 】



【 図 8 】

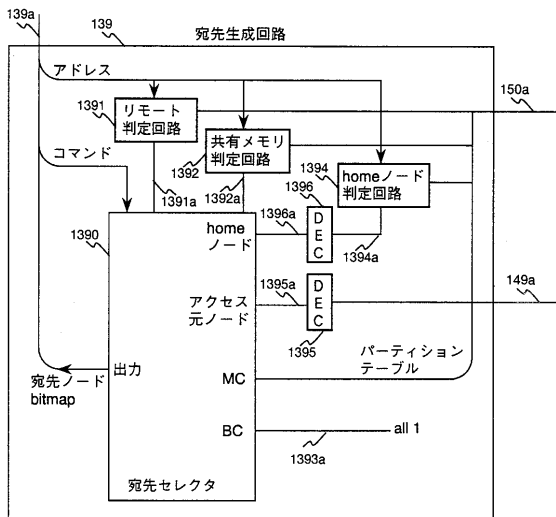


【 図 10 】



【 図 1 1 】

(図 1 1)



【 図 1 2 】

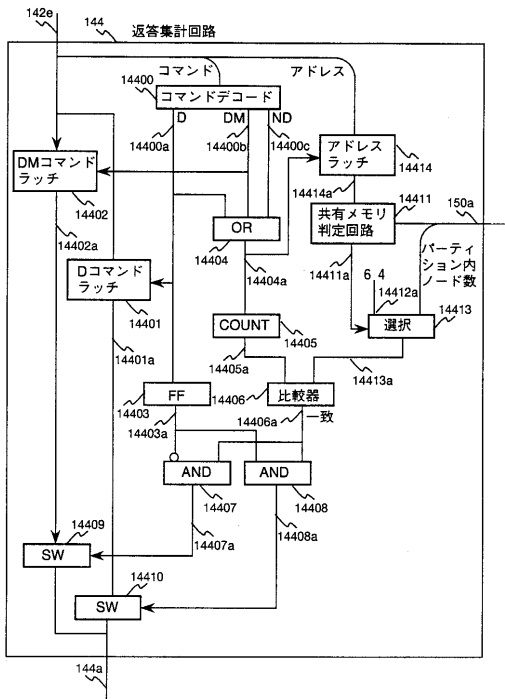
(図 1 2)

宛先セクタの出力

アクセス コマンド	アクセスアドレス	
	内部	リモート
F	アクセスアドレスに応じて分岐 ・ローカル領域→MC ・共有領域 →BC	
F I		
I		
WB	homeノード	
D、DM	アクセス元ノード	
ND	アクセス元ノード	

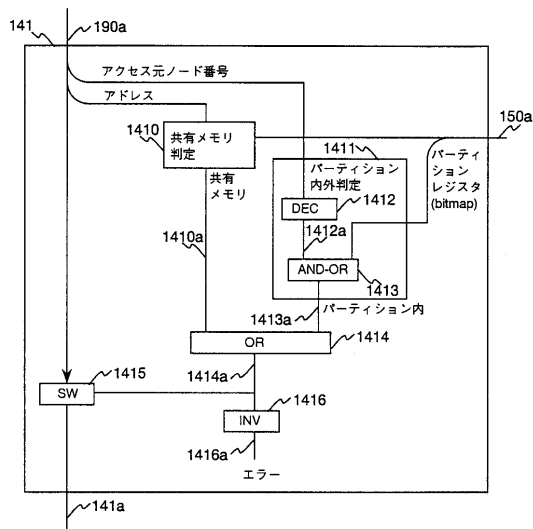
【 図 1 3 】

(図 1 3)



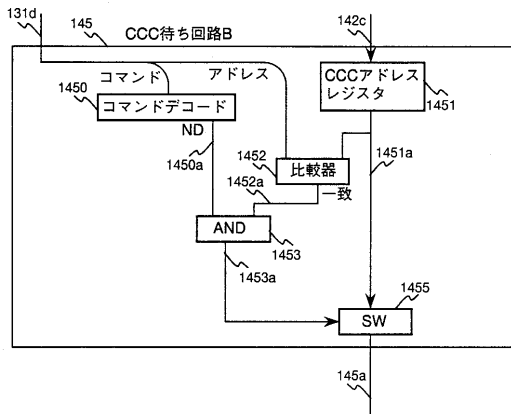
【 図 1 4 】

(図 1 4)



【 図 1 5 】

(図 1 5)



フロントページの続き

- (72)発明者 庄内 亨
東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内
- (72)発明者 大河内 俊夫
東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内
- (72)発明者 明石 英也
東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所 中央研究所内

審査官 鳥居 稔

- (56)参考文献 特開平01-109464(JP,A)
特開平06-274461(JP,A)
特開平07-210520(JP,A)
特開平07-160581(JP,A)
特開平08-016469(JP,A)
特開平08-016470(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G06F 12/08
G06F 15/16-177