

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号
特許第4718012号
(P4718012)

(45) 発行日 平成23年7月6日 (2011.7.6)

(24) 登録日 平成23年4月8日 (2011.4.8)

(51) Int.Cl. F I

G O 6 F 12/08 (2006.01)

G O 6 F 12/08 5 3 1 B

G O 6 F 12/08 5 0 1 C

G O 6 F 12/08 5 5 1 C

G O 6 F 12/08 5 7 5

請求項の数 10 (全 27 頁)

(21) 出願番号	特願2000-590062 (P2000-590062)	(73) 特許権者	591016172
(86) (22) 出願日	平成11年8月26日 (1999.8.26)		アドバンスト・マイクロ・ディバイズ・
(65) 公表番号	特表2002-533813 (P2002-533813A)		インコーポレイテッド
(43) 公表日	平成14年10月8日 (2002.10.8)		ADVANCED MICRO DEVI
(86) 国際出願番号	PCT/US1999/019856		CES INCORPORATED
(87) 国際公開番号	W02000/038070		アメリカ合衆国、94088-3453
(87) 国際公開日	平成12年6月29日 (2000.6.29)		カリフォルニア州、サニibel、ピー・
審査請求日	平成18年7月12日 (2006.7.12)		オウ・ボックス・3453、ワン・エイ・
(31) 優先権主張番号	09/217,699		エム・ディ・プレイス、メイル・ストップ
(32) 優先日	平成10年12月21日 (1998.12.21)		・68 (番地なし)
(33) 優先権主張国	米国 (US)	(74) 代理人	100064746
(31) 優先権主張番号	09/217,212		弁理士 深見 久郎
(32) 優先日	平成10年12月21日 (1998.12.21)	(74) 代理人	100085132
(33) 優先権主張国	米国 (US)		弁理士 森田 俊雄

最終頁に続く

(54) 【発明の名称】 メモリキャンセルメッセージを用いたシステムメモリ帯域幅の節約およびキャッシュコヒーレンシ維持

(57) 【特許請求の範囲】

【請求項 1】

マルチプロセッシングコンピュータシステムであって、
相互接続構造を介して相互接続される複数の処理ノードを含み、前記複数の処理ノードは、
指定されたメモリ位置からデータを読み出す第1の読出動作を開始するよう構成される第1の処理ノードと、
第2の処理ノードとを含み、前記第2の処理ノードは、前記第1の読出動作に応答して、前記第2の処理ノードに連結されたメモリ内の前記指定されたメモリ位置からデータを読み出す第2の読出動作を開始するよう構成され、前記第2の処理ノードは、さらに、前記第1の読出動作に
10 応答してプローブを発行するよう構成され、前記マルチプロセッシングコンピュータシステムはさらに、
第3の処理ノードを含み、前記第3の処理ノードは、前記第2の処理ノードからのプローブを受取るように、かつ、前記プローブに
20 応答して、前記指定されたメモリ位置に対応し、前記第3の処理ノード内に記憶された、変更されたデータを検出するよう連結され、前記第3の処理ノードは、前記指定されたメモリ位置の変更されたコピーを前記第3の処理ノード内に検出すると、前記第2の処理ノードにメモリキャンセル応答を送信するよう構成され、前記メモリキャンセル応答は、前記第2の処理ノードに前記第2の読出動作のさらなる処理を打切らせ、
前記第2の処理ノードは、前記第2の読出動作の間に読出された前記データを、前記第

1の処理ノードに第1の読出応答を送信することにより転送するよう構成され、前記メモリキャンセル応答は、前記第2の処理ノードが前記メモリキャンセル応答を前記第1の読出応答の送信前に受信した場合に、前記第2の処理ノードに前記第1の読出応答の送信をキャンセルさせ、

前記第3の処理ノードは前記メモリキャンセル応答と並行に第2の読出応答を送信するよう構成され、前記第2の読出応答は前記第1の処理ノードに送信される、マルチプロセッシングコンピュータシステム。

【請求項2】

前記第2の処理ノードは、前記指定されたメモリ位置が前記第3の処理ノード内にキャッシュされているか否かに拘らず、プローブコマンドを送信するよう構成される、請求項1に記載のマルチプロセッシングコンピュータシステム。

【請求項3】

前記第2の読出動作の間に読出された前記データのサイズは、前記第1の読出動作のタイプに依存する、請求項1に記載のマルチプロセッシングコンピュータシステム。

【請求項4】

前記第2の読出応答は、前記第3の処理ノード内にキャッシュされた前記指定されたメモリ位置の前記変更されたコピーを含むデータバケットを含み、

前記第2の処理ノードは、前記第3の処理ノードから前記メモリキャンセル応答を受信すると、前記第1の処理ノードにtarget done応答を送信するよう構成され、前記target done応答は、前記第1の読出応答が送信されるか否かに拘らず送信される、請求項1に記載のマルチプロセッシングコンピュータシステム。

【請求項5】

前記第1の処理ノードは、前記target done応答および前記第2の読出応答を受信すると、前記第2の処理ノードにsource doneメッセージを送信するよう構成される、請求項4に記載のマルチプロセッシングコンピュータシステム。

【請求項6】

相互接続構造を介して相互接続される複数の処理ノードを含むマルチプロセッシングコンピュータシステムにおいて、前記複数の処理ノードは第1の処理ノードと、第2の処理ノードと、第3の処理ノードとを含み、前記第2の処理ノードに関連のメモリ内のメモリ位置の内容を選択的に読出するための方法であって、

前記第1の処理ノードによる前記メモリ位置の前記内容を読出す第1の読出動作を開始するステップと、

前記第1の読出動作に回答して、前記第2の処理ノードによる第2の読出動作をさらに開始するステップとを含み、前記第2の処理ノードは、前記第2の読出動作の間に前記第2の処理ノードに連結されたメモリ内の前記メモリ位置の前記内容を読出し、前記第2の処理ノードは、さらに、前記第1の読出動作に回答してプローブを発行し、前記第2の読出動作は前記第2の処理ノードから前記第1の処理ノードへの第1の読出応答を含み、前記第1の読出応答は前記メモリ位置の前記内容に対する第1のデータバケットを含み、方法はさらに、

前記第3の処理ノードが、前記第2の処理ノードからのプローブを受取り、前記プローブに回答して、前記メモリ位置に対応し、前記第3の処理ノード内に記憶された、変更されたデータを検出し、前記第3の処理ノードが、前記第3の処理ノード内に前記メモリ位置の変更されたコピーを検出すると、第2の処理ノードにメモリキャンセル応答を送信するステップと、前記メモリキャンセル応答が、前記第2の処理ノードに前記第2の読出動作のさらなる処理を打切らせるステップと、

前記メモリキャンセル応答が、前記第2の処理ノードが前記メモリキャンセル応答を前記第1の読出応答の送信前に受信した場合に、前記第2の処理ノードに前記第1の読出応答の送信をキャンセルさせるステップと、

前記第3の処理ノードが前記メモリキャンセル応答と並行に第2の読出応答を送信するステップとを含み、前記第2の読出応答は前記第1の処理ノードに送信される、方法。

10

20

30

40

50

【請求項 7】

前記第 1 のデータパケットのサイズは、前記第 1 の読出動作のタイプに依存する、請求項 6 に記載の方法。

【請求項 8】

前記第 2 の読出応答は、前記第 3 の処理ノード内にキャッシュされた前記メモリ位置の前記変更されたコピーを含む第 2 のデータパケットを含む、請求項 6 に記載の方法。

【請求項 9】

前記第 2 の処理ノードが、前記第 3 の処理ノードから前記メモリキャンセル応答を受信すると、前記第 1 の処理ノードに target done 応答を送信するステップをさらに含み、前記 target done 応答は、前記第 1 の読出応答が送信されるか否かに拘らず送信される、請求項 8 に記載の方法。

10

【請求項 10】

前記第 1 の処理ノードが、前記 target done 応答および前記第 2 の読出応答を受信すると、前記第 2 の処理ノードに source done メッセージを送信するステップをさらに含む、請求項 9 に記載の方法。

【発明の詳細な説明】**【0001】****【発明の背景】****1. 技術分野**

この発明は広くはコンピュータシステムに関し、より特定のには、マルチプロセッシング演算環境を達成するためのメッセージ通信方式に関する。

20

【0002】**2. 関連技術分野の背景**

一般的には、パーソナルコンピュータ (PC) およびその他の種類のコンピュータシステムは、メモリにアクセスするために共用バスシステムを中心に設計されてきた。1 つ以上のプロセッサおよび 1 つ以上の入力 / 出力 (I/O) 装置が、共用バスを介してメモリに結合される。I/O 装置は I/O ブリッジを介して共用バスに結合される場合もあり、該 I/O ブリッジは共用バスと I/O 装置との間の情報の転送を管理する。プロセッサは典型的には、直接またはキャッシュ階層構造を介して、共用バスに結合される。

【0003】

残念ながら、共用バスシステムはいくつかの欠点を有する。たとえば、共用バスには多数の装置が装着されることから、バスは典型的には比較的低い周波数で動作される。さらに、共用システムを介したシステムメモリ読出および書込サイクルは、プロセッサ内のキャッシュが関連するか、または 2 つ以上のプロセッサが関連する情報転送よりも、かなり長い時間を必要とする。共用バスシステムの他の欠点は、より多くの装置に対するスケーラビリティの欠如である。上述のように、帯域幅が固定される (そして、もし付加的な装置の追加によってバスの動作可能周波数が減じられると、減少し得る)。バスに (直接的にまたは間接的に) 装着された装置の帯域幅要件が、一旦バスの利用可能な帯域幅を超えると、装置はバスへのアクセスを試みたときにしばしばストールし得る。限られたシステムメモリ帯域幅を節約する機構を提供しない限り、全体的な性能は減じられるであろう。

30

40

【0004】

ノンキャッシュシステムメモリに対しアドレス指定された書込または読出動作は、2 つのプロセッサの間の、またはプロセッサとその内部キャッシュとの間の同様の動作よりも、より多くのプロセッサクロックサイクルをとる。バス帯域幅への制限は、システムメモリへの読出または書込のための長いアクセス時間とあいまって、コンピュータシステム性能に悪影響を及ぼす。

【0005】

上述の問題のうち 1 つまたはいくつかには、分散メモリシステムを用いて対処し得る。分散メモリシステムを用いるコンピュータシステムは、複数のノードを含む。2 つ以上のノードがメモリに接続され、それらのノードは何らかの好適な相互接続を用いて相互接続さ

50

れる。たとえば、ノードの各々は専用ラインを用いて他のノードに互いに接続されることが出来る。これに代えて、ノードの各々は固定された数の他のノードに接続され、トランザクションは、第1のノードから1つ以上の中間ノードを介して、第1のノードに直接接続されていない第2のノードに経路制御されてもよい。メモリアドレス空間は、各々のノードのメモリにわたって割当てられる。

【0006】

ノードはさらに、1つ以上のプロセッサを含み得る。プロセッサは典型的には、メモリから読出したデータのキャッシュブロックをストアするキャッシュを含む。さらに、ノードはプロセッサの外部の1つ以上のキャッシュを含み得る。プロセッサおよび/またはノードは、他のノードからアクセスされるキャッシュブロックをストアし得るために、ノード内のコヒーレンスを維持するための機構が望まれる。

10

【0007】

EP-A-0 379 771は、デジタルコンピュータのためのメモリ制御システムを開示するが、ここでは、要求されたデータの変更された状態のものが別の関連のCPUキャッシュ内に利用可能である場合に、関連するメインメモリを備えたシステム制御ユニット(SCU)の読出が、関連の中央演算ユニット(CPU)に应答して自動的に打切られる。

C. A. プリート(C. A. Prete)による「密結合マルチプロセッサシステムのためのRTSキャッシュメモリ設計(RTS Cache Memory Design for a Tightly Coupled Multiprocessor System)」IEEE Micro. US, IEEE Inc. New York Vol.11 No.2、1991年4月11日、pp.16-19、40-52、は、縮小状態遷移(Reduced State Transitions)として知られる、マルチプロセッサシステムにおけるキャッシュメモリのためのコヒーレンスプロトコルを開示する。読出または書込動作の際にキャッシュミスが起こると、キャッシュはまずキャッシュコピーを置換え、次いで要求された動作を実行する。置換えの段階は、ヴィクティムキャッシュブロックを選択するステップと、ヴィクティムコピーに関連するメモリブロックを更新するステップと、読出ブロックトランザクションにより、要求されたメモリブロックを読出すステップと、からなる。

20

上に概略を述べた問題は、ここに説明するコンピュータシステムによってほとんどが解決される。コンピュータシステムは多数の処理ノードを含むことができ、そのうち2つ以上は分散メモリシステムを形成し得る別々のメモリに結合し得る。処理ノードはキャッシュを含むことができ、コンピュータシステムは、キャッシュと分散メモリシステムとの間のコヒーレンスを維持し得る。

30

【0008】

この発明の第1の局面によると、コンピュータシステムが提供され、該コンピュータシステムは読出トランザクションによりアドレス指定されたデータの変更されたコピーを保持していることに应答して、該読出トランザクションに対応するメモリキャンセル应答を送信するよう構成される第1の処理ノードを含み、該第1の処理ノードは、(i)第2の処理ノードからプローブを受取り、(ii)該読出トランザクションによってアドレス指定されたデータの変更されたコピーを検出し、(iii)該変更されたコピーの検出に应答して該メモリキャンセル应答を送信するよう構成される、処理ノードであり、該コンピュータシステムはさらに該第2の処理ノードを含み、該第2の処理ノードはトランザクションのターゲットノードを含み、かつシステムメモリの少なくとも一部に結合され、該システムメモリの少なくとも一部は、読出トランザクションによってアドレス指定されたデータに対応する記憶位置を含み、該第2の処理ノードは該第1の処理ノードから該メモリキャンセル应答を受け取るよう結合され、該第2の処理ノードは、該メモリキャンセル应答に应答して、該記憶位置への読出サイクルのさらなる処理を打ち切るよう構成される。

40

一実施例においては、処理ノードは複数のデュアル単方向リンクを介して相互接続される。単方向リンク対の各々は、処理ノードのうちの2つを接続するコヒーレントなリンク構造を形成する。単方向リンク対の一方のリンクは、第1の処理ノードから信号を、その単方向リンク対を介して接続された第2の処理ノードに送る。単方向リンク対の他方のリン

50

クは、信号の逆のフローを運ぶ。すなわち、信号を第2の処理ノードから第1の処理ノードへ送る。こうして、単方向リンクの各々は、パケット化情報転送のために設計されたポイントツーポイント相互接続を形成する。2つの処理ノード間の通信は、システム内の1つ以上の残りのノードを介して経路制御されることもある。

【0009】

処理ノードの各々は、メモリバスを介してそれぞれのシステムメモリに結合されることができる。メモリバスは双方向であってもよい。処理ノードの各々は、少なくとも1つのプロセッサコアを含み、かつ選択によりそれぞれのシステムメモリと通信するためのメモリコントローラを含み得る。1つ以上のI/Oブリッジを介したさまざまなI/O装置との接続性を可能にするため、1つ以上の処理ノードに他のインターフェースロジックを含んでもよい。

10

【0010】

一実施例においては、1つ以上のI/Oブリッジを、1組の非コヒーレントなデュアル単方向リンクを介してそれぞれの処理ノードに結合し得る。これらのI/Oブリッジは、この非コヒーレントなデュアル単方向リンクの組を介してそれらのホストプロセッサと通信するが、これは2つの直接リンクされたプロセッサがコヒーレントなデュアル単方向リンクを介して互いと通信するのとはほぼ同じ方法である。

【0011】

プログラム実行の間のある時点で、キャッシュ内にメモリデータのダーティコピーを持つ処理ノードは、その変更されたデータを含むキャッシュブロックを捨てることができる。一実施例においては、その処理ノード（ソースノードとも呼ばれる）はヴィクティムブロックコマンドをキャッシュされたダーティなデータと併せて第2の処理ノードに送信するが、該第2の処理ノードはすなわち、キャッシュされたデータのための対応するメモリ位置を有するシステムメモリの一部に結合されたものである。この第2の処理ノード（ターゲットノードとも呼ばれる）は、応答してターゲット終了メッセージを送信処理ノードに送り、メモリ書込サイクルを開始して受取ったデータを関連のノンキャッシュメモリに転送し、対応するメモリ位置の内容を更新する。もし送信処理ノードが、ヴィクティムブロックコマンドを送った時間と、ターゲット終了メッセージを受取った時間との間で無効化プローブに出会えば、送信ノードはターゲットノード、すなわち第2の処理ノードにメモリキャンセル応答を送り、メモリ書込サイクルのさらなる処理を打切る。これはシステムメモリ帯域幅を節約するという効果をもたらし、ノンキャッシュメモリに書込まれるべきデータが失効している場合、時間がかかるメモリ書込動作を回避し得る。

20

30

【0012】

メモリキャンセル応答は、ヴィクティムブロック書込動作の間のキャッシュコヒーレンスを維持し得るが、特に、ヴィクティムブロックの宛先であるメモリ位置の内容を読み出すための第3の処理ノード（ヴィクティムブロックを送ったソースノード以外のもの）からの読出コマンドの後に、ヴィクティムブロックがターゲットノード（すなわち、第2の処理ノード）に到着する状況において、コヒーレンスを維持し得る。読出コマンドは、そのメモリ位置から読出したデータを変更するという第3の処理ノードの意図を明らかにし得る。したがって、ターゲットノードは応答して、ソースノードを含むシステム内の処理ノードの各々に無効化プローブを伝送し得る。後から到着したヴィクティムブロックは、最新のデータを含み得ず、かつターゲットノードメモリ内の対応するメモリ位置にコミットする必要がないために、ソースノードがターゲット終了応答を受取ったときに、ソースノードはターゲットノードにメモリキャンセル応答を送る。さらに、ターゲット終了応答は無効化プローブの介入の後で受取られるために、ソースノードからのメモリキャンセル応答はこうして処理ノードの間のキャッシュコヒーレンスを維持する助けをする。

40

【0013】

一実施例においては、第1の処理ノードが第2の処理ノードに読出コマンドを送って、第2の処理ノードに関連する指定されたメモリ位置からデータを読み出すと、第2の処理ノードは、応答してシステム内のすべての残りの処理ノードにプローブコマンドを送信する。

50

指定されたメモリ位置のキャッシュされたコピーを有する処理ノードの各々は、そのキャッシュされたデータに関連するキャッシュタグを更新してデータの現在のステータスを反映させる。プローブコマンドを受取った処理ノードの各々は次いで、処理ノードがデータのキャッシュされたコピーを有するかどうかを示すプローブ応答を送る。処理ノードが指定されたメモリ位置のキャッシュされたコピーを有する場合には、その処理ノードからのプローブ応答はキャッシュされたデータの状態、すなわち変更、共用などをさらに含む。

【0014】

プローブコマンドを受取ると、すべての残りのノードは、指定されたメモリ位置のキャッシュされたコピーがもしあれば、上述のようにそのステータスをチェックする。ソースノードとターゲットノード以外の処理ノードが、指定されたメモリ位置のキャッシュされたコピーで、かつ変更された状態のものを見出した場合には、その処理ノードは応答してターゲットノード、すなわち第2の処理ノードにメモリキャンセル応答を送る。このメモリキャンセル応答は、第2の処理ノードにさらなる読出コマンドの処理を打切らせ、かつまだ読出応答を送っていないければ、読出応答の送信を中止させる。それでも他のすべての残りの処理ノードは、それらのプローブ応答を第1の処理ノードに送る。変更されたキャッシュされたデータを有する処理ノードは、その変更されたデータをそれ自体の読出応答を介して第1の処理ノードに送る。プローブ応答と読出応答とを含むメッセージ通信方式はこうして、システムメモリ読出動作の間にキャッシュコヒーレンスを維持する。

【0015】

メモリキャンセル応答はさらに、第2の処理ノードがそれ以前に読出応答を第1の処理ノードに送ったかどうかにかかわらず、第2の処理ノードがターゲット終了応答を第1の処理ノードに送信するようにさせる。第1の処理ノードは、すべての応答、すなわちプローブ応答、ターゲット終了応答、および変更されたキャッシュされたデータを有する処理ノードからの読出応答を待ち、その後で第2の処理ノードにソース終了応答を送ることにより、データ読出サイクルを完了させる。この実施例においては、メモリキャンセル応答は、要求されたデータの変更されたコピーが異なった処理ノードにおいてキャッシュされたときに時間のかかるメモリ読出動作を打切らせることにより、システムメモリ帯域幅を節約し得る。処理ノードとシステムメモリとの間の比較的低速のシステムメモリバスが関与する同様のデータ伝送よりも、高速デュアル単方向リンクを介した2つの処理ノード間のデータ伝送が実質的に速いことが観察されたとき、こうしてデータ転送レイテンシの減少が達成される。

【0016】

以下の図面と併せて、以下の好ましい実施例の詳細な説明を考察することにより、この発明はよりよく理解されるであろう。

【0017】

【発明の実施の形態】

図1は、マルチプロセッシングコンピュータシステム10の一実施例を示す。図1の実施例においては、コンピュータシステム10はいくつかの処理ノード12A、12B、12C、および12Dを含む。処理ノードの各々は、処理ノード12A-12Dにそれぞれ含まれるメモリコントローラ16A-16Dを介して、それぞれのメモリ14A-14Dに結合される。さらに、処理ノード12A-12Dは、インターフェイスロジックとしても知られる、1つ以上のインターフェイスポート18を含んで処理ノード12A-12Dの間で通信し、かつ処理ノードと対応するI/Oブリッジとの間でも通信する。たとえば、処理ノード12Aは、処理ノード12Bと通信するためのインターフェイスロジック18Aと、処理ノード12Cと通信するためのインターフェイスロジック18Bと、さらに別の処理ノード(図示せず)と通信するための第3のインターフェイスロジック18Cを含む。同様に、処理ノード12Bはインターフェイスロジック18D、18E、および18Fを含み、処理ノード12Cはインターフェイスロジック18G、18H、および18Iを含み、処理ノード12Dはインターフェイスロジック18J、18K、および18Lを含む。処理ノード12Dは、インターフェイスロジック18Lを介して結合されてI/O

ブリッジ20と通信する。他の処理ノードは同様の様式で他のI/Oブリッジと通信し得る。I/Oブリッジ20はI/Oバス22に結合される。

【0018】

処理ノード12A-12Dを相互接続するインターフェイス構造は、1組のデュアル単方向リンクを含む。デュアル単方向リンクの各々は、パケットベースの1対の単方向リンクとして実現化されて、コンピュータシステム10内のどの2つの処理ノード間でも高速パケット化情報転送を達成する。単方向リンクの各々は、パイプライン化され分割されたトランザクションによる相互接続として見る事ができる。単方向リンク24の各々は、1組のコヒーレントな単方向ラインを含む。こうして、単方向リンク対の各々は、第1の複数のバイナリパケットを担持する1つの送信バスと、第2の複数のバイナリパケットを担持する1つの受信バスとを含む、と見る事ができる。バイナリパケットの内容は第1に、要求される動作の種類と、動作を開始する処理ノードとに依存する。デュアル単方向リンク構造の一例は、リンク24Aおよびリンク24Bである。単方向ライン24Aを用いてパケットを処理ノード12Aから処理ノード12Bに送信し、ライン24Bを用いてパケットを処理ノード12Bから処理ノード12Aに送信する。ライン24C-24Hの他の組を用いて、図1に示すようにそれらの対応する処理ノードの間のパケットを送信する。

10

【0019】

同様のデュアル単方向リンク構造を用いて、処理ノードとその対応のI/O装置、またはグラフィック装置、もしくは処理ノード12Dに関して示すI/Oブリッジとの間の相互接続を行ない得る。デュアル単方向リンクは、処理ノード間の通信のためにキャッシュコヒーレント様式で動作するか、または処理ノードと外部I/O、またはグラフィック装置、もしくはI/Oブリッジとの間の通信のために、非コヒーレント様式で動作し得る。一方の処理ノードから他方へ送信されるべきパケットは、1つ以上の残りのノードを通過し得ることに留意されたい。たとえば、処理ノード12Aによって処理ノード12Dに送信されるパケットは、図1の構成内の処理ノード12Bまたは処理ノード12Cのいずれをも通過し得る。好適な経路制御アルゴリズムのいずれを用いることもできる。コンピュータシステム10の他の実施例は、図1に示すものよりもより多くの、またはより少ない処理ノードを含み得る。

20

【0020】

処理ノード12A-12Dは、メモリコントローラおよびインターフェイスロジックに加えて、1つ以上のプロセッサコア、内部キャッシュメモリ、バスブリッジ、グラフィックスロジック、バスコントローラ、周辺装置コントローラなどの他の回路素子を含み得る。概略的には、処理ノードは少なくとも1つのプロセッサを含み、選択により、メモリおよび所望の他のロジックと通信するためのメモリコントローラを含む。さらに、処理ノード内の回路素子の各々は、処理ノードによって行なわれる機能に依拠して1つ以上のインターフェイスポートに結合されることが出来る。たとえばある回路素子は、I/Oブリッジを処理ノードに接続するインターフェイスロジックのみを結合し、他の回路素子は2つの処理ノードを接続するインターフェイスロジックのみを結合し得る。他の組合せは、所望のように容易に実現し得る。

30

40

【0021】

メモリ14A-14Dは、いずれかの好適なメモリ装置を含み得る。たとえば、メモリ14A-14Dは、1つ以上のRAMBUS DRAM(RDRAM)、シンクロナスDRAM(SDRAM)、スタティックRAMなどを含み得る。コンピュータシステム10のメモリアドレス空間は、メモリ14A-14Dの間で分割される。処理ノード12A-12Dの各々はメモリマップを含むことができ、該メモリマップを用いて、どのアドレスがどのメモリにマッピングされているかを判断し、よって、ある特定のアドレスに対するメモリ要求がどの処理ノード12A-12Dに経路制御されるべきかを判断する。一実施例においては、コンピュータシステム10内のアドレスに対するコヒーレンシ点は、アドレスに対応するバイトをストアしているメモリに結合された、メモリコントローラ16A-

50

16Dである。言い換えると、メモリコントローラ16A - 16Dは、対応するメモリ14A - 14Dへのメモリアクセスの各々を、キャッシュコヒーレントな様式で起こることを確実にすることを担当している。メモリコントローラ16A - 16Dは、メモリ14A - 14Dにインターフェイスするための制御回路を含み得る。さらに、メモリコントローラ16A - 16Dは、メモリ要求を待ち行列として管理するための、要求キューを含み得る。

【0022】

一般的には、インターフェイスロジック18A - 18Lは、1つの単方向リンクからのパケットを受取り、かつ別の単方向リンクに送信されるべきパケットをバッファするための、さまざまなバッファを含み得る。コンピュータシステム10は、パケットを転送するための好適なフロー制御であればいずれでも用い得る。たとえば一実施例においては、送信インターフェイスロジック18の各々は、送信インターフェイスロジックが接続されたリンクの他端の受信インターフェイスロジック内に、いくつかの種類のバッファのカウントをストアする。インターフェイスロジックは、受信インターフェイスロジックがパケットをストアするフリーのバッファを有さない限り、パケットを送信しない。パケットを次に経路制御することにより受信バッファが解放されると、受信インターフェイスロジックは送信インターフェイスロジックにメッセージを送り、バッファが解放されたことを示す。そのような機構は、「クーポンに基づく」システムと呼べる。

【0023】

次に図2は、処理ノード12Aおよび12Bのブロック図を示し、処理ノード12Aおよび12Bを接続するデュアル単方向リンク構造のより詳細な一実施例を例示する。図2の実施例においては、ライン24A（単方向リンク24A）は、クロックライン24AAと、制御ライン24ABと、コマンド/アドレス/データバス24ACとを含む。同様に、ライン24B（単方向リンク24B）は、クロックライン24BAと、制御ライン24BBと、コマンド/アドレス/データバス24BCとを含む。

【0024】

クロックラインは、対応する制御ラインおよびコマンド/アドレス/データバスに対するサンプルポイントを示すクロック信号を送信する。特定の一実施例においては、データ/制御ビットはクロック信号のエッジの各々（すなわち立上がりエッジおよび立下がりエッジ）で送信される。したがって、クロックサイクルごとに、ラインごとに2つのデータビットを送信し得る。ラインごとに1ビットを送信するために使用される時間は、ここでは「ビット時間」と呼ぶ。上述の実施例は、クロックサイクルごとに2つのビット時間を含む。パケットは2つ以上のビット時間で伝送し得る。コマンド/アドレス/データバスの幅に依拠して、多数のクロックラインを用い得る。たとえば32ビットコマンド/アドレス/データバスに対しては2つのクロックラインを用い得る（コマンド/アドレス/データバスの半分では一方のクロックラインが参照され、残りの半分のコマンド/アドレス/データバスと制御ラインとは他方のクロックラインが参照される）。

【0025】

制御ラインは、コマンド/アドレス/データバスに送信されたデータが、ビット時間の制御パケットか、またはビット時間のデータパケットであるかを示す。制御ラインはアサートされて制御パケットを示し、デアサートされてデータパケットを示す。ある制御パケットは、後にデータパケットが続くことを示す。データパケットは、対応する制御パケットのすぐ後に続き得る。一実施例においては、他の制御パケットがデータパケットの送信に割込むおそれがある。そのような割込みは、データパケットの送信の間に制御ラインをいくつかのビット時間アサートし、かつ制御ラインがアサートされている間にビット時間の制御パケットを送信することにより行なわれる可能性がある。データパケットに割込む制御パケットは、データパケットが後に続くことを示さないおそれがある。

【0026】

コマンド/アドレス/データバスは、データ、コマンド、応答、およびアドレスビットを送信するための1組のラインを含む。一実施例においては、コマンド/アドレス/データ

10

20

30

40

50

バスは、8、16、または32のラインを含み得る。処理ノードまたはI/Oブリッジの各々は、設計選択にしたがってサポートされる数のラインのうちのいずれかを用い得る。他の実施例は、所望の他のサイズのコマンド/アドレス/データバスをサポートし得る。

【0027】

一実施例によると、コマンド/アドレス/データバスラインおよびクロックラインは、反転データを担持し得る（すなわち、論理1はライン上の低電圧として表わされ、論理0が高電圧として表わされる）。これに代えて、これらのラインは非反転データを担持してもよい（論理1はライン上の高電圧として表わされ、論理0は低電圧として表わされる）。好適な正および負論理の組合せもまた実現化し得る。

【0028】

図3から図7は、コンピュータシステム10の一実施例に従った、キャッシュコヒーレントな通信（すなわち処理ノード間の通信）に用いられる例示的なパケットを示す。図3から図6は制御パケットを示し、図7はデータパケットを示す。他の実施例は異なったパケット定義を用い得る。制御パケットおよびデータパケットは集合的にバイナリパケットとも呼ぶ。パケットの各々は、「ビット時間」の見出しの下に列挙される一連のビット時間で示される。パケットのビット時間は、リストされたビット時間順序に従って送信される。図3から図7は、8ビットコマンド/アドレス/データバス実現化のためのパケットを示す。したがって、（7から0まで番号が付与された）8ビットの制御情報またはデータ情報は、ビット時間の各々の間に8ビットコマンド/アドレス/データバス上を送信される。図中、いずれの値も付与されていないビットは、所与のパケットのために予約されているか、またはパケット特定情報を伝送するために用いられるかのいずれかであり得る。

【0029】

図3は情報パケット（infoパケット）30を示す。情報パケット30は、8ビットリンク上の2つのビット時間を含む。この実施例においては、コマンド符号化はビット時間1の間に送信され、かつコマンドフィールドCMD[5:0]で示す、6ビットを含む。例示的なコマンドフィールド符号化を図8に示す。図4、図5、図6に示す他方の制御パケットの各々は、ビット時間1の間に同じビット位置においてコマンド符号化を含む。メッセージがメモリアドレスを含まないときに、情報パケット30を用いてこのメッセージを処理ノード間で送信し得る。

【0030】

図4はアドレスパケット（addressパケット）32を示す。アドレスパケット32は、8ビットリンク上の8つのビット時間を含む。コマンド符号化は、DestNodeフィールドで示す宛先ノード番号の一部と併せて、ビット時間1の間に送信される。宛先ノード番号の残りとソースノード番号（SrcNode）とは、ビット時間2の間に送信される。ノード番号はコンピュータシステム10内の処理ノード12A-12Dのうちの1つを明確に識別し、かつ用いられてパケットをコンピュータシステム10を介して経路制御する。さらに、パケットのソースは、ビット時間2および3の間に送信されるソースタグ（SrcTag）を割当て得る。ソースタグは、ソースノードによって開始される特定のトランザクションに対応するパケットを識別する（すなわち、特定のトランザクションに対応するパケットの各々は、同一のソースタグを含む）。こうして、たとえばSrcTagフィールドが7ビット長さであれば、対応するソースノードはシステム内で進行する間に最大128（ 2^7 ）の異なったトランザクションを有し得る。システム内の他のノードからの応答は、応答内のSrcTagフィールドを介して対応のトランザクションと関連付けられる。ビット時間4から8までを用いて、アドレスフィールドAddr[39:0]で示すトランザクションによって影響されたメモリアドレスを送信する。アドレスパケット32を用いて、トランザクション、たとえば読出または書込トランザクションを開始し得る。

【0031】

図5は、応答パケット（responseパケット）34を示す。応答パケット34は、コマンド符号化、宛先ノード番号、ソースノード番号、およびアドレスパケット32と同様のソースタグを含む。SrcNode（ソースノード）フィールドは好ましくは、応答パケットの生成

10

20

30

40

50

を促すトランザクションを発信したノードを識別する。一方、DestNode（宛先ノード）フィールドは、応答パケットの最終的なレシーバである処理ノードを、すなわちソースノードまたはターゲットノード（後に説明）を識別する。さまざまな種類の応答パケットが付加的な情報を含み得る。たとえば、図 1 1 A を参照して後に説明する読出応答パケットは、以下のデータパケットで提供される読出データの量を示し得る。後に図 1 2 を参照して説明するプローブ応答は、要求されたキャッシュブロックに対してヒットが検出されたかどうかを示し得る。一般的に、応答パケット 3 4 は、トランザクションを行なう間にアドレスの送信を必要としないコマンドに対して用いられる。さらに、応答パケット 3 4 を用いて肯定応答パケットを送信してトランザクションを終了させることができる。

【 0 0 3 2 】

図 6 は、コマンドパケット（command パケット）3 6 の例を示す。上述のように、単方向リンクの各々はパイプライン化され、分割されたトランザクション相互接続であって、トランザクションはソースノードによってタグ付けされ、応答は任意の所与の時間にも、パケットの経路制御に依存して順不同でソースノードに戻ることができる。ソースノードは、コマンドパケットを送信してトランザクションを開始する。ソースノードはアドレスマッピングテーブルを含み、ターゲットノード番号（TgtNode フィールド）をコマンドパケットに入れて、コマンドパケット 3 6 の宛先である処理ノードを識別する。コマンドパケット 3 6 は、C M D フィールド、SrcNode フィールド、SrcTag フィールド、および Addr フィールドを有するが、これらはアドレスパケット 3 2（図 4）を参照に説明され示されたものと同様である。

【 0 0 3 3 】

コマンドパケット 3 6 の 1 つの際立った特徴は、Count フィールドの存在である。キャッシュ不可能な読出または書込動作においては、データのサイズはキャッシュブロックのサイズよりも小さくあり得る。こうして、たとえば、キャッシュ不可能な読出動作は、システムメモリまたは I / O 装置からのちょうど 1 バイトまたは 1 クワッドワード（6 4 ビット長さ）だけのデータを必要とし得る。この種類のサイズ指定された読出または書込動作は、Count フィールドの助けによって容易となる。この例においては、Count フィールドは 3 ビット長さで示す。したがって、所与のサイズ指定されたデータ（バイト、クワッドワードなど）は、最高 8 回まで送信されることができる。たとえば、8 ビットリンクにおいては、Count フィールドの値が 0（バイナリ 0 0 0）である場合、コマンドパケット 3 6 は 1 つのビット時間でのちょうど 1 バイトだけのデータの転送を示す。一方、Count フィールドの値が 7（バイナリ 1 1 1）である場合、クワッドワード、すなわち 8 バイトが、合計で 8 ビット時間の間に伝送されることができる。C M D フィールドは、いつキャッシュブロックが伝送されたのかを識別し得る。この場合、Count フィールドは固定値を有し、キャッシュブロックが 6 4 バイトサイズである場合 7 であるが、これは 8 クワッドワードがキャッシュブロックを読出または書込するために伝送されなければならないためである。8 ビットワイドの単方向リンクの場合においては、6 4 のビット時間にわたる 8 つの完全なデータパケット（図 7）の伝送を必要とし得る。好ましくは、データパケット（図 7 を参照して後に説明）は、書込コマンドパケットまたは読出応答パケット（後に説明）の直後に続き、データバイトはアドレスの昇順で転送されることができる。単一のバイトまたはクワッドワードのデータ転送は、自然に整地されたそれぞれ 8 または 6 4 バイト境界をまたがらない。

【 0 0 3 4 】

図 7 は、データパケット（data パケット）3 8 を示す。データパケット 3 8 は、図 7 の実施例において、8 ビットリンク上の 8 つのビット時間を含む。データパケット 3 8 は、6 4 バイトのキャッシュブロックを含み得るが、この場合キャッシュブロック転送を完了させるために（8 ビットリンク上の）6 4 のビット時間がかかるであろう。他の実施例では、キャッシュブロックのサイズを所望により別に定義し得る。さらに、コマンドパケット 3 6（図 6）を参照して上に説明したように、キャッシュ不可能な読出および書込に対するキャッシュブロックサイズよりも小さなサイズでデータを送信し得る。キャッシュブ

10

20

30

40

50

ックサイズより小さなデータを送信するためのデータパケットは、より少ないビット時間しか必要としない。

【 0 0 3 5 】

図 3 から図 7 は、8 ビットリンクのためのパケットを示す。1 6 および 3 2 ビットリンクのためのパケットは、図 3 から図 7 に示す連続的なビット時間を連結することにより形成し得る。たとえば、1 6 ビットリンク上のパケットのビット時間 1 は、8 ビットリンク上のビット時間 1 および 2 の間に送信される情報を含み得る。同様に、3 2 ビットリンク上のパケットのビット時間 1 は、8 ビットリンク上のビット時間 1 から 4 までの間に送信される情報を含み得る。以下の式 (1) および式 (2) は、8 ビットリンクに対するビット時間における、1 6 ビットリンクのビット時間 1 および 3 2 ビットリンクのビット時間 1 の構成を示す。

【 0 0 3 6 】

【 数 1 】

$$BT1_{16}[15:0] = BT2_8[7:0] \parallel BT1_8[7:0] \quad (1)$$

$$BT1_{32}[31:0] = BT4_8[7:0] \parallel BT3_8[7:0] \parallel BT2_8[7:0] \parallel BT1_8[7:0] \quad (2)$$

【 0 0 3 7 】

図 8 は、コンピュータシステム 1 0 内のデュアル単方向リンク構造の 1 つの例示的な実施例に対して用いられるコマンドを示すテーブル 4 0 を示す。テーブル 4 0 は、コマンドの各々に割当てられたコマンド符号化 (C M D フィールド) を示すコマンド符号化列、コマンドの名前を示すコマンド列、およびどのコマンドパケット 3 0 - 3 8 (図 3 から図 7) がそのコマンドに対して用いられるかを示すパケットタイプ列を含む。図 8 におけるコマンドのいくつかに対する簡単な機能の説明を以下に示す。

【 0 0 3 8 】

読出トランザクションは、Rd (Sized) , RdBlk , RdBlkSまたはRdBlkModのコマンドのうち、1 つを用いて開始される。サイズ指定された読出コマンドである Rd (Sized) は、キャッシュ不可能な読出のために、またはサイズの合ったキャッシュブロック以外のデータの読出のために用いられる。読出されるべきデータ量は、Rd (Sized) コマンドパケット内に符号化される。キャッシュブロックの読出には、以下の場合以外に RdBlk コマンドを用いることができる。すなわち、(i) キャッシュブロックの書込可能なコピーを所望する場合。この場合は RdBlkMod コマンドを用い得る。または (ii) キャッシュブロックのコピーを所望するが、ブロックを変更する意図があるとは分らない場合。RdBlkS コマンドを用いて、ある種のコヒーレントな方式 (たとえばディレクトリに基づくコヒーレントな方式) をより効率化できる。RdBlkS コマンドに応答して、ターゲットノードはキャッシュブロックをソースノードに共用状態で返し得る。一般的に、適切な読出コマンドはソースノードから送信されて、ソースノードから要求されたキャッシュブロックに対応するメモリを有するターゲットノードへの読出トランザクションを開始する。

【 0 0 3 9 】

ソースノードにストアされた書込不可能または読出専用状態のキャッシュブロックへの書込許可を得るために、ソースノードは ChangeToDirty パケットを送信し得る。ChangeToDirty コマンドによって開始されるトランザクションは、ターゲットノードがデータを返さないという点を除いて、読出と同様に動作し得る。もしソースノードがキャッシュブロック全体を更新する意図があるのであれば、ValidateBlk コマンドを用いて、ソースノードにストアされていないキャッシュブロックへの書込許可を得ることができる。そのようなトランザクションに対してはソースノードへデータは転送されないが、それ以外では読出トランザクションと同様に動作する。好ましくは、ValidateBlk および ChangeToDirty コマンドは、メモリのみに向けられ、かつコヒーレントなノードによってのみ生成されることができる。

10

20

30

40

50

【 0 0 4 0 】

InterruptBroadcast、InterruptTarget、およびIntrResponseパケットを用いて、それぞれ割込をブロードキャストし、特定のターゲットノードに割込を送り、かつ割込に応答し得る。CleanVicBlkコマンドを用いて、（たとえば、ディレクトリに基づくコヒーレントな方式のために）クリーンな状態のキャッシュブロック（ヴィクティムブロック）がノードから捨てられたことを伝えることができる。TgtStartコマンドはターゲットによって用いられて、（たとえば、後のトランザクションの順序付けのために）トランザクションが開始したことを示す。エラーコマンドを用いて、エラー表示を送信する。

【 0 0 4 1 】

図 9、図 1 3 および図 1 4 に、コンピュータシステム 1 0 内の処理ノードが指定されたメモリ位置の読出を試みる時のパケットのフローのいくつかの例を示す。指定されたシステムメモリ位置または対応のシステムメモリ位置は、例示のためにのみ、ターゲット処理ノード 7 2 に関連のシステムメモリ 4 2 1 内にあると想定する。システムメモリ 4 2 1 は、ターゲット処理ノード 7 2 の一部であるか、またはここに示すようにターゲットノード 7 2 の外部にあってもよい。さらに、メモリ読出トランザクションの間に、指定されたメモリ位置のコピーが既にターゲットノード 7 2 の内部または外部キャッシュメモリに存在する可能性がある。いずれにしても、ソースノード 7 0 が関連の指定されたメモリ位置を読出すために読出コマンドをターゲットノード 7 2 に送信するときはいつでも、パケットのフローは同じままである。いずれの処理ノード 1 2 A - 1 2 D（図 1）もソースノードまたはターゲットノードとして機能し得ることに留意されたい。ソースノードでもターゲットノードでもないノードは残りのノードと呼ばれるが、ここではノード 7 4 および 7 6 である。図 9、図 1 3、図 1 4 において、理解を助けるためにのみ、同じ番号を用いてソースノード、ターゲットノードおよび残りのノードを識別する。これは図 9 におけるソースノード 7 0 が図 1 3 における同じソースノードであることを表示しない。

【 0 0 4 2 】

上述のように、図 1 におけるいずれの処理ノードも、特定のトランザクションに依拠してソースノード、ターゲットノードまたは残りのノードとして機能し得る。図 9、図 1 3 および図 1 4 の構成は例示のためにのみ示し、これらは処理ノード 1 2 A - 1 2 D の間の同様の実際の接続を示すものではない。すなわち、残りのノード、たとえばノード 7 6、またはターゲットノード 7 2 は、ソースノード 7 0 に直接接続されないかもしれない。したがって、付加的なパケット経路制御が生じ得る。さらに、図 9、図 1 3 および図 1 4 の構成は、図 1 における回路トポロジを参照して説明される。2 つ以上の処理ノードの間の他の相互接続が企図可能であり、これらのさまざまな相互接続で図 9、図 1 3、および図 1 4 のパケット転送方式を容易に実現化し得ることが理解される。矢印は、従属性と、矢印によって結合されるそれぞれのノードの間で送信されるべきパケットとを示す。一般的に外部に向かう矢印は、対応の incoming 従属関係のすべてが発生するまで、生じない。これを図 9、図 1 3、および図 1 4 に示す動作を参照して以下に詳述する。

【 0 0 4 3 】

図 9 を参照して、読出トランザクションの間のパケット 4 2 のフローを示すが、該読出トランザクションは上述のように Rd(Sized) またはブロック読出（RdBlk、RdBlkS、RdBlkMod）である。ソースノード 7 0 内のプロセッサ（図示せず）は、適切な読出コマンドをターゲットノード 7 2 内のメモリコントローラ（図示せず）に送る。典型的なコマンドパケットは、図 6 を参照して既に説明した。ソースプロセッサから読出コマンドを受取ると、応答してターゲットメモリコントローラは、以下の 2 つの動作を行なう。（1）RdResponse（読出応答）パケットをメモリ 4 2 1 から要求されたデータと併せてソースノード 7 0 に送る。（2）コンピュータシステム 1 0 内のすべての処理ノードに Probe / Src コマンドをブロードキャストする。一般的には、Probe / Src コマンド（より簡単には、プローブコマンド）は、キャッシュブロックがノード内に含まれているかどうか判断するためのそのノードへの要求であり、かつもしキャッシュブロックがそのノードにストアされていれば、そのノードが取るべき動作を表示する。一実施例においては、パケットが 1 つ以上の宛先

に対してブロードキャストされると、パケットを最初に受信する受信ノードのルータは、そのノードでパケットを終了させ、隣接する処理ノードにそのパケットのコピーを再生成し送信し得る。

【 0 0 4 4 】

上述のこれら 2 つの動作の正確な実行の順序は、ターゲットノード 7 2 内のさまざまな内部バッファにおける未完了の動作のステータスに依拠し得る。好ましくは、コンピュータシステム 1 0 内の処理ノードの各々は、コマンドパケット、さまざまなコマンドパケットに関連のデータパケット（たとえば、メモリ書込コマンド）、プローブ、応答パケット（たとえば、ProbeResp、SrcDone、TgtDone、MemCancel）、および読出応答（RdResponseパケットおよびその関連のデータパケットの両方を含む）をストアするためのいくつかのバッファを含む。データバッファの各々は、たとえば 6 4 バイトサイズのキャッシュブロックのための記憶装置を含み得る。これに代えて、設計要件に基づいて他の便利な記憶容量のいずれをも実現化し得る。

【 0 0 4 5 】

上述のバッファを用いた 2 つの処理ノード間のバイナリパケットのフローは、上述の「クーポンに基づく」システムを実現化することにより、制御し得る。この実現化においては、送信ノードは、受信ノードの各種のバッファに対するカウンタを含み得る。システムリセットの際に送信ノードはそのカウンタをクリアすることができ、リセット信号がデアサートされたときには、受信ノードは情報パケットを（C M D フィールドがNopコマンドを識別する、図 3 に示すものと同様のフォーマットによって）送信ノードに送り、それが各種の利用可能なバッファをいくつか有するかを示す。送信ノードがパケットを受信ノードに送ると、これは関連のカウンタをデクリメントし、特定のカウンタが値 0 に到達すると、送信ノードプロセッサは関連のバッファへのパケットの送信を停止する。レシーバがバッファを解放すると、これは別の情報パケットをトランスミッタに送り、トランスミッタは関連のカウンタをインクリメントする。トランスミッタは、レシーバがコマンドバッファおよびデータバッファの両方を利用可能にしない限り、メモリ書込動作を開始し得ない。

【 0 0 4 6 】

再び図 9 に戻ると、ターゲットノード 7 2 内のメモリコントローラは、Probe / Srcコマンドをシステム内の他のノードに送信して、これらのノード内のキャッシュブロックの状態を変化させることと、キャッシュブロックの更新されたコピーを有するノードに、キャッシュブロックをソースノードに送らせることとにより、コヒーレンシを維持する。方式は、プローブコマンド内で受信ノードを識別する表示を用いてプローブ応答を受信する。ここで、Probe / Srcコマンド（プローブコマンド）は、残りのノード 7 4、7 6 のそれぞれに、ProbeResp（プローブ応答）をソースノードに送信させる。プローブ応答は、動作が起こったことを示し、かつもしキャッシュブロックがノードによって変更されていればデータの送信を含み得る。もしプローブされたノードが読出データの更新されたコピー（すなわちダーティデータ）を有していれば、図 1 3 に関して後に説明するように、ノードはRdResponse（読出応答）パケットとダーティデータとを送信する。Probe / Srcコマンドは、（ターゲットノード 7 2 を含む）所与の処理ノード内のキャッシュコントローラによって受取られ、ProbeRespおよびRdResponseは、そのキャッシュコントローラによって生成されることができる。一般的には、関連のキャッシュを有する処理ノード内のキャッシュコントローラは、Probe / Srcコマンドに応答してプローブ応答パケットを生成し得る。一実施例においては、処理ノードがキャッシュを有さないとき、その処理ノードはプローブ応答パケットを生成し得ない。

【 0 0 4 7 】

一旦（残りのノード 7 4 および 7 6 からの）プローブ応答と、（ターゲットノード 7 2 からの）要求されたデータを備えたRdResponseとがソースノードにおいて受取られると、ソースノードプロセッサはSrcDone（ソース終了）応答パケットをトランザクション終了の肯定応答としてターゲットノードメモリコントローラ（図示せず）に送信する。読出動作の各々の間の処理ノード間のコヒーレンシを維持するために、ソースノードは（残りのノ

10

20

30

40

50

ードからの)すべてのプローブ応答をも受取るまで、ターゲットノード72からRdResponseを介して受取ったデータを使用し得ない。ターゲットノードがSrcDone応答を受取ったとき、これは(ソースノード70から受取った)読出コマンドをそのコマンドバッファキューから取除き、次いでこれは同様に指定されたメモリアドレスに対してコマンドへの応答を開始し得る。

【0048】

送られたコマンドに依拠してプローブ応答を異なった受信ノードへ経路制御する柔軟性を与えることにより、コヒーレンシの維持を比較的効率的な態様で行ない得る(たとえば、処理ノード間で最も少ない数のパケット送信を用いる)一方で、さらにコヒーレンシが維持されることを確実にする。たとえば、トランザクションのターゲットまたはソースがトランザクションに対応するプローブ応答を受取るべきであることを示すプローブコマンドをも含み得る。プローブコマンドは、トランザクションのソースを読出トランザクションに対する受信ノードとして特定し得る(それによりダーティデータをストアしていたノードからソースノードへダーティデータが引き渡される)。一方で、(トランザクションのターゲットノードのメモリ内でデータが更新される)書込トランザクションに対しては、プローブコマンドはトランザクションのターゲットを受信ノードとして特定し得る。こうして、ターゲットは書込データをいつメモリにコミットするか判断し、かつ書込データとマージされるべきダーティデータのいずれかを受取り得る。

【0049】

図10から図12は、プローブコマンド、読出応答およびプローブ応答パケットのそれぞれの一実施例を示す。図10Aのプローブコマンドパケット44は、図6に示す一般的なコマンドパケットとはやや異なる。CMDフィールドは、受信ノードにその応答をソースノード70へ送信することを要求するProbe/Srcコマンドとして、プローブを識別する。上述のように特定の他のトランザクションにおいてはターゲットノード72はプローブコマンドに対する応答の受信側であり得るが、CMDフィールドはまたこれらの場合でもそのように示すであろう。さらに、関連する経路制御に依拠して、ソースノード70またはターゲットノード72のいずれか、またはこれらの両方が、システム内の他の残りのノードよりも前かまたは同時にProbe/Srcコマンドを受取る可能性もある。プローブコマンドのSrcNodeおよびTgtNodeフィールドは、ソースノードとターゲットノードとをそれぞれ識別し、ソースノードキャッシュコントローラがプローブコマンドに応答することを防ぐであろう。SrcTagフィールドは、図4を参照に先に説明したものと同様に機能する。DM(データ移動)ビットは、このプローブコマンドに応答してデータ移動が要求されているかどうかを示す。たとえば、DMビットがクリアされていれば、それはいずれのデータ移動もないことを示す。一方で、もしDMビットがセットされていれば、プローブコマンドが残りのノード74または76のうちの1つの中の、内部(外部)キャッシュ内のダーティブロックまたは共用ノダーティブロックをヒットした場合に、データ移動が要求される。

【0050】

上述のように、ソースノードからの読出コマンドは、サイズ指定された読出コマンド[Rd(sized)]またはブロック読出コマンド[RdBlk、RdBlkS、またはRdPlkMmd]であり得る。どちらの種類の読出コマンドも好ましくはデータ移動を要求し、よってDMビットはターゲットノードのメモリコントローラによってセットされてデータ移動要求を示し得る。異なった実施例においては、DMビットはクリアされている場合に、データ移動を示し、DMビットはセットされている場合に、いずれのデータ移動もないことを示す。

【0051】

NextStateフィールド46(図10B)は、Probeビットがあった場合、すなわち、1つ以上の残りのノードが、プローブコマンドAddrフィールドによって識別される指定されたメモリ位置のキャッシュコピーを有する場合に、起こるべきステートトランザクションを示す2ビットフィールドである。図10BにNextStateフィールド46に対する1つの例示的な符号化を示す。ブロック読出コマンドの間では、NextStateフィールドは1であり、よってメモリデータのキャッシュコピーを有する残りのノードは、Probe/Srcコマンドを

10

20

30

40

50

受信するとそのコピーを共用としてマークする。一方で、サイズ指定された読出コマンドの間では、NextStateフィールドは0であり、よって、いずれの残りのノードも、メモリ421からのデータのキャッシュコピーを有する場合であっても、対応のキャッシュタグを変える必要はない。特定の他のターゲットメモリトランザクション（たとえば特定の書込動作）においても、対応する残りのノード内のキャッシュされたデータを、値2でNextStateフィールド46によって示されるように、無効としてマークすることが望ましいであろう。

【0052】

こうしてプローブコマンドは、このNextStateフィールドを介してメモリ読出動作の間のシステム処理ノード間のキャッシュコピーレンシを維持し得る。ターゲットノードキャッシュコントローラは、ターゲットノードメモリコントローラによるProbe/Srcコマンドブロードキャストの受信の際、およびターゲットノード（内部または外部）キャッシュメモリ内の要求されたデータの発見の際に、プローブ応答パケットを読出応答パケットと併せて送信し得る。後に説明するように、ソースノードは、RdResponseおよびProbeRespパケットによって供給される情報によって要求されたデータに関連のキャッシュタグを更新する。このようにしてソースノードは、（対応するキャッシュタグを介して）これが要求されたデータの排他的または共用コピーを有するかどうかを表示し得る。ターゲットノードキャッシュコントローラからのプローブ応答パケットは、たとえばターゲットノードだけが要求されたデータのコピーをそのキャッシュ内に有し、他のいずれの残りのノードも要求されたデータのキャッシュコピーを有さない状況において、助けになり得る。しかしながら、ターゲットノードは、ターゲットノードがソースによって要求されたデータをそのキャッシュ内に有するとき、そのキャッシュ状態を自動的に更新するよう構成されてもよく、したがってこれはターゲットノードキャッシュからデータをソースへ送る。

【0053】

図11Aを参照すると、RdResponseパケット48に対する例示的な符号化を示す。ターゲットノード72内のメモリコントローラ（図示せず）は、サイズ指定された読出コマンドまたはブロック読出コマンドのいずれであっても、読出コマンドの各々に応答してRdResponseをソースノード70に送るよう構成されることができる。これに代えて上述のように、ターゲットノードキャッシュコントローラ（図示せず）は、要求されたデータがターゲットノード内にキャッシュされている場合に、適切な読出応答パケットを送るよう構成されてもよい。典型的には、RdResponseパケット48の後には要求されたデータを含むデータパケット38（図7）が続く。サイズ指定された読出動作に対するデータパケットは、最も低いアドレスのデータが最初に返され、残りのアドレスのデータが昇順に返されるよう構成されてもよい。しかしながら、キャッシュブロック読出に対するデータパケットは、要求されたクワッドワード（64ビット）が最初に返され、残りのキャッシュブロックはインタリーブラッピングを用いて返されるよう構成されてもよい。

【0054】

RdResponseパケット48内のCountフィールドは、読出トランザクションを開始する読出コマンド内のCountフィールド（たとえば図6を参照）と同一である。Typeフィールドは元の読出要求のサイズを符号化し、かつCountフィールドと併せて、データパケットのサイズの合計を表示する。Typeフィールドはバイナリ値0または1のいずれかをとり得る。一実施例においては、Typeフィールドは0であるとき、バイトサイズのデータが転送されるべきであることを示し得る。Typeフィールドが1であるとき、クワッドワード（64ビット）のデータが転送されるべきであることを示し得る。一方、Countフィールドは、Typeフィールドによって示されるそのサイズのデータが、リンクをわたって何回転送されるべきであることを示し得る。こうして、CountフィールドとTypeフィールドとは組合わされて、転送されるべきデータの合計のサイズを判断し得る。たとえば、8ビット単方向リンクをわたるサイズ指定された読出動作の間、ダブルワードの転送のためにはTypeフィールドは0であって、Countフィールドは3でなければならない[バイナリでは011]。

【0055】

10

20

30

40

50

RdResponseパケット 4 8 内のRespNodeフィールドは、読出応答パケットが向けられるべきノードを識別する。SrcNodeフィールドは、トランザクションを開始したノード、すなわちソースノード 7 0 を識別する。読出動作の間、RespNodeおよびSrcNodeフィールドは、同一のノード、すなわちソースノード 7 0 を識別するであろう。図 1 3 を参照して後に説明するように、キャッシュ内に（ターゲットメモリ 4 2 1 内の）アドレス指定されたメモリ位置のダーティコピーを有する残りのノードのうちの 1 つによって、RdResponseが生成されるであろう。データ移動を要求するプローブに回答してノードによって読出応答 4 8 が生成されたことを示すために、プローブビットがセットされることができる。クリアされたProbeビットは、メモリコントローラ（図示せず）、またはターゲットノード 7 2 のキャッシュコントローラ（図示せず）のいずれからRdResponse 4 8 が来たのかを示し得る。

10

【 0 0 5 6 】

Tgtビットは、CMD [5 : 0] フィールド内のビット位置 [0] のビットである。一実施例においては、Tgtビットはセットされている場合に、RdResponse 4 8 がターゲットノード 7 2 内のメモリコントローラ（図示せず）に（たとえば、ある書込トランザクションの間に）宛先決めされていることを示し得る。一方で、Tgtビットはクリアされている場合に、RdResponse 4 8 がソースノード 7 0 に宛先決めされていることを示し得る。こうしてTgtビットは、どのようにデータフローがノード内で内部的に管理されているかを識別する。ある実施例においては、Tgtビットを省いてもよい。

【 0 0 5 7 】

20

図 1 1 B 内のテーブル 5 0 は、Probeビット、Tgtビット、Typeフィールド、およびCountフィールド間の関係の 1 つの例を示す。ここで示すように、RdResponse 4 8 がターゲットノード 7 2 のキャッシュコントローラ（図示せず）またはメモリコントローラ（図示せず）から来ている場合はいつでも、Probeビットはクリアされている。一実施例においては、ターゲットノード 7 2 は、（たとえばサイズ指定された読出動作の間）キャッシュブロックサイズよりも小さなデータを供給し得る。TypeフィールドとCountフィールドとは共にソースノード 7 0 に転送されるべきデータのサイズを特定し得る。後に説明するように、残りのノード（ノード 7 4 またはノード 7 6 ）のうちの 1 つがソースノード 7 0 にRdResponseパケットを送るとき、転送されることのできるサイズのデータはキャッシュブロックのみである。この状況において（キャッシュブロックサイズが 6 4 バイトであると想定すると）6 4 バイトのデータ転送を達成するためには、Countフィールドは 7（バイナリ 1 1 1）であり、Typeフィールドは 1 でなくてはならない。

30

【 0 0 5 8 】

図 1 2 を参照して、ProbeRespパケット 5 2 の例を示す。一般的に、関連のキャッシュメモリを有する処理ノード（1 つ以上の残りのノードまたはターゲットノード 7 2）は、MissまたはHitNotDirtyを示してProbeRespパケットをソースノード 7 0 に向けることにより、Probe / Srcコマンドに回答する。しかしながら、もし回答するノードが要求されたデータの変更されたキャッシュされたコピーを有すれば、これは代わりに、後に説明するようにRdResponseを送信する。CMDフィールド、RespNodeフィールド、SrcNodeフィールド、およびSrcTagフィールドについては、既に 1 つ以上の制御パケットを参照して先に説明した。一実施例においては、ヒットビットがセットされていると、回答するノードがアドレス指定されたメモリ位置の変更されていないキャッシュされたコピーを有することを（ソース処理ノード 7 2 に）示す。別の実施例においては、クリアされたヒットビットが同様の表示を示し得る。こうして、ソースノード 7 0 は、ターゲットノード 7 2 から受取ったデータのブロックを（そのキャッシュ内に）どのようにマークするかについての必要な情報を得る。たとえば、もし残りのノード 7 4 または 7 6 のうちの 1 つが変更されていない（すなわちクリーンである）アドレス指定されたメモリ位置のコピーを有すれば、ソースノード 7 0 はターゲットメモリコントローラ（図示せず）から受取ったデータブロックをクリーン / 共用であるとマークするであろう。一方で、もしこれがサイズ指定された読出動作であれば、ソースノード 7 0 は、読出したデータがキャッシュブロックよりもサイ

40

50

ズが小さいために、受取ったデータに関連するそのキャッシュタグを変える必要はない。これは残りのノードを参照した先の説明（図10B）と極めて似ている。

【0059】

図13に、残りのノードの1つ（ここではノード76）がそのキャッシュ内にターゲットメモリ位置の変更されたコピー（すなわちダーティコピー）を有する場合の packets のフローの例、すなわち配置54を示す。上述のように、ターゲットノードメモリコントローラ（図示せず）は、ソースノード70から読出コマンドを受取ると、Probe / SrcCommand（プローブコマンド）とRdResponseとを送る。ここで、ターゲットノード72は関連のキャッシュメモリを有すると想定し、よって、ターゲットノードキャッシュコントローラ（図示せず）は上述のようにソースノード70にプローブ応答を送る。ターゲットノード72もまた要求されたデータのキャッシュされたコピーを有する場合には、ターゲットノードキャッシュコントローラ（図示せず）は上述のように、要求されたデータと併せて読出応答 packet をも送る。関連のキャッシュがない場合には、ターゲットノード72はプローブ応答 packet を送り得ない。

10

【0060】

プローブコマンド packet および読出応答 packet の一実施例は、それぞれ図10Aおよび図11Aに関して先に説明した。しかしながら、図13の実施例においては、応答ノード76はプローブコマンドに回答して、そのキャッシュコントローラを介して2つの packet を送るよう構成される。すなわち、RdResp packet をソースノード70内のプロセッサに送り、MemCancel 応答をターゲットノードメモリコントローラ（図示せず）に送る。残りのノード76からの読出応答の後に、プローブコマンド packet 内のDM（データ移動）ビット（図10A）から要求されて、変更されたキャッシュブロックを含むデータ packet が続く。図11Aを参照して先に説明したように、非ターゲットノードからのRdResponseは、そのProbeビットをセットにしてデータブロックのソースがターゲットノード72ではないことを示し得る。応答ノード76からのこのRdResponse packet を介して、ソースノード70は受取ったデータのキャッシュブロックの状態を（その内部キャッシュ内に）変更 / 共有としてマークするための表示を得る。

20

【0061】

残りのノード76からのRdResponse packet は、読出コマンドがサイズ指定された読出トランザクションを識別した場合でも、対応するキャッシュブロックの全体を（変更された状態で）含む。異なった実施例においては、応答非ターゲットノード（ここではノード76）は、要求されたデータのみを直接ソースノードに送るよう構成されてもよい。この実施例においては、ソースノードに転送されるべきデータのサイズは、プローブコマンドの一部として符号化してもよい。さらなる別の実施例においては、応答ノード76は要求されたデータのみをターゲットノード72内のメモリコントローラ（図示せず）に送ってもよく、その後で、ターゲットノードメモリコントローラはデータをソースノード70に返す。

30

【0062】

応答ノード76からのMemCancel（メモリキャンセル）応答は、ターゲット処理ノード72のメモリコントローラにソースノード70からの読出コマンドのさらなる処理を打ち切らせる。言い換えると、MemCancel 応答は、ターゲットノードメモリコントローラからのRdResponse packet（および要求されたデータ）の送信をキャンセルし、かつターゲットノードメモリコントローラによる先行のメモリ読出サイクルさえもキャンセルする効果を有するが、該メモリ読出サイクルとは、もしターゲットノード読出応答バッファからのRdResponse packet の解放の前、またはメモリ読出サイクルの完了の前に、ターゲットノードメモリコントローラがMemCancel 応答をそれぞれ受取っていればソース70からの読出コマンドに回答して開始されているであろうものである。こうしてMemCancel 応答は、2つの主要な目的を達成する。（1）システムメモリ（たとえば、メモリ421）が失効したデータを有する場合、比較的長いメモリアクセスを可能な限りなくすことにより、システムメモリバス帯域幅を節約する。これはまたコヒーレントなリンク上の不必要なデータ

40

50

転送をも減じる。(2) 処理ノード間で最新のキャッシュデータの転送を可能にすることにより、マルチプロセッシングコンピュータシステムにおいてさまざまな処理ノード間のキャッシュコピーレンシを維持する。

【0063】

図1の回路構成に含まれる経路制御のために、応答ノード76からのMemCancel応答パケットのターゲットノード72への到達は、ターゲットノードメモリコントローラの読出応答パケットの送信または比較的長いメモリ読出サイクルを打ち切らせるのに間に合わないおそれがあることに留意されたい。そのような状況においては、ターゲット処理ノード72は、読出応答送信またはシステムメモリ読出サイクルに遅すぎる場合には、遅れて到着したMemCancel応答を単に無視し得る。トランザクションが打ち切られる正確な時点は、回路構成、実現化された経路制御、オペレーティングソフトウェア、さまざまな処理ノードを構成するハードウェアなどに依存し得る。ソースノードが、ターゲットノードメモリコントローラからRdResponseを受取るとき、これはこのRdResponse(およびその関連のデータパケット)を単に無視し、代わりに、残りのノード76からRdResponseパケットと併せて供給されたキャッシュブロックからの要求されたデータを受取る。

10

【0064】

MemCancel応答を受取ると、ターゲットノードメモリコントローラは、TgtDone(ターゲット終了)応答をソース処理ノード70に送信する。TgtDone応答は、それ以前にターゲットノードがRdResponseパケット(および要求されたデータ)をソースノード70に送ったかどうかにかかわらず、送信される。もしターゲットノードメモリコントローラがそれ以前にRdResponseパケットを送っていなければ、これはRdResponseパケット(および要求されたデータ)の送信をキャンセルし、代わりに、TgtDone応答をソースノード70に送る。TgtDone応答は、ソースノード70にキャッシュブロックフィルのソースを伝える機能を果たす。TgtDone応答の存在は、ターゲットノードメモリ421またはターゲットノード内部キャッシュ(図示せず)が要求されたデータの失効した状態のものを有することを、ソースノードに示し、よってソースノード70は残りのノードのうちの1つ(たとえばノード74または76)からのキャッシュブロックの変更されたコピーを待たねばならない。

20

【0065】

ソースノードプロセッサは、TgtDone応答の受取の前に、応答ノード76からのRdResponseパケットと併せて送信される変更されたキャッシュブロックを用い得る。しかしながら、ソースノード70は、SrcDone応答を送る前にそのソースタグ(図6の読出コマンドパケット内のSrcTagフィールド)を再使用し得ないが、これは読出コマンドパケットによって開始されたトランザクション、すなわち読出動作は、読出トランザクションの開始によって生成されたすべての応答をソースノード70が受取るまで、完了しないおそれがあるためである。したがってソースノード70は、(送信されていれば)ターゲットノード72からのRdResponse、ターゲットノードからのTgtDone応答、および他の残りのノードからの他の応答のいずれか(図14を参照して後に説明)を受取るまで待機し、その後でターゲットノードメモリコントローラへのSrcDone応答を送信する。図9を参照した説明と同様に、図13におけるSrcDone応答は、ソースノードによって開始されたメモリ読出トランザクションの完了の信号を、ターゲットノードに送る。ターゲットノード72がRdResponseとTgtDone応答とを送信すると、ソースノードはこれら両方の応答を待ってから、SrcDone応答を介した読出トランザクションの完了に対して肯定応答しなければならないであろう。SrcDone応答はこうして、メモリ読出トランザクションの間のキャッシュブロックのフィル-プローブ順序の維持を助けるが、これは要求されたデータのソースがターゲットノードメモリコントローラであるか、ターゲットノード内部(または外部)キャッシュであるか、または要求されたデータを含むキャッシュブロックのダーティコピーを有する残りのノードのうちの1つであるかにかかわらない。

30

40

【0066】

図14を参照すると、ソースノード70によって開始されるメモリ読出トランザクション

50

に関するパケットフロー構成 56 を示す。この実施例は 1 つ以上の残りのノード、すなわちノード 74 および 76 を示し、残りのノードのうちの 1 つである 76 はそのキャッシュ内に要求されたデータを含むメモリブロックのダーティ（変更された）コピーを有すると想定する。図 14 に示すさまざまなコマンドおよび応答パケットは、図 9 から図 13 を参照して先に説明したものと同様である。ソースプロセッサは、ノード 76 から RdResponse と併せて受取ったデータを、システム内の他の残りのノードのすべて（ここではノード 74 のみ）からのプローブ応答をも受取るまで、使用し得ない。図 13 を参照して説明したように、ソースノードは開始されたトランザクション、すなわちメモリ読出動作が SrcDone 応答の送信によって完全に確立されるまで、SrcTag を再使用し得ない。応答ノード 76 からの RdResponse、すべての残りの処理ノードからのプローブ応答、ターゲットノード 72 からの TgtDone 応答、および（既に送信されていれば）ターゲットノードからの RdResponse が、ソースノード 70 によって受取られたときに SrcDone 応答が送信される。（図 9、図 13、図 14、図 15 A および図 15 B の）SrcDone および TgtDone 応答は、こうして用いられてコマンドと応答との間のエンドツーエンド肯定応答を提供する。

【0067】

図 15 A は、ダーティヴィクティムブロック書込動作の間のパケット 58 の例示的なフローを示す。ダーティヴィクティムブロックとは一般的には、ヴィクティムブロック書込動作を発信する処理ノード、すなわちソースノード 70 内にあるキャッシュ（図示せず）から排除された、変更されたキャッシュブロックであって、好適なキャッシュブロック置き換えアルゴリズムのいずれかに従って置き換えられる。置き換えのためにダーティヴィクティムブロックが選択されると、ここではターゲットノード 72 に関連するメモリ 421 である、対応のシステムメモリ内に VicBlk コマンドを用いてライトバックされる。メモリライトバック動作は、VicBlk パケットを用いて開始され、その後に変更されたヴィクティムキャッシュブロックを含むデータパケットが続く。VicBlk コマンドに対してはプローブは必要ではない。したがって、ターゲットメモリコントローラが受取ったヴィクティムブロックデータをメモリ 421 にコミットするよう準備されると、ターゲットメモリコントローラは TgtDone パケットをソースノードプロセッサに送る。ソースノードプロセッサは、SrcDone パケットで応答してデータがコミットされるべきことを表示するか、または MemCancel パケットで応答してデータが VicBlk コマンドの送信と TgtDone パケットの受信との間で（たとえば、介入するプローブに回答して）無効化されたことを表示する。

【0068】

ソースノード 70 は、ヴィクティムブロックがシステムメモリ 421 内の適切なメモリ位置に書込まれるためにターゲットノードメモリコントローラ（図示せず）によって受取られるまで、該ヴィクティムブロックを所有することに留意されたい。ターゲットノード 72 は、受取ったヴィクティムブロックをそのコマンドデータバッファ内に配置し、ソースノードプロセッサ（図示せず）に TgtDone 応答を送り返して、ヴィクティムブロックの受取を表示し得る。ソースノード 70 は、TgtDone 応答を受取るまで、ヴィクティムブロック内に含まれるデータを含む他のトランザクションを処理し続ける。

【0069】

図 15 B を参照すると、TgtDone 応答の前のソースノード 70 による無効化プローブの受取を示す、パケット 59 の詳細なフローを示す。上述のように、制御またはデータパケットのターゲットノードへの引渡しは、システム 10 内に含まれる経路制御に依存する。図 15 B に示すように、ソースノード 70 からの VicBlk コマンドおよびヴィクティムブロックデータパケット（図 15 B においてデータ - 1 と示す）は、残りのノード 74 のうちの 1 つを含む経路を通過し得る。システム 10 内のパケット伝搬にかかわる時間は、一般的には経路内に介在する処理ノードの数と、介入するノードが受取ったコマンドおよびデータパケットを経路上の他の処理ノードに、またはこの場合のようにターゲットノード 72 に、伝送するためにかかる時間とに依存する。

【0070】

図 15 B は、ソースノード 70 が VicBlk コマンドとヴィクティムブロック（データ - 1）

10

20

30

40

50

とを送信した後であって、このVicBlkコマンドがターゲットノードメモリコントローラ（図示せず）によって受取られる前に、残りの処理ノードのうちの１つ（ここではノード７６）がRdBlkModコマンドをターゲットノード７２に送る１つの例を示す。ノード７６からのRdBlkModコマンドは、ソースノード７０からのヴィクティムブロックに対する宛先である、メモリ４２１内の同一のメモリ位置を特定し得る。先に簡単に説明したように、キャッシュブロックの書込可能なコピーが所望である場合にはRdBlkModコマンドを用い得る。RdBlkModコマンドは、読出コマンドの１種であることから、RdBlkModコマンド実行の間に図９から図１４を参照して示し説明したさまざまな信号フローパターンが生じ得る。

【００７１】

RdBlkModコマンドに回答して、図９を参照して説明したように、ターゲットノード７２はプローブコマンドパケット（図１０Ａ）をソースノード７０と他の残りのノード７４とに送信し得る。ソースノード７０は、要求されたデータすなわちヴィクティムブロックと併せて読出応答パケット（図１１Ａ）を送ることにより、プローブコマンド（無効化プローブとしても知られる）に回答するが、これは（i）ソースノードは指定されたメモリ位置の変更されたコピー（すなわちヴィクティムブロック）を有し、かつ（ii）ソースノードは、先行して送信したヴィクティムブロックのターゲットノードによる受取と受認とを表示する、ターゲットノード７２からのターゲット終了応答をまだ受取っていないためである。ソースノード７０はまた、図１３を参照して先に説明したように、メモリキャンセル応答（図１５Ｂには図示せず）をターゲットノード７２に送ってもよい。他の残りのノード７４から読出コマンドのソースノード７６へのプローブ応答は、より明確にするため図１５Ｂには示さない。

【００７２】

処理ノード７６はまた、ソースノード７０からの受取ったヴィクティムブロックを変更し、かつ変更されたデータ（データ－２）をターゲットノード７２に送信して、データ－２をシステムメモリ４２１内の対応するメモリ位置にコミットさせることができる。図１５に示す状況においては、VicBlkコマンドおよびもとのヴィクティムブロック（データ－１）は、ターゲットノードが変更されたヴィクティムブロック（データ－２）を受取った後で、ターゲットノードに到着する。もとのヴィクティムブロック（データ－１）を受取ると、ターゲットノードメモリコントローラ（図示せず）は、ターゲット終了応答をソースノード７０に送信して、ヴィクティムブロックデータパケット（データ－１）の受認を表示する。ターゲットノード７２はデータ送信事象の履歴を追跡し得ないために、ターゲットノード７２が、変更されたヴィクティムブロック（データ－２）を含むメモリ位置に、（より早く送信されたが）遅れて到着した、失効したヴィクティムブロック（データ－１）を上書きすることを防ぐことが望ましい。この場合ソースノード７０は、ソースノードがターゲットノード７２からターゲット終了応答を受取ったときに、SrcDone応答の代わりにMemCancel応答をターゲットノードメモリコントローラに送信する。ソースノード７０からのMemCancel応答はこうして、ターゲットノード７２が失効したデータ（データ－１）を共通のメモリ位置に上書きすることを防ぐ。

【００７３】

一般的に、ソースノード７０は、ソースノードが無効化プローブをTgtDoneメッセージを受取る前であって、VicBlkコマンドおよびヴィクティムブロックデータパケットを送信した後であればいつでも、ターゲットノードからのTgtDoneメッセージに回答してメモリキャンセルメッセージ（MemCancel）を送る。メモリキャンセル応答はこうして、たとえば、ソースノード７０以外の処理ノード（ここではノード７６）が、図１５Ｂに示されるように、ソースノード７０によって先立って送られたヴィクティムブロック（データ－１）内に含まれるデータを変更する意図を表示する場合の、システム内のさまざまな処理ノード間のキャッシュコヒーレンスを維持する。メモリキャンセル応答はまた、システムメモリ４２１にコミットされるべきデータがもはや有効でない場合に、ターゲットノードメモリコントローラが長いメモリ書込動作を開始することを防ぐことにより、システムメモリ帯域幅を節約し得る。

【0074】

これに代えてソースノードプロセッサは、TgtDone応答が無効化プローブに先立って受取られた場合、図15Bに点線の矢印によって示すように、SrcDoneパケットをターゲットノードメモリコントローラに送ってもよい。言い換えると、ソースノードは、ヴィクティムブロックがまだ有効であれば、TgtDoneメッセージを受取った後にSrcDone応答をターゲットノードに送ってもよい。図15Bに示す状況においては、ソースノードは、プローブコマンドがターゲット終了応答の後に到着したときに、読出応答パケットの代わりにプローブ応答パケット(図12)を送ってもよいが、これはソースノードがソース終了メッセージを送ることによりヴィクティムブロックをターゲットノードに解放すると、もはやヴィクティムブロックを有し得ないためである。SrcDone応答は、ソースノードプロセッサによって開始されたダーティヴィクティムブロック(データ-1)書込動作の終了信号を送る。メモリキャンセル応答は必要ではないが、これはたとえば、変更されたヴィクティムブロック(データ-2)を含む同じメモリ位置への後の書込動作が、先行の(よって失効した)ヴィクティムブロック(データ-1)を正確に上書きするためである。処理ノード間のキャッシュコヒーレンシはこうして適切に維持される。

10

【0075】

ヴィクティムブロックコマンド(VicBlk)はシステムメモリに対してのみ向けられ、かつコヒーレントな処理ノード(すなわち図1における処理ノード12A-12Dのうちの1つ)によってのみ生成されることができ、たとえばI/Oブリッジ20によっては生成されないことに留意されたい。SrcDoneおよびTgtDone応答とは、上述のようにコマンドと応答とのエンドツーエンド肯定応答を提供するために用いられる。

20

【0076】

最後に、図16Aはメモリ書込動作に含まれるトランザクション(サイズ指定された読出またはブロック読出動作)に対する例示的なフローチャート60を示す。さらに、図16Bはダーティヴィクティムブロック書込動作に関連するトランザクションに対する例示的なフローチャート62を示す。図16Aおよび図16Bのフローチャート内のさまざまなブロックに関連の詳細のすべては、図9から図15Bを参照に先に説明した。(コマンドパケットと応答パケットとを含む)さまざまな制御パケットとデータパケットとが、図3から図8および図10から図12において例示的な実施例を用いて示された。システムは同様の目的のために他の制御およびデータパケットを実現化し得るが、異なったフォーマットおよび符号化を使用する。図1のシステム構成におけるコマンドおよび応答パケットを含むこのメッセージ通信方式は、他のシステム構成においても実現化し得る。

30

【0077】

先の説明は、マルチプロセッシングコンピュータシステム環境におけるキャッシュコヒーレントなデータ転送通信方式を開示する。データ転送通信方式は、ターゲット処理ノードにより遅いシステムメモリバスでの比較的長いメモリ読出または書込動作を打切らせることにより、システムメモリ帯域幅を節約し得る。コマンドと応答とのエンドツーエンド肯定応答は、マルチプロセッシングシステムを通してキャッシュコヒーレンシを維持し得る。

【0078】

この発明はさまざまな変形および代替形に対処するものであるが、その特定の実施例は本明細書と図面中に例示の目的でのみ示された。しかしながら、図面と詳細な説明とはこの発明の範囲を開示された特定の形に限定するものではなく、反対に、すべてのそのような変形、等価物および代替物を、前掲の特許請求の範囲に規定されるこの発明の精神および範囲に入れるものであることを理解されたい。

40

【0079】

【産業的用途】

この発明は一般的にコンピュータシステムに適する。

【図面の簡単な説明】

【図1】 コンピュータシステムの一実施例のブロック図である。

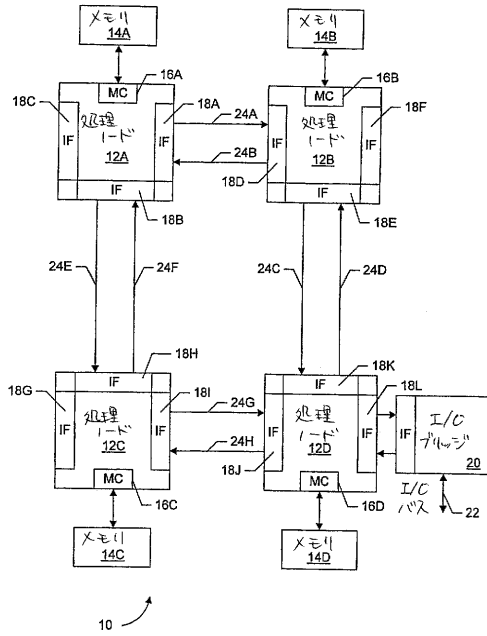
50

- 【図 2】 図 1 からの、1 対の処理ノードの間の相互接続の一実施例の詳細な図である。
- 【図 3】 情報パケットの一実施例のブロック図である。
- 【図 4】 アドレスパケットの一実施例のブロック図である。
- 【図 5】 応答パケットの一実施例のブロック図である。
- 【図 6】 コマンドパケットの一実施例のブロック図である。
- 【図 7】 データパケットの一実施例のブロック図である。
- 【図 8】 図 1 のコンピュータシステムにおいて用い得る例示的なパケットタイプを示すテーブルである。
- 【図 9】 メモリ読出動作に対応するパケットのフローの例を示す図である。
- 【図 10 A】 プローブコマンドパケットの一実施例のブロック図である。
- 【図 10 B】 図 10 A のプローブコマンドパケットにおける nextState フィールドに対する符号化の一実施例のブロック図である。
- 【図 11 A】 読出応答パケットの一実施例のブロック図である。
- 【図 11 B】 一実施例において図 11 A の読出応答パケットの Probe、Tgt および Type フィールドの関係を示す図である。
- 【図 12】 プローブ応答パケットの一実施例のブロック図である。
- 【図 13】 メモリキャンセル応答にかかわる、パケットのフローの例を示す図である。
- 【図 14】 プローブコマンドとメモリキャンセル応答とを組合せるメッセージ通信方式を示すパケットのフローの例を示す図である。
- 【図 15 A】 ヴィクティムブロック書込動作の間のパケットのフローの例を一般的に示す図である。
- 【図 15 B】 ヴィクティムブロック書込動作の間の無効化プローブとメモリキャンセル応答とを示すパケットのフローを詳細に示す図である。
- 【図 16 A】 メモリ読出動作に含まれるトランザクションに対する例示的なフローチャートの図である。
- 【図 16 B】 ヴィクティムブロック書込動作に含まれるトランザクションに対する例示的なフローチャートの図である。

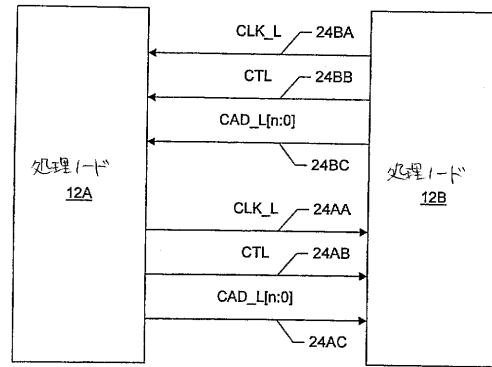
10

20

【 図 1 】



【 図 2 】



【 図 3 】

ビット番号	7	6	5	4	3	2	1	0
1			CMD[5:0]					
2								

【 図 4 】

$\frac{N}{N_{\text{tag}} - 0.5 \frac{N}{N_{\text{tag}}}}$	7	6	5	4	3	2	1	0
1	DestNode [1:0]	CMD[5:0]						
2	SrcTag [1:0]	SrcNode[3:0]					DestNode [3:2]	
3			SrcTag[6:2]					
4	Addr[7:0]							
5	Addr[15:8]							
6	Addr[23:16]							
7	Addr[31:24]							
8	Addr[39:32]							

【 図 6 】

ビット時間	7	6	5	4	3	2	1	0
1	TgtNode [1:0]	CMD[5:0]						
2	SrcTag [1:0]	SrcNode[3:0]					TgtNode [3:2]	
3	Count			SrcTag[6:2]				
4	Addr[7:0]							
5	Addr[15:8]							
6	Addr[23:16]							
7	Addr[31:24]							
8	Addr[39:32]							

【 図 5 】

$C_{i-1} = \{c_{i-1}^0, c_{i-1}^1, c_{i-1}^2, c_{i-1}^3\}$	7	6	5	4	3	2	1	0
1	DestNode [1:0]	CMD[5:0]						
2	SrcTag [1:0]	SrcNode[3:0]				DestNode [3:2]		
3				SrcTag[6:2]				
4								

【図 7】

ビット時間	7	6	5	4	3	2	1	0
1	Data [7:0]							
2	Data [15:8]							
3	Data [23:16]							
4	Data [31:24]							
5	Data [39:32]							
6	Data [47:40]							
7	Data [55:48]							
8	Data [63:56]							

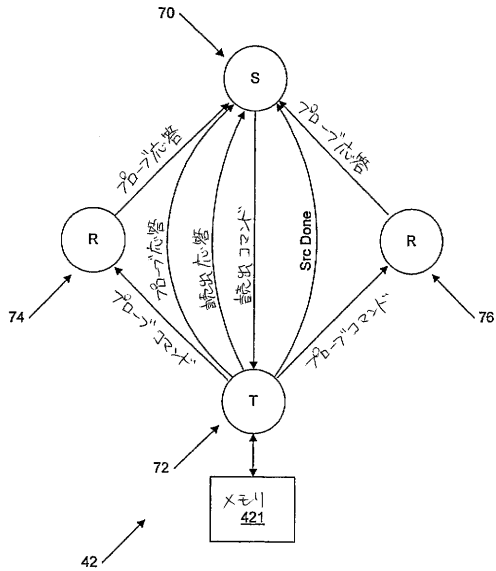
38

【図 8】

CMD ⁴⁴	コマンド	パケットタイプ ⁴⁵
000000	Nop	Info
000001	Interrupt Broadcast	Command
000011	Reserved	-
000100	Probe/Src	Command/Address
000101	Probe/Tgt	Command/Address
00011x	Reserved	-
001000	RdBlkS	Command/Address
001001	RdBlk	Command/Address
001010	RdBlkMod	Command/Address
001011	ChangeToDirty	Command/Address
001100	ValidateBlk	Command/Address
001110	CleanVicBlk	Command/Address
001101	Interrupt Target	Command
001111	VicBlk	Command/Address/Data
01xxxx	Read(Sized)	Command/Address
10xxxx	Write(Sized)	Command/Address/Data
11000x	RdResponse	Rd Response
11001x	Reserved	-
11010x	ProbeResp	Response
11011x	Reserved	-
111000	SrcDone	Response
111001	MemCancel	Response
111010	TgtStart	Response
111011	Reserved	-
11110x	TgtDone	Response
111110	IntrResponse	Response
111111	Error	Info

40

【図 9】



【図 10 A】

ビット時間	7	6	5	4	3	2	1	0
1	TgtNode [1:0]		CMD[5:0]					
2	SrcTag [1:0]		SrcNode[3:0]				TgtNode [3:2]	
3	NextState [1:0]		DM	SrcTag[6:2]				
4	Addr[7:0]							
5	Addr[15:8]							
6	Addr[23:16]							
7	Addr[31:24]							
8	Addr[39:32]							

44

【図 10 B】

Next State [1:0]	ネクストステート
0	変化なし
1	クイック → 安用 ダミー → 安用/ダミー
2	無効

46

【図 11A】

C ₁ 時間	7	6	5	4	3	2	1	0
1	RespNode [1:0]		CMD[5:0]					
2	SrcTag [1:0]		SrcNode[3:0]				TgtNode [3:2]	
3	Count			SrcTag[6:2]				
4	Rsv						Probe	Type

48

【図 11B】

Probe	Probe	Type, Count
0	X	ターゲットの符号化長
1	0	フェーズプロビ: Typeは1, Countは7がある
1	1	フェーズプロビ: Typeは1, Countは7がある

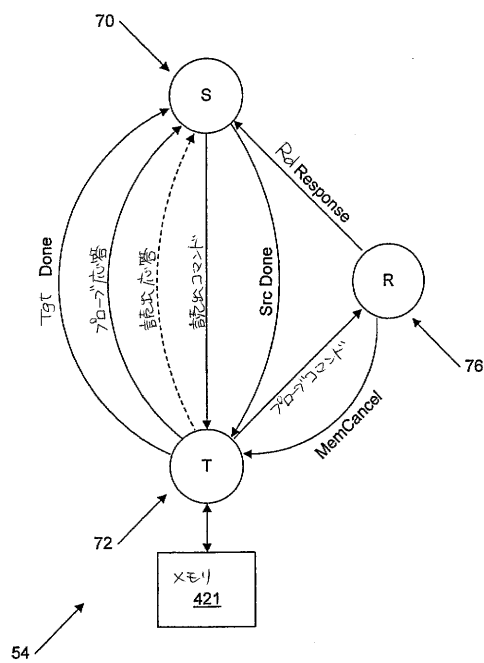
50

【図 12】

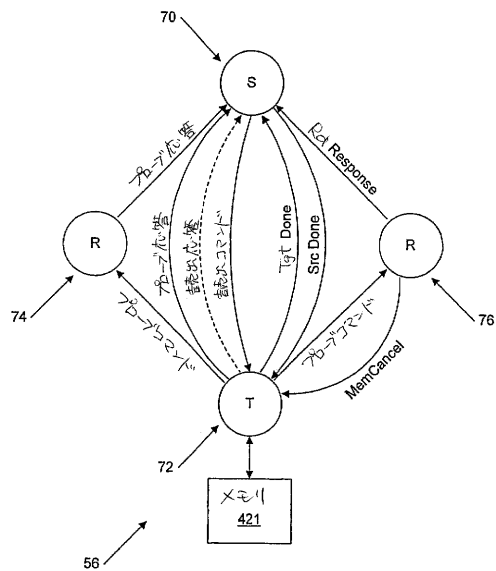
ビット時間	0	1	2	3	4	5	6	7
1	CMD[5:0]		RespNode [3:2]		SrcNode[3:0]		SrcTag[6:2]	
	RespNode [1:0]		SrcTag [1:0]		Rsv		Hit	
2								
3								
4								
5								
6								
7								

52

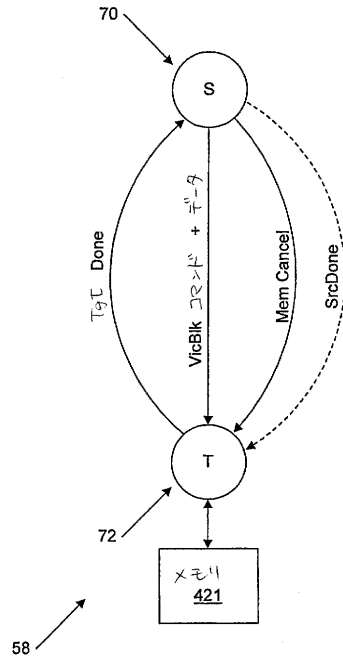
【図 13】



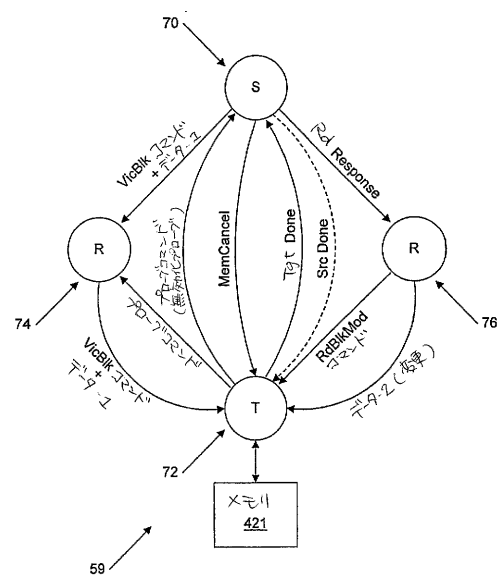
【図 14】



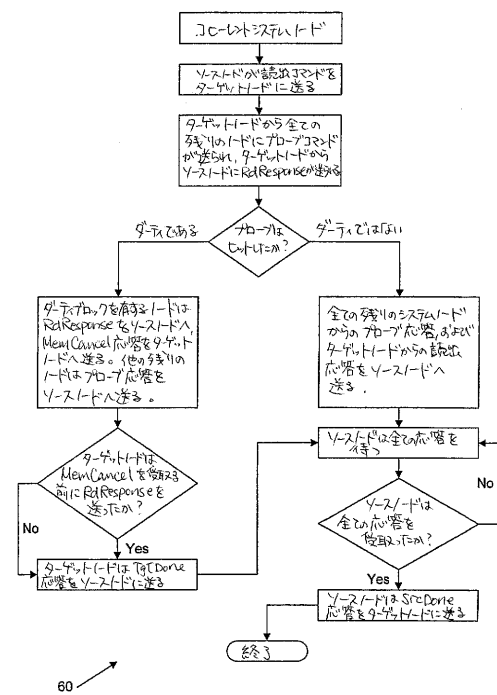
【図 15 A】



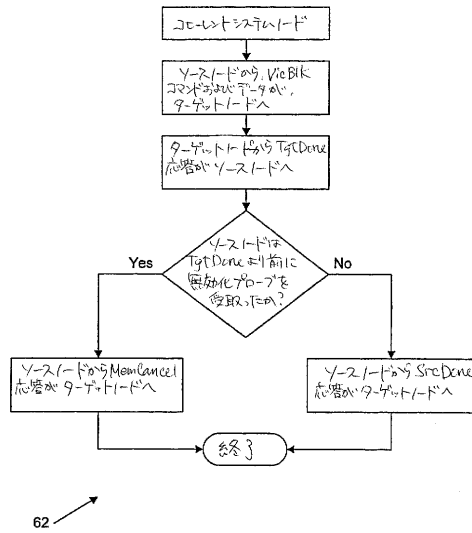
【図 15 B】



【図 16 A】



【図 16 B】



フロントページの続き

(31)優先権主張番号 09/217,649

(32)優先日 平成10年12月21日(1998.12.21)

(33)優先権主張国 米国(US)

(31)優先権主張番号 09/370,970

(32)優先日 平成11年8月10日(1999.8.10)

(33)優先権主張国 米国(US)

(74)代理人 100083703

弁理士 仲村 義平

(74)代理人 100091409

弁理士 伊藤 英彦

(74)代理人 100096781

弁理士 堀井 豊

(74)代理人 100096792

弁理士 森下 八郎

(72)発明者 ケラー, ジェイムズ・ビィ

アメリカ合衆国、9 4 3 0 3 カリフォルニア州、パロ・アルト、イリス・ウェイ、2 1 0

審査官 高瀬 勤

(56)参考文献 特開平06-110844(JP,A)

特開平06-274461(JP,A)

特表2002-533812(JP,A)

特開2000-132531(JP,A)

特開2000-076130(JP,A)

特開平10-105464(JP,A)

特開平08-016474(JP,A)

特開平06-314239(JP,A)

特開平05-134991(JP,A)

特開平02-205963(JP,A)

特開昭60-200351(JP,A)

特表平06-509671(JP,A)

特開平04-113444(JP,A)

特開昭63-223948(JP,A)

特開昭58-064846(JP,A)

特開平09-091262(JP,A)

特開平05-210584(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 12/08