

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4712974号  
(P4712974)

(45) 発行日 平成23年6月29日 (2011.6.29)

(24) 登録日 平成23年4月1日 (2011.4.1)

(51) Int.Cl.

F I

G 0 6 F 12/08 (2006.01)

G 0 6 F 12/08 5 3 1 B

G 0 6 F 12/08 5 0 1 Z

G 0 6 F 12/08 5 5 1 C

G 0 6 F 12/08 5 7 5

請求項の数 14 (全 23 頁)

(21) 出願番号 特願2000-590061 (P2000-590061)  
 (86) (22) 出願日 平成11年8月26日 (1999.8.26)  
 (65) 公表番号 特表2002-533812 (P2002-533812A)  
 (43) 公表日 平成14年10月8日 (2002.10.8)  
 (86) 国際出願番号 PCT/US1999/019471  
 (87) 国際公開番号 W02000/038069  
 (87) 国際公開日 平成12年6月29日 (2000.6.29)  
 審査請求日 平成18年7月12日 (2006.7.12)  
 (31) 優先権主張番号 09/217,367  
 (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)

(73) 特許権者 591016172  
 アドバンスト・マイクロ・ディバイシズ・  
 インコーポレイテッド  
 ADVANCED MICRO DEVI  
 CES INCORPORATED  
 アメリカ合衆国、94088-3453  
 カリフォルニア州、サニベイ、ピー・  
 オウ・ボックス・3453、ワン・エイ・  
 エム・ディ・プレイス、メイル・ストップ  
 ・68 (番地なし)  
 (74) 代理人 100064746  
 弁理士 深見 久郎  
 (74) 代理人 100085132  
 弁理士 森田 俊雄

最終頁に続く

(54) 【発明の名称】 コヒーレンシ維持のための柔軟なプローブ／プローブ応答経路制御

(57) 【特許請求の範囲】

【請求項 1】

装置であって、

プローブを受取るよう結合される第1のノードと、前記第1のノードに結合される第2のノードおよび第3のノードと、各前記ノードに結合されたメモリとを含み、前記第1のノードはキャッシュを含み、前記プローブは、前記第1のノードが、前記メモリから読み出されたデータのキャッシュブロックのコピーを前記キャッシュにストアしているかどうかを判断する要求であり、前記キャッシュブロックが前記第1のノードにストアされている場合に前記第1のノードによって行なわれる動作の第1の表示を含んでおり、前記プローブは、ソースノードとして機能する前記第3のノードによってターゲットノードとして機能する前記第2のノードに送信される要求にตอบสนองして、前記ターゲットノードによって生成され、前記第1のノードはプローブにตอบสนองしてプローブ応答を生成するよう構成され、前記第1のノードは、プローブ内の第2の表示にตอบสนองしてターゲットノードおよびソースノードのいずれの1つがプローブ応答を受取るべきかを選択するよう構成され、前記第1のノードは、前記第2の表示にตอบสนองして、ターゲットノードまたはソースノードのいずれか1つにプローブ応答を送信するよう構成されており、前記要求が読出ならば、前記第2の表示は前記ソースノードを示し、前記要求が書込ならば、前記第2の表示は前記ターゲットノードを示す、装置。

【請求項 2】

前記プローブは、ターゲットノードを示す第1のノード番号と、ソースノードを示す第

10

20

2のノード番号とを含み、前記第1のノードは、プローブ内の前記第2の表示に応答するプローブ応答を経路制御するために前記第1のノード番号および前記第2のノード番号のいずれか1つを選択するよう構成される、請求項1に記載の装置。

【請求項3】

前記プローブ応答は、前記第1のノードがプローブによって識別されたデータの修正されたコピーをストアしている場合に、データを伴いうる、請求項1に記載の装置。

【請求項4】

前記プローブは、前記修正されたコピーがプローブに응答して送信されるべきか否かを示すデータ移動表示を含み、前記プローブ応答は、データ移動表示が修正されたデータを送信すべきことを示す場合に、データを伴いうる、請求項3に記載の装置。

10

【請求項5】

前記プローブは、前記第1のノードがデータをキャッシュしている場合における、前記第1のノードにおいて前記プローブによって識別されるデータの次の状態を識別するネクストステートフィールドを含み、前記第1のノードは、ネクストステートフィールドに응答して前記第1のノードにおけるデータの状態を変更するよう構成される、請求項1に記載の装置。

【請求項6】

方法であって、

第1のノードにおいてプローブを受取るステップを含み、前記第1のノードは、第2のノードおよび第3のノードに結合されており、各前記ノードはメモリに結合されており、前記第1のノードはキャッシュを含み、前記プローブは、ソースノードとして機能する前記第3のノードによってターゲットノードとして機能する前記第2のノードに送信される要求に응答して、前記ターゲットノードによって生成され、前記プローブは、前記第1のノードが、前記メモリから読み出されたデータのキャッシュブロックのコピーを前記キャッシュにストアしているかどうかを判断する要求であり、前記キャッシュブロックが前記第1のノードにストアされている場合に前記第1のノードによって行なわれる動作の第1の表示を含んでおり、

20

方法はさらに、

前記プローブに응答して前記第1のノードにおいてプローブ応答を生成するステップと、

30

プローブ内の第2の表示に응答して、ターゲットノードおよびソースノードのいずれの1つがプローブ応答を受取るべきかを選択するステップと、

前記第2の表示に응答して、ターゲットノードまたはソースノードのいずれか1つにプローブ応答を送信するステップとを含み、前記要求が読出ならば、前記第2の表示は前記ソースノードを示し、前記要求が書込ならば、前記第2の表示は前記ターゲットノードを示す、方法。

【請求項7】

前記プローブは、ターゲットノードを示す第1のノード番号と、ソースノードを示す第2のノード番号とを含み、方法はさらに、

プローブ内の前記第2の表示に응答するプローブ応答を経路制御するために前記第1のノード番号および前記第2のノード番号のいずれか1つを選択するステップを含む、請求項6に記載の方法。

40

【請求項8】

前記第1のノードがプローブによって識別されたデータの修正されたコピーをストアしている場合に、プローブ応答とともにデータを送信しうるステップをさらに含む、請求項6に記載の方法。

【請求項9】

前記プローブは、前記修正されたコピーがプローブに응答して送信されるべきか否かを表示するデータ移動表示を含み、前記データを送信しうるステップは、データ移動表示が修正されたコピーを送信すべきことを示す場合に、プローブ応答とともにデータを送信し

50

うるものである、請求項 8 に記載の方法。

【請求項 10】

前記ブローブは、前記第 1 のノードがデータをキャッシュしている場合における、前記第 1 のノードにおいて前記ブローブによって識別されるデータの次の状態を識別するネクストステートフィールドを含み、方法はさらに、

ネクストステートフィールドにตอบสนองして前記第 1 のノードにおけるデータの状態を変更するステップを含む、請求項 6 に記載の方法。

【請求項 11】

第 1 のノードに結合されて、情報を転送するためのトランザクションが向けられるターゲットノードとして機能するノードであって、前記第 1 のノードはキャッシュを含み、各前記ノードにはメモリが結合されており、

前記ターゲットノードに結合されて、トランザクションを開始するソースノードとして機能するノードによって生成される要求を受信するための手段を含み、

前記受信するための手段は、前記要求にตอบสนองしてブローブを生成するように構成されており、

前記ブローブは、当該ブローブを受信する第 1 のノードが、前記メモリから読み出されたデータのキャッシュブロックのコピーを前記キャッシュにストアしているかどうかを判断するための要求であり、前記キャッシュブロックが前記第 1 のノードにストアされている場合に前記第 1 のノードによって行なわれる動作の第 1 の表示を含んでおり、前記ブローブは、当該ブローブへのตอบสนองを受信するための受信ノードを指定する手段によっ

て生成される第 2 の表示を含み、

前記受信ノードは、前記トランザクションが書込であることにตอบสนองしてターゲットノードとなり、

前記受信ノードは、前記トランザクションが読出であることにตอบสนองしてソースノードとなる、ノード。

【請求項 12】

前記ブローブは、当該ブローブとしてパケットを識別するコマンドフィールドを有するパケットを含み、

前記第 2 の表示は、前記コマンドフィールドに含まれる、請求項 11 に記載のノード。

【請求項 13】

前記ブローブは、前記ターゲットノードを識別するターゲットノードフィールドと、前記ソースノードを識別するソースノードフィールドとを有する第 1 のパケットを含む、請求項 11 に記載のノード。

【請求項 14】

前記要求を受信するための手段は、前記キャッシュブロックがストアされるメモリと通信するように構成されたメモリコントローラを含み、

前記メモリコントローラは、前記メモリにアクセスする要求を選択することに対応して、前記ブローブを生成するように構成されている、請求項 11 に記載のノード。

【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、より特定的には、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。

【0002】

2. 関連技術分野の背景

クメール・A (Kumar A) 他による「共有メモリハイパーキューブマイクロプロセッサのための、効率的でスケラブルなキャッシュコヒーレンス方式 (Efficient and Scalable Cache Coherence Schemes for Shared Memory Hypercube Microprocessor)」Proceedings of the Supercomputing Conference, US Los Alamitos, IEEE, Vol.Conf.7, 1994, pp

10

20

30

40

50

.498-507は、前掲の請求項1のプリアンブルに記載の特徴を有しキャッシュコヒーレンシが維持されるコンピュータシステムと、前掲の請求項10のプリアンブルの特徴を有しコンピュータシステム内でキャッシュコヒーレンシを維持するための方法とを開示する。

EP-A-0 817 076は、ローカルおよびグローバルアドレススペースと、多数のアクセスモードとを用いる、マルチプロセッシングコンピュータシステムを開示する。ノード内のプロセッサは、ノード間通信を要求するトランザクションを開始し得る。要求するノードからホームノードへ要求が送られると、ホームノードは、要求されたデータのキャッシュコピーを有する従属ノードのいずれかに、読出および/または無効化要求を送る。システムは、キャッシュコヒーレンシを維持するためのグローバルコヒーレンシプロトコルを実現化する。

10

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、共用バスとI/O装置との間の情報の転送を管理するI/Oブリッジを介して共用バスに結合される場合もある一方、プロセッサは典型的には、直接共用バスに結合されるか、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。多数の装着物は、バス上で信号を駆動する装置に容量的な高負荷をもたらし、多数の装着点は、比較複雑な高周波数に対する伝送ラインモデルをもたらし、したがって、周波数は低く留められ、共用バスで利用できる帯域幅も同様に比較的低い。低帯域幅は共用バスに付加的な装置を装着するのに障壁になり得るが、これは使用できる帯域幅によって性能が制限されるおそれがあるためである。

20

#### 【0004】

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。全体的な性能が減じられるであろう。

30

#### 【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でき、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】

40

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シを維持するための機構が望まれる。

#### 【0007】

#### 【発明の開示】

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。この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1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

50

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シを維持し得る。特に、このコンピュータシステムは柔軟なプローブコマンド/応答経路制御方式を実現化し得る。

【0008】

一実施例においては、この方式はプローブコマンド内の表示を用い、これはプローブ応答を受信する受信ノードを識別する。一般的にはプローブコマンドとは、キャッシュブロックがノード内にストアされているか判断するためのそのノードへの要求であり、かつ、キャッシュブロックがそのノードにストアされていた場合にそのノードが行なうべき動作の表示である。プローブ応答は、動作が行なわれたことを表示し、かつ、キャッシュブロックがノードによって変更されていた場合にはデータの送信を含み得る。送られたコマンドに依存してプローブ応答を異なった受信ノードに経路制御することへ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（たとえば最少の数のパケット送信を処理ノード間で用いて）行ない得る一方で、それでもコヒーレンシが維持されることを確実にする。

10

【0009】

たとえば、ターゲットまたはトランザクションの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ことを表示するプローブコマンドを含み得る。プローブコマンドは、トランザクションのソースを讀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するノードからダーティデータがソースノードに引渡される）。一方で、（データがトランザクションのターゲットノードでメモリ内に更新される）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。このようにして、ターゲットはいつ書込データをメモリにコミットするか判断し、書込データとマージされるべきいかなるダーティデータをも受取ることができる。

20

【0010】

概略的には、コンピュータシステムが企図される。コンピュータシステムは、第1の処理ノードと第2の処理ノードとを含み得る。第1の処理ノードは、要求を送信することによりトランザクションを開始するよう構成し得る。第2の処理ノードは、第1の処理ノードからの要求を受取るよう結合されて、要求に応答してプローブを生成するよう構成し得る。プローブは、プローブに対する応答を受けるための受信ノードを指定する表示を含み得る。さらに、第2の処理ノードはトランザクションのタイプに応答して、表示を生成するよう構成し得る。

30

【0011】

コンピュータシステム内でのコヒーレンシを維持するための方法もまた企図される。ソースノードからの要求はターゲットノードに送信される。プローブは要求に応答してターゲットノード内で生成される。プローブ内の表示を介して、プローブに対する応答のための受信ノードが指定される。プローブに対するプローブ応答は、受信ノードに経路制御される。

【0012】

この発明の他の目的と利点とは、以下の詳細な説明を読み、添付の図面を参照することにより、より明らかとなるであろう。

40

【0013】

この発明はさまざまな変形と代替形に対処するものであるが、その特定の実施例を図面において例示の目的で示し、以下に詳述する。しかしながら、図面とその詳細な説明とは開示される発明を特定の形に限定することを意図せず、反対に、すべての変形、等価物、および代替例は前掲の特許請求の範囲に規定されるこの発明の範囲に入ることを意図する。

【0014】

【発明の実施の形態】

例示的なコンピュータシステムの実施例

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。他の実施例

50

が可能であり企図される。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 12 A、12 B、12 C、および 12 D を含む。処理ノードの各々は、処理ノード 12 A - 12 D にそれぞれ含まれるメモリコントローラ 16 A - 16 D を介して、それぞれのメモリ 14 A - 14 D に結合される。さらに、処理ノード 12 A - 12 D は、処理ノード 12 A - 12 D の間の通信のために用いるインターフェイスロジックを含む。たとえば、処理ノード 12 A は、処理ノード 12 B と通信するためのインターフェイスロジック 18 A と、処理ノード 12 C と通信するためのインターフェイスロジック 18 B 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 18 C を含む。同様に、処理ノード 12 B はインターフェイスロジック 18 D、18 E、および 18 F を含む、処理ノード 12 C はインターフェイスロジック 18 G、18 H、および 18 I を含む、処理ノード 12 D はインターフェイスロジック 18 J、18 K、および 18 L を含む。処理ノード 12 D は、インターフェイスロジック 18 L 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。

#### 【0015】

処理ノード 12 A - 12 D は、ノード通信相互処理のためのパケットベースのリンクを実現化する。この実施例においては、リンクは単方向ラインの組として実現化される（たとえば、ライン 24 A はパケットを処理ノード 12 A から処理ノード 12 B へ送信するために用いられ、ライン 24 B はパケットを処理ノード 12 B から処理ノード 12 A へ送信するために用いられる）。ライン 24 C - 24 H の他の組は、図 1 に示すようにパケットを他の処理ノード間で通信するために用いられる。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 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 つ以上の中間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 12 A から処理ノード 12 D に送信されるパケットは、図 1 に示すように、処理ノード 12 B または処理ノード 12 C のいずれかを通過し得る。いかなる好適な経路制御アルゴリズムをも用い得る。コンピュータシステム 10 の他の実施例は、図 1 に示す実施例よりもより多いかまたはより少ない処理ノードを含み得る。

#### 【0016】

処理ノード 12 A - 12 D は、メモリコントローラと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の間で通信するためのメモリコントローラを含む。この開示では「ノード」という用語も使用し得る。ノードという用語は、「処理ノード」を意味することを意図する。

#### 【0017】

メモリ 14 A - 14 D は、何らかの好適なメモリ装置を含み得る。たとえば、メモリ 14 A - 14 D 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 のアドレス空間は、メモリ 14 A - 14 D の間で分割される。処理ノード 12 A - 12 D 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 12 A - 12 D 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 16 A - 16 D である。言い換えると、メモリコントローラ 16 A - 16 D は、対応するメモリ 14 A - 14 D 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 16 A - 16 D は、メモリ 14 A - 14 D 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 16 A - 16 D 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。

#### 【0018】

一般的には、インターフェイスロジック 18A - 18L は、リンクからのパケットを受取り、かつ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、インターフェイスロジック 18 の各々は、その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レシーバ内に、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。

10

#### 【0019】

次に図 2 は、処理ノード 12A および 12B のブロック図を示し、それらの間のリンクの 1 実施例を詳細に例示する。図 2 の実施例においては、ライン 24A は、クロックライン 24AA と、制御ライン 24AB と、制御/アドレス/データバス 24AC とを含む。同様に、ライン 24B は、クロックライン 24BA と、制御ライン 24BB と、制御/アドレス/データバス 24BC とを含む。

#### 【0020】

クロックラインは、制御ラインおよび制御/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送信し得る。制御/アドレス/データバスの幅に依存して、多数のクロックラインを用い得る。たとえば 32 ビット制御/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（制御/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、残りの半分の制御/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）。一般的には、「パケット」とは 2 つの処理ノード 12A - 12D の間の通信である。1 つ以上のパケットが「トランザクション」を形成し得るが、これは一方の処理ノードから他方への情報の転送である。トランザクションを形成するパケットは、ソースノード（転送を要求する開始するノード）からのターゲットノード（トランザクションが向けられるノード）へのトランザクションを開始する要求パケットと、コヒーレンスを維持するために他の処理ノード間で送信されるパケットと、データパケットと、トランザクションを終了させる肯定応答パケットとを含み得る。

20

30

#### 【0021】

制御ラインは、制御/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。

40

#### 【0022】

制御/アドレス/データバスは、データ/制御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、制御/アドレス/データバスは、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、所望により他のサイズ

50

の制御／アドレス／データバスをサポートし得る。

【 0 0 2 3 】

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）。

【 0 0 2 4 】

図 3 から図 6 は、コンピュータシステム 10 の一実施例に従って用いられる、例示的なパケットを示す。図 3 から図 5 は制御パケットを示し、図 6 はデータパケットを示す。他の実施例は、所望により異なったパケット定義を用い得る。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 6 は、8 ビット制御／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、ビット時間の各々は、7 から 0 まで番号が付与された 8 つのビットを含む。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送信するために用いられるかのいずれかであり得る。

【 0 0 2 5 】

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 6 ビットを含む。図 4、および図 5 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。

【 0 0 2 6 】

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、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）。ビット時間 4 から 8 までを用いて、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 は、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）を開始するのに用いられるだけでなく、トランザクションを実行する過程において、トランザクションによって影響を受けるメモリアドレスを担持するコマンドについて、コマンドを送信し得る。

【 0 0 2 7 】

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。さまざまな種類の応答パケットは付加的な情報を含み得る。たとえば、読出応答パケットは、後に続く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。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 34 は、トランザクションを行なう間にトランザクションによって影響されるメモリ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 34 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 2 8 】

図 6 は、データパケット（data パケット）36 を示す。データパケット 36 は、図 6 の実

10

20

30

40

50



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36は、送信されるデータの量に依存して、異なった数のビット時間を含み得る。たとえば、一実施例においてはキャッシュブロックは64バイトを含み、したがって8ビットリンク上の64のビット時間を含む。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、キャッシュ不可能な読出および書込に対しては、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。キャッシュブロ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を用いる。

【0029】

図3から図6は、8ビットリンクのためのパケットを示す。16および32ビットリンクのためのパケットは、図3から図6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、16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、32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ビットリンクによるビット時間における、16ビットリンクのビット時間1および32ビットリンクのビット時間1の構成を示す。

【0030】

【数1】

$$BT_{16}[15:0] = BT_8[7:0] \parallel BT_8[7:0] \dots (1)$$

$$BT_{32}[31:0] = BT_{16}[15:0] \parallel BT_{16}[15:0] \parallel BT_{16}[15:0] \parallel BT_{16}[15:0] \dots (2)$$

【0031】

図7は、コンピュータシステム10内のリンクの1つの例示的な実施例によって用いられるコマンドを示すテーブル38を示す。他の実施例も可能であり企図される。テーブル38は、コマンドの各々に割当てられたコマンド符号化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30-34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。

【0032】

読出トランザクションは、ReadSized, 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ReadSized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、ReadSized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である場合。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。この場合はRdBlkSコマンドを用いることができる。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。一般的に、適切な読出コマンドはソースから送信されて、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。ターゲットノード内のメモリコントローラは、システム内の他のノードにProbe/Srcコマンドを送信して、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、およびキャッシュブロックの更新されたコピーを含む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により、コヒーレンスを維持する。Probe/Srcコマンドを受取ったノードの各々は、ProbeRespパケットをソースノードに送信する。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、そのノードはRdResponseパケットおよびダーティデータを送信する。ダーティデータを送信するノードはまた、ターゲットノードによる要求された読出データの送信をキャンセルしようと試みて、ターゲットノードにMemCancelパケットをも送信し得る。さらに、ターゲットノード内のメモリコントローラは、RdResponseパケットとその後続くデータパケット内のデータとを用いて要求された読出データを送信する。ソースノードがプローブされたノードからRdResponseパケットを受取れば、その読出されたデータが用いられる。そうでなければ、ターゲットノードからのデータが用いられる。一旦ソースノードにおいてプローブ応答および読出されたデータの各々が受取られると、ソースノードはトランザクションの終了の肯定応答として、ターゲットノードにSrcDone応答パケットを送信する。

10

**【 0 0 3 3 】**

書込トランザクションはWrSizedまたはVicBlkコマンドを用いて開始され、その後に関連のデータパケットが続く。WrSizedコマンドは、キャッシュ不可書込またはサイズの合わないキャッシュブロックのデータの書込に用いられる。WrSizedコマンドに対するコヒーレンシの維持のために、ターゲットノード内のメモリコントローラはシステム内の他のノードの各々にProbe/Tgtコマンドを送信する。Probe/Tgtコマンドに応答して、プローブされたノードの各々はターゲットノードにProbeRespパケットを送信する。もしプローブされたノードがダーティデータをストアしていれば、プローブされたノードはRdResponseパケットおよびダーティデータで応答する。このようにして、WrSizedコマンドによって更新されたキャッシュブロックは、メモリコントローラに返されて、WrSizedコマンドによって提供されるデータとマージされる。メモリコントローラは、プローブされたノードの各々からプローブ応答を受取ると、ソースノードにTgtDoneパケットを送信してトランザクションの終了の肯定応答を提供する。ソースノードはSrcDone応答パケットで応答する。

20

**【 0 0 3 4 】**

ノードによって変更され、かつノード内のキャッシュで置き換えられたヴィクティムキャッシュブロックは、VicBlkコマンドを用いてメモリに返される。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ビクティムブロックデータをメモリにコミットする準備ができたとき、ターゲットメモリコントローラはビクティムブロックのソースノードにTgtDoneパケットを送信する。ソースノードは、データがコミットされるべきことを示すSrcDoneパケットか、またはデータが(たとえば介入するプローブに)応答して)VicBlkコマンドの送信とTgtDoneパケットの受信との間で無効化されたことを示すMemCancelパケットのいずれかで応答する。

30

**【 0 0 3 5 】**

ソースノードにストアされた書込不可能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、ソースノードによって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。

40

**【 0 0 3 6 】**

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、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。

50

## 【 0 0 3 7 】

プローブ / プローブ 応答 経路 制御

図 8 は、プローブ パケット 4 0 の一実施例のブロック図を示す。異なった態様でプローブ パケットを構成し、代替的な、同様の、または代用の情報を有する他の実施例も可能であり企図される。プローブ パケット 4 0 は、図 4 に示すアドレス パケット 3 2 の 1 種である。図 8 に示すように、プローブ パケット 4 0 はコマンド フィールド ( 図 8 における C M D [ 5 : 0 ] )、ターゲット ノード フィールド ( 図 8 における TgtNode [ 1 : 0 ] および TgtNode [ 3 : 2 ] )、ソース ノード フィールド ( 図 8 における SrcNode [ 3 : 0 ] )、ソース タグ フィールド ( 図 8 における SrcTag [ 1 : 0 ] および SrcTag [ 6 : 2 ] )、データ 移動 フィールド ( 図 8 における D M )、ネクスト ステート フィールド ( NextState [ 1 : 0 ] )、および アドレス フィールド ( 図 8 における ビット 時間 4 - 8 にわたる Addr [ 3 9 : 0 ] ) を含む。

10

## 【 0 0 3 8 】

コマンド フィールドは符号化されてパケット 4 0 をプローブ パケットとして識別する。たとえば、図 7 に示す Probe / Src および Probe / Tgt に対する符号化を用い得る。一般的には、トランザクションのターゲット ノードはプローブ コマンドを生成して、トランザクションによって影響されるキャッシュ ブロックのコヒーレンシを維持する。開始されるトランザクションの種類に基づいて、プローブ コマンドの組の中の 1 つがターゲット ノードによって選択される。選択されたプローブ コマンドは、プローブ コマンドに対する応答のための受信 ノードを識別する。多数の使用できる受信 ノードのうちの 1 つに、柔軟にプローブ 応答を経路制御することにより、( 予め定められた受信 ノードにプローブ 応答を経路制御することとは対照的に )、コヒーレンシの維持は ( たとえば、ノード間で送信されるパケットの数の観点から ) 効率的であって、正確な結果を導く態様で行なうことができる。

20

## 【 0 0 3 9 】

たとえば、キャッシュ ブロックをトランザクションのソース ノードに転送させるトランザクションは、プローブ 応答 ( データを転送する応答を含む ) をソース ノードに向けることにより、コヒーレンシを維持し得る。ソース ノードは、他のノード ( プローブされたノード ) の各々からの応答とターゲット ノードからの応答とを待ち、次いでノード内に送信されたキャッシュ ブロックを確立し、かつターゲット ノードに肯定応答 パケットを送信してトランザクションを終了させる。トランザクションを終了させる前にプローブされたノードからのプローブ 応答を待つことにより、ノードの各々は終了の前にそのトランザクションに対する正しいコヒーレンシ状態を確立し得る。一方、WrSized トランザクション ( キャッシュ ブロックに満たない更新をする ) は、WrSized コマンドと関連のデータとをターゲット ノードに送信することにより開始される。ターゲット ノードは、ターゲット ノードのメモリにデータをコミットするために、ターゲット ノードはソース ノードの代わりにプローブ 応答を受信し得る。特に、WrSized トランザクションによって更新されたバイトを含むキャッシュ ブロックがプローブ ノード内でダーティであれば、ターゲット ノードはプローブに 応答してダーティなデータを受信し、そのデータを WrSized コマンドに関連のデータとマージし得る。一旦プローブ 応答が受取られると、書込動作によって命令されるコヒーレンシ状態はプローブされたノード内で確立され、書込データはメモリにコミットし得る。したがってこの実施例においては、2 つの種類のプローブ コマンドがサポートされ、どのノードがプローブ 応答を受取るべきかをプローブされたノードに対して示す ( たとえばソース またはターゲット ノード )。他の実施例は、所望により、付加的な受信 ノードを示す付加的なプローブ コマンドをサポートし得る。

30

40

## 【 0 0 4 0 】

1 つの例示的な実施例においては、ターゲット ノード内のメモリ コントローラは処理されるべき要求のキューからトランザクションを選択することができ、かつその選択に 応答してプローブ コマンドを生成し得る。プローブ コマンドはコンピュータ システム内のプローブされたノードにブロードキャストされることができ、コマンド フィールドは、プローブ 応答がターゲット ノードまたはソース ノードに経路制御されるべきかどうかを示す。た

50

例えば、図 7 に示すテーブル 3 8 において例示するように、コマンドフィールドのビット 0 は受信ノードを示し得る（バイナリ 0 はソースノードを、バイナリ 1 はターゲットノードを示す）。ターゲットノードフィールドはトランザクションのターゲットノードを識別し、ソースノードフィールドはトランザクションのソースノードを識別する。プローブコマンドもまたターゲットノード内のキャッシュにブロードキャストされ、ターゲットノードはプローブ応答を提供し得ることに留意されたい。言い換えると、ターゲットノードもまたプローブされたノードであり得る。

#### 【 0 0 4 1 】

さらに、データ移動フィールド（一実施例においては、たとえば 1 ビット）を用いて、プローブされたノード内でキャッシュブロックがダーティであった場合にプローブに回答してデータが返されるべきかどうかを表示し得る。データ移動フィールドがいずれの移動も示さなければ（たとえばクリアな状態）、ダーティデータを持つノードはプローブ応答を宛先指定された受信ノードに戻すが、ダーティデータを受信ノードに送信しない。データ移動フィールドが移動を示せば（たとえばセットされた状態）、ダーティデータを持つノードは、ダーティデータを含む読出応答をソースノードに返す。ネクストステートフィールドは、（プローブされたノードがキャッシュブロックのコピーをストアしている場合）プローブされたノードにおいてキャッシュブロックの次の状態がどうなるかを示す。一実施例においては、ネクストフィールドの符号化は以下のテーブル 1 に示すとおりである。

#### 【 0 0 4 2 】

一実施例によると、以下のテーブル 2 は例示的なトランザクションタイプと対応のプローブコマンドを示す。

#### 【 0 0 4 3 】

#### 【表 1】

テーブル 1: ネクストステートフィールド符号化

| 符号化 | 意味                             |
|-----|--------------------------------|
| 00  | 変化なし                           |
| 01  | 共用:<br>クリーン→共用<br>ダーティ→共用／ダーティ |
| 10  | 無効                             |
| 11  | 予約                             |

テーブル 2: トランザクションタイプおよびプローブコマンド

| トランザクションタイプ   | プローブコマンド  | ネクストステート | データ移動 |
|---------------|-----------|----------|-------|
| ReadSized     | Probe/Src | 変化なし(00) | Y(1)  |
| Block Reads   | Probe/Src | 共用(01)   | Y(1)  |
| ChangetoDirty | Probe/Src | 無効(10)   | Y(1)  |
| ValidateBlk   | Probe/Src | 無効(10)   | N(0)  |
| WrSized       | Probe/Tgt | 無効(10)   | Y(1)  |
| VicBlk        | None      | -        | -     |

#### 【 0 0 4 4 】

一般的には、「プローブ」または「プローブコマンド」という用語は、プローブノードに

対する要求を指し、これによりプローブによって（たとえばアドレスフィールドを介して）定義されるキャッシュブロックがプローブされたノードによってストアされているかどうか判断し、かつプローブされたノードからの所望の応答を示す。応答は、キャッシュブロックに対する異なったコヒーレンシ状態を確立すること、および/またはキャッシュブロックを受信ノードに転送することを含む。プローブされたノードは「プローブ応答」を用いて応答するが、これは所望のコヒーレンシ状態変化が行なわれたことを受信ノードに肯定応答する。プローブ応答は、プローブされたノードがダーティデータをストアしていた場合、データをも含み得る。

#### 【 0 0 4 5 】

図 9 は、プローブ応答パケット 4 2 の一実施例のブロック図を示す。応答パケットを異なった態様で構成し、代替的な、同様の、または代用の情報を有する他の実施例も可能であり企図される。プローブ応答パケット 4 2 は、図 5 に示す応答パケット 3 4 の 1 種である。図 8 に示すプローブパケットと同様に、プローブ応答パケット 4 2 はコマンドフィールド (CMD [ 5 : 0 ])、ソースノードフィールド (SrcNode [ 3 : 0 ])、およびソースタグフィールド (SrcTag [ 1 : 0 ]) および SrcTag [ 6 : 2 ]) を含み得る。さらに、プローブ応答パケット 4 2 は、応答ノードフィールド (図 9 における RespNode [ 3 : 2 ]) および RespNode [ 1 : 0 ]) およびヒット表示 (図 9 における Hit) を含み得る。

#### 【 0 0 4 6 】

プローブされたノードは、プローブパケット 4 0 で受信されたアドレスによって識別されたキャッシュブロックに対して、そのキャッシュ内を探索する。キャッシュブロックのコピーを見出すと、プローブされたノードはネクストステートフィールドにより規定されていることにしたがってキャッシュブロックの状態を変化させる。さらに、キャッシュがキャッシュブロックのダーティコピーをストアしており、かつプローブパケット 4 0 のデータ移動フィールドがキャッシュブロックが転送されるべきであることを示すと、プローブされたノードはキャッシュからダーティデータを読み出す。

#### 【 0 0 4 7 】

プローブされたノードは、プローブ応答パケット 4 2 を以下の場合において指定された受信ノードに経路制御する。すなわち、( i ) ダーティデータがない場合、または ( ii ) ダーティデータおよびデータ移動フィールドがいずれの移動も示さない場合。この実施例においてより特定的には、プローブされたノードはプローブパケット 4 0 のターゲットノードフィールド (指定された受信ノードがターゲットノードである場合) か、またはプローブパケット 4 0 のソースノードフィールド (指定された受信ノードがソースノードである場合) のいずれかを讀出し、結果として生じるノード ID をプローブ応答パケット 4 2 の応答ノードフィールドにストアする。さらに、ヒット表示を用いて、プローブノードがキャッシュブロックのコピーを保持しているかどうかを示す。たとえば、ヒット表示は、セットされている場合に、プローブされたノードがキャッシュブロックのコピーを保持していることを示し、そしてクリアされている場合に、プローブされたノードがキャッシュブロックのコピーを保持していないことを示すビットを含み得る。ヒットビットは、以下の場合にクリアされている。すなわち、( i ) プローブされたノード内のキャッシュにおいてキャッシュブロックのコピーが見出されなかった場合、または ( ii ) プローブされたノードにおけるキャッシュ内でキャッシュブロックのコピーが見出されたが、プローブコマンドに応答して無効化されていた場合、である。

#### 【 0 0 4 8 】

図 7 のテーブル 3 8 に示す一実施例においては、プローブ応答パケット 4 2 のコマンドフィールドは ProbResp を示すことができ、さらにプローブ応答が Probe/Src コマンドへの応答であるか、または Probe/Tgt コマンドへの応答であることをさらに示すことができる。より特定的には、コマンドフィールドのビット 0 はクリアされてプローブ応答が Probe/Src コマンドへの応答であることを示し、セットされてプローブ応答が Probe/Tgt コマンドへの応答であることを示す。この表示は受信ノードによって用いられて、プローブ応答が受信ノード内のメモリコントローラに経路制御されるべきか (ビット 0 がセットされている

10

20

30

40

50

場合)、または受信ノード内のキャッシュに対するフィル受信ロジックに対して経路制御されるべきか(ビット0がクリアされている場合)を判断し得る。

【0049】

図10は、読出応答パケット44の一実施例のブロック図を示す。読出応答パケットを異なった態様で構成し、代替的な、同様の、または代用の情報を有する他の実施例も可能であり企図される。応答パケット44は図5に示す応答パケット34の1種である。図9に示すプローブ応答パケットと同様に、読出応答パケット44は、コマンドフィールド(CMD[5:0])、ソースノードフィールド(SrcNode[3:0])、ソースタグフィールド(SrcTag[1:0]およびSrcTag[6:2])、および応答ノードフィールド(RespNode[3:2]およびRespNode[1:0])を含む。さらに、読出応答パケット44は、

10

【0050】

データ移動が要求されることを示すプローブコマンドに対してダーティデータのヒットを検出した、プローブされたノードは、読出パケット44を用いてプローブ応答およびダーティデータを送信し得る。カウントフィールドを用いて送信されるデータの量を示し、タイプフィールドはカウントがバイトまたはクワッドワード(8バイト)で測定されるかを示す。プローブに対して応答するキャッシュブロック転送に対しては、カウントフィールドは8(バイナリ「111」として符号化される)を示し、タイプフィールドはクワッドワード(たとえば、一実施例においてはタイプフィールドにおいてセットされたビット)を示す。プローブフィールドを用いて、読出応答パケット44がターゲットノードまたはプローブされたノードのいずれから送信されているかを示す。たとえば、プローブフィールドは、セットされている場合に読出応答パケット44がプローブされたノードから転送されたことを示し、クリアされている場合に読出応答パケット44がターゲットノードから送信されたことを示す、ビットを含み得る。したがって、プローブされたノードは読出応答パケット44を用いる場合にプローブビットをセットする。

20

【0051】

読出応答パケット44は、データパケットが後に続くことを示すパケットタイプである。したがって、プローブされたノードは読出応答パケット44の後にデータパケットでキャッシュブロックを送信する。プローブされたノードは、読出応答パケット44を、上述のようにプローブ読出パケット42を経路制御するのと同様の態様で経路制御する(たとえば、プローブされたノードは(指定された受信ノードがターゲットノードであれば)プローブパケット40のターゲットノードフィールドを読出すか、または(指定された受信ノードがソースノードであれば)プローブパケット40のソースノードフィールドを読出し、結果として生じるノードIDを読出応答パケット44の応答ノードフィールドにストアする)。

30

【0052】

図7のテーブル38において示される一実施例においては、読出応答パケット44のコマンドフィールドはRdResponseを示すことができ、さらに読出応答がProbe/Srcコマンドへの応答であるか、またはProbe/Tgtコマンドへの応答であるかを示し得る。より特定的には、コマンドフィールドのビット0はクリアされて読出応答がProbe/Srcコマンドへの応答であることを示し、セットされて読出応答がProbe/Tgtコマンドへの応答であることを示す。この表示を受信ノードによって用いて、読出応答が受信ノード内のメモリコントローラに経路制御されるべきか(ビット0がセットされている場合)、または受信ノード内のキャッシュに対するフィル受信ロジックに経路制御されるべきか(ビット0がクリアされている場合)を判断し得る。

40

【0053】

図11は、例示的な読出ブロックトランザクションに対応する1組のノードの間のパケットフローを示す図である。ソースノード50、ターゲットノードメモリコントローラ52、および1組のプローブされたノード54A-54Nを示す。図11にはパケットの(時

50

系列での順序を左から右に示す。言い換えると、RdBlkパケットはソースノード50からターゲットノードメモリコントローラ52に送信され、その後でターゲットノードメモリコントローラ52はProbe/Srcパケットをプローブノード54A - 54Nに送り、以下も同様である。パケットの時間順序を示すために、図11ではソースノード50およびターゲットメモリコントローラ52を2箇所に示す。同様に、図12および図13では特定のブロックを1箇所以上に示す。ソースノード50、ターゲットノードメモリコントローラ52を含むターゲットノード、プローブされたノード54A - 54Nの各々は、図1に示す処理ノード12A - 12Dと同様の処理ノードを含み得る。

#### 【0054】

ソースノード50はRdBlkパケットをターゲットノードメモリコントローラ52に送信して、読出ブロックトランザクションを開始する。次いでターゲットノードメモリコントローラ52は処理されるべきRdBlkパケットを選択する。ターゲットノードメモリコントローラ52はProbe/Srcパケットを生成し、パケットをプローブされたノード54A - 54Nにブロードキャストする。さらに、ターゲットノードメモリコントローラ52は、ターゲットノードメモリコントローラ52が結合されているメモリ14A - 14Dからの読出を開始する。メモリ14A - 14Dからの読出が完了すると、ターゲットノードメモリコントローラ52はデータを含むRdResponseパケットを生成し、パケットをソースノード50に送信する。

#### 【0055】

プローブされたノード54A - 54Nの各々はそのキャッシュ内を探索して、RdBlkパケットによって読出されたキャッシュブロックがその中にストアされているかどうか判断する。ヒットが検出されると、対応するプローブされたノード54A - 54Nは、ターゲットノードメモリコントローラ52から受取ったプローブパケット内のネクストステートフィールドに従ってキャッシュブロックの状態を更新する。さらに、プローブされたノード54A - 54Nの各々は(Probe/Srcパケットが既に受信されているために)、ProbeRespパケットをソースノード50に経路制御する。この例においては、いずれのプローブされたノード54A - 54Nもキャッシュブロックのダーティコピーをストアしていない。

#### 【0056】

ソースノード50は、プローブされたノード54A - 54NからのProbeRespパケットおよびターゲットメモリコントローラ52からのRdResponseパケットを待機する。一旦これらのパケットが受取られると、ソースノード50はSrcDoneパケットをターゲットメモリコントローラ52に送信し、トランザクションを終了する。

#### 【0057】

次いで図12は、第2の例示的な読出ブロックトランザクションを例示する図を示す。ソースノード50、ターゲットメモリコントローラ52、プローブされたノード54A - 54Nが示される。ソースノード50、ターゲットメモリコントローラ52を含むターゲットノード、およびプローブされたノード54A - 54Nの各々は、図1に示す処理ノード12A - 12Dと同様の処理ノードを含み得る。

#### 【0058】

図11に示す例と同様に、ソースノード50はRdBlkパケットをターゲットメモリコントローラ52に送信する。ターゲットノードメモリコントローラ52はProbe/Srcパケットをプローブされたノード54A - 54Nに送信し、かつRdResponseパケットをソースノード50に送信し得る。

#### 【0059】

図12の例においては、プローブされたノード54Aは読出ブロックトランザクションによってアクセスされたキャッシュブロックに対してダーティデータを検出する。したがって、プローブされたノード54Aは(Probe/Srcコマンドにより命令されて)ソースノード50に、RdResponseパケットとプローブされたノード54Aの内部キャッシュから読出されたダーティキャッシュブロックとを送信する。したがって、一実施例においては、プローブされたノード54AはMemCancelパケットをターゲットノードメモリコントローラ

10

20

30

40

50

5 2 に送信し得る。ターゲットノードメモリコントローラ 5 2 が RdResponse パケットをソースノード 5 0 に送信する前に、MemCancel パケットがターゲットノードメモリコントローラ 5 2 に到達すると、ターゲットノードメモリコントローラ 5 2 は RdResponse パケットを送信しない。よって、ターゲットノードメモリコントローラ 5 2 からソースノード 5 0 への「RdResponse」と表示された線は、その選択的な性質を示すために点線で表わされる。MemCancel メッセージに応答して、ターゲットノードメモリコントローラ 5 2 は TgtDone パケットをソースノード 5 0 に送信する。

【 0 0 6 0 】

プローブされたノード 5 4 B - 5 4 N は、この実施例においてはダーティデータを検出せず、よって ProbeResp パケットをソースノード 5 0 に経路制御する。一旦ソースノード 5 0 が TgtDone、RdResponse、および ProbeResp パケットを受取ると、ソースノード 5 0 は SrcDone パケットをターゲットメモリコントローラ 5 2 に送信して読出ブロックトランザクションを終了させる。

【 0 0 6 1 】

図 1 3 は、例示的なサイズ指定された書込トランザクションを示す図である。ソースノード 5 0、ターゲットノードメモリコントローラ 5 2、およびプローブされたノード 5 4 A - 5 4 N を示す。ソースノード 5 0、ターゲットノードメモリコントローラ 5 2 を含むターゲットノード、およびプローブされたノード 5 4 A - 5 4 N の各々は、図 1 に示す処理ノード 1 2 A - 1 2 D と同様の処理ノードを含み得る。

【 0 0 6 2 】

ソースノード 5 0 は、WrSized パケットおよび書込まれるべきデータをターゲットノードメモリコントローラ 5 2 に送信することにより、サイズ指定された書込トランザクションを開始する。サイズ指定された書込トランザクションは、キャッシュブロックの一部を更新するが、キャッシュブロックの残りの部分は更新しないために、ターゲットノードメモリコントローラ 5 2 はプローブされたノード 5 4 A - 5 4 N から（もし存在すれば）ダーティキャッシュブロックを収集する。さらに、キャッシュブロックのクリーンなコピーはプローブされたノード 5 4 A - 5 4 N において無効化され、コヒーレンスを維持する。ターゲットメモリコントローラ 5 2 は、処理されるべきサイズ指定された書込トランザクションを選択する際に、Probe/Tgt パケットをプローブされたノード 5 4 A - 5 4 N に送信する。プローブされたノード 5 4 A - 5 4 N は、（Probe/Tgt パケットが受取られているために）ターゲットノードメモリコントローラ 5 4 に、（いずれのダーティデータも検出されなければ）ProbeResp パケットか、または（ダーティデータが検出されると）RdResponse パケットのいずれかを返す。一旦ターゲットノードメモリコントローラ 5 2 が、プローブされたノード 5 4 A - 5 4 N からの応答を受取ると、ターゲットノードメモリコントローラ 5 2 は TgtDone パケットをソースノード 5 0 に送信し、これはサイズ指定された書込トランザクションを終了させる SrcDone パケットで応答する。

【 0 0 6 3 】

ターゲットノードメモリコントローラ 5 2 が、プローブされたノード 5 4 A - 5 4 N のうちの 1 つからダーティなキャッシュブロックを受取ると、ターゲットノードメモリコントローラ 5 2 は、ダーティキャッシュブロックと WrSized データパケット内のソースノード 5 0 によって提供されたバイトとのマージを実行する。マージを実行するためのいかなる好適な機構をも用い得る。たとえば、ターゲットノードメモリコントローラ 5 2 はデータをマージし、単一ブロック書込を行なってメモリを更新してもよい。これに代えて、ダーティブロックが第 1 にメモリに書込まれ、次いでソースノード 5 0 によって提供されたバイトの書込が続いてもよい。

【 0 0 6 4 】

この説明は、ノード間で通信されるパケットについて記載する一方で、コマンド、応答および他のメッセージを送信するためのいかなる好適な機構をも用い得ることに留意されたい。

【 0 0 6 5 】

10

20

30

40

50



図 1 4 は、処理のためのトランザクションの選択にตอบสนองする、メモリコントローラ 1 6 A - 1 6 D の一実施例の一部の動作を例示するフローチャートを示す。特に、プローブを生成するメモリコントローラ 1 6 A - 1 6 D の部分が示される。他の実施例も可能であり企図される。図 1 4 において、理解を容易にするためにステップを特定の順序に従って示すが、いかなる好適な順序をも用い得る。さらに、ステップはデザイン選択によって所望のように、メモリコントローラ 1 6 A - 1 6 D 内で並列ハードウェアを用いてステップを並行して行なってもよい。

【 0 0 6 6 】

メモリコントローラは、選択されたトランザクションがWrSizedトランザクションであるかどうか判断する（判断ブロック 6 0）。もし選択されたトランザクションがWrSizedトランザクションであれば、メモリコントローラはProbe/Tgtパケットをプローブノードの各々に送信する（ステップ 6 2）。そうでなければ、メモリコントローラは選択されたトランザクションがVicBlkまたはCleanVicBlkトランザクションであるかどうか判断する（判断ブロック 6 4）。選択されたトランザクションがVicBlkまたはCleanVicBlkトランザクションであれば、プローブパケットは生成されない。しかしながら、選択されたトランザクションがWrSized、VicBlk、またはCleanVicBlkでなければ、Probe/Srcパケットはプローブされたノードに送信される（ステップ 6 6）。

【 0 0 6 7 】

図 1 5 は、プローブパケットにตอบสนองするプローブされたノードの一実施例の動作を示すフローチャートである。他の実施例も可能であり企図される。図 1 5 において、理解を容易にするためにステップを特定の順序に従って示すが、いかなる好適な順序をも用い得る。さらに、設計選択によって所望のように、プローブされたノード内で並列ハードウェアを用いてステップを並行して行なってもよい。

【 0 0 6 8 】

プローブされたノードはそのキャッシュを探索して、プローブによって表示されたキャッシュブロックがその中にストアされているかどうか判断し、もしキャッシュブロックが見出されればその状態を判断する。キャッシュブロックがダーティな状態で見出されると（判断ブロック 7 0）、プローブされたノードはRdResponseパケットを生成する。プローブされたノードはRdResponseパケットの後にキャッシュからダーティデータを読み出して、データパケットとして送信する（ステップ 7 2）。一方で、キャッシュブロックが見出されないか、またはダーティな状態になれば、プローブされたノードはProbeRespパケットを生成する（ステップ 7 4）。さらに、キャッシュブロックの状態はプローブパケットのネクストステートフィールドに特定されるように更新される。

【 0 0 6 9 】

プローブされたノードは受取ったプローブパケットを調べる（判断ブロック 7 6）。プローブパケットがProbe/Srcパケットであれば、プローブされたノードは上で生成された応答を、Probe/Srcパケットに表示されるソースノードに経路制御する（ステップ 7 8）。言い換えると、プローブされたノードは応答パケット内のRespNodeフィールドを、Probe/SrcパケットのSrcNodeフィールド内の値にセットする。一方、プローブパケットがProbe/Tgtパケットであれば、プローブされたノードは上で生成された応答を、Probe/Tgtパケットに表示されるターゲットノードに経路制御する（ステップ 8 0）。言い換えると、プローブされたノードは応答パケット内のRespNodeフィールドを、Probe/TgtパケットのTgtNodeフィールド内の値にセットする。

【 0 0 7 0 】

図 1 6 は、例示的な処理ノード 1 2 A の一実施例のブロック図を示す。他の実施例も可能であり企図される。図 1 6 の実施例においては、処理ノード 1 2 A はインターフェイスロジック 1 8 A、1 8 B、1 8 C およびメモリコントローラ 1 6 A を含む。さらに、処理ノード 1 2 A は、プロセッサコア 9 2、キャッシュ 9 0、コヒーレンシ管理ロジック 9 8、および選択により第 2 のプロセッサコア 9 6 および第 2 のキャッシュ 9 4 を含む得る。インターフェイスロジック 1 8 A - 1 8 C は互いに結合され、さらにコヒーレンシ管理ロジ

10

20

30

40

50

ック 9 8 に結合される。プロセッサコア 9 2 および 9 6 は、それぞれキャッシュ 9 0 および 9 4 に結合される。キャッシュ 9 0 および 9 4 は、コヒーレンシ管理ロジック 9 8 に結合される。コヒーレンシ管理ロジック 9 8 はメモリコントローラ 1 6 A に結合される。

【 0 0 7 1 】

一般的に、コヒーレンシ管理ロジック 9 8 は、メモリコントローラ 1 6 A によって処理のために選択されたトランザクションにตอบสนองしてプローブコマンドを生成し、かつ処理ノード 1 2 A によって受取られるプローブコマンドにตอบสนองするよう構成される。コヒーレンシ管理ロジック 9 8 は、処理のために選択されたトランザクションのタイプに依存して、Probe/Src コマンドまたは Probe/Tgt コマンドのいずれかをブロードキャストする。さらに、コヒーレンシ管理ロジック 9 8 は、受取ったプローブコマンドによって特定されるキャッシュブロックに対してキャッシュ 9 0 および 9 4 を探索し、適切なプローブตอบสนองを生成する。さらに、Probe/Tgt コマンドを生成する場合には、コヒーレンシ管理ロジック 9 8 は、Probe/Tgt コマンドにตอบสนองして返されたプローブตอบสนองを収集し得る。キャッシュ 9 0 および 9 4 は、処理ノード 1 2 A 内で開始された読出要求からのデータの受取を管理するフィルロジックを含むか、またはコヒーレンシ管理ロジック 9 8 内にフィルロジックが含まれてもよい。コヒーレンシ管理ロジック 9 8 はさらに、メモリコントローラ 1 6 A に対する非コヒーレント要求を経路制御するよう構成されてもよい。一実施例においては、プロセッサ 9 2 とプロセッサ 9 6 とは、キャッシュ 9 0、キャッシュ 9 4、およびコヒーレンシ管理ロジック 9 8 をバイパスして、特定のキャッシュ不可および/または非コヒーレントメモリ要求に対して直接メモリコントローラ 1 6 A にアクセスし得る。

【 0 0 7 2 】

キャッシュ 9 0 および 9 4 は、データのキャッシュブロックをストアするよう構成された高速キャッシュメモリを含む。キャッシュ 9 0 および 9 4 は、それぞれのプロセッサコア 9 2 および 9 6 内に統合されてもよい。これに代えて、キャッシュ 9 0 および 9 4 は、所望のようにプロセッサコア 9 2 および 9 6 に、バックサイドキャッシュ構成またはインライン構成で結合されてもよい。さらに、キャッシュ 9 0 および 9 4 はキャッシュ階層構造として実現化されてもよい。プロセッサコア 9 2 および 9 6 に（階層内で）より近いキャッシュは、所望であればプロセッサコア 9 2 および 9 6 に統合されてもよい。

【 0 0 7 3 】

プロセッサコア 9 2 および 9 6 は、予め定められた命令セットに従って命令を実行するための回路を含む。たとえば、x 8 6 命令セットアーキテクチャを選択し得る。これに代えて、Alpha、PowerPC、または他のいかなる命令セットアーキテクチャを選択してもよい。一般的に、プロセッサコアはデータおよび命令のためにキャッシュにアクセスする。キャッシュミスが検出されると、読出要求が生成され、ミスの発生したキャッシュブロックがマッピングされているノード内のメモリコントローラに送信される。

【 0 0 7 4 】

分散メモリシステムを用いて、コンピュータシステム 1 0 の特定の実施例を説明してきたが、分散メモリシステムを用いない実施例も、ここで説明した柔軟なプローブ/プローブตอบสนอง経路制御を用い得ることに留意されたい。そのような実施例が企図される。

【 0 0 7 5 】

上の開示を完全に理解すると、当業者においてはいくつもの変更および変形が明らかとなるであろう。前掲の特許請求の範囲がすべてのそのような変更および変形を包含することが理解されることを意図する。

【 0 0 7 6 】

【産業上の用途】

この発明は一般的にコンピュータシステムに適用可能である。

【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。

【図 2】 図 1 に示す 1 対の処理ノードの間の相互接続の一実施例を強調した、ブロック図である。

10

20

30

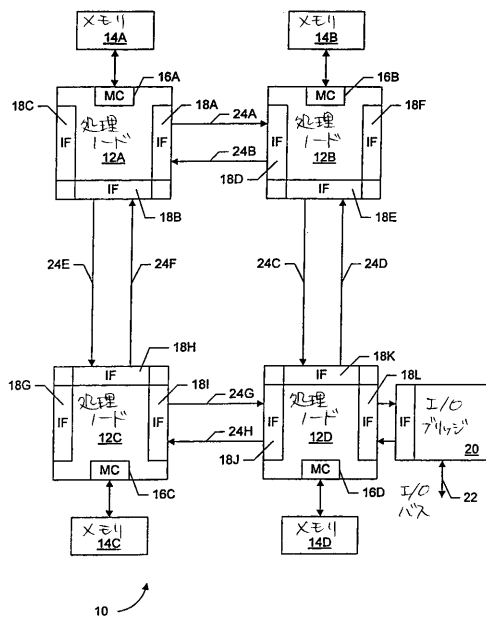
40

50

- 【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】 データパケットの一実施例のブロック図である。
- 【図 7】 コンピュータシステムの一実施例によって用い得る、例示的なパケットタイプを示すテーブルである。
- 【図 8】 プローブパケットの一実施例のブロック図である。
- 【図 9】 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 10】 読出応答パケットの一実施例のブロック図である。
- 【図 11】 読出ブロックトランザクションに対応するパケットの例示的なフローを示す図である。
- 【図 12】 読出ブロックトランザクションに対応するパケットの第 2 の例示的なフローを示す図である。
- 【図 13】 サイズ指定された書込トランザクションに対応するパケットの例示的なフローを示す図である。
- 【図 14】 メモリコントローラの一実施例の動作を示すフローチャートの図である。
- 【図 15】 プローブパケットを受取る処理ノードの一実施例の動作を示すフローチャートの図である。
- 【図 16】 処理ノードの一実施例のブロック図である。

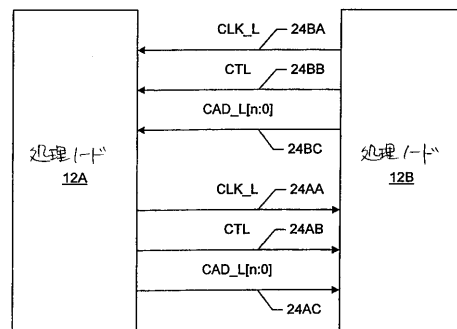
10

【図 1】



10

【図 2】



【図 3】

| ビット時間 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0        |
|-------|---|---|---|---|---|---|---|----------|
| 1     |   |   |   |   |   |   |   | CMD[5:0] |
| 2     |   |   |   |   |   |   |   |          |

30

【図 4】

| ビット時間 | 7                 | 6 | 5            | 4           | 3 | 2 | 1                 | 0 |
|-------|-------------------|---|--------------|-------------|---|---|-------------------|---|
| 1     | DestNode<br>[1:0] |   | CMD[5:0]     |             |   |   |                   |   |
| 2     | SrcTag<br>[1:0]   |   | SrcNode[3:0] |             |   |   | DestNode<br>[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]       |   |              |             |   |   |                   |   |

32

【図 5】

| ビット時間 | 7                 | 6 | 5            | 4           | 3 | 2                 | 1 | 0 |
|-------|-------------------|---|--------------|-------------|---|-------------------|---|---|
| 1     | DestNode<br>[1:0] |   | CMD[5:0]     |             |   |                   |   |   |
| 2     | SrcTag<br>[1:0]   |   | SrcNode[3:0] |             |   | DestNode<br>[3:2] |   |   |
| 3     |                   |   |              | SrcTag[6:2] |   |                   |   |   |
| 4     |                   |   |              |             |   |                   |   |   |

34

【図 6】

| 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] |   |   |   |   |   |   |   |

36

【図 7】

| CMD 符号 | コマン                 | パケットタイプ  |
|--------|---------------------|----------|
| 000000 | Nop                 | Info     |
| 000001 | Interrupt Broadcast |          |
| 00001x | Reserved            | -        |
| 000100 | Probe/Src           | Address  |
| 000101 | Probe/Tgt           | Address  |
| 00011x | Reserved            | -        |
| 001000 | RdBlkS              | Address  |
| 001001 | RdBlk               | Address  |
| 001010 | RdBlkMod            | Address  |
| 001011 | ChangetoDirty       | Address  |
| 001100 | ValidateBlk         | Address  |
| 001101 | CleanVicBlk         | Address  |
| 001110 | Interrupt Target    |          |
| 001111 | VicBlk              | Address  |
| 01xxxx | ReadSized           | Address  |
| 10xxxx | WrSized             | Address  |
| 11000x | RdResponse          | 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     |

38

【図 8】

| ビット時間 | 7                  | 6 | 5            | 4           | 3 | 2 | 1                | 0 |
|-------|--------------------|---|--------------|-------------|---|---|------------------|---|
| 1     | TgtNode<br>[1:0]   |   | CMD[5:0]     |             |   |   |                  |   |
| 2     | SrcTag<br>[1:0]    |   | SrcNode[3:0] |             |   |   | TgtNode<br>[3:2] |   |
| 3     | NextState<br>[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]        |   |              |             |   |   |                  |   |

40

【図 9】

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

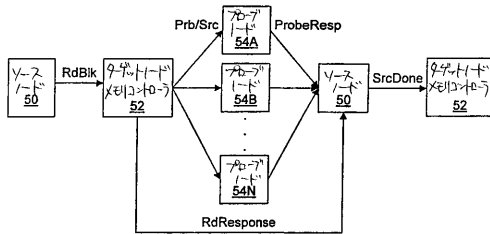
42

【図 10】

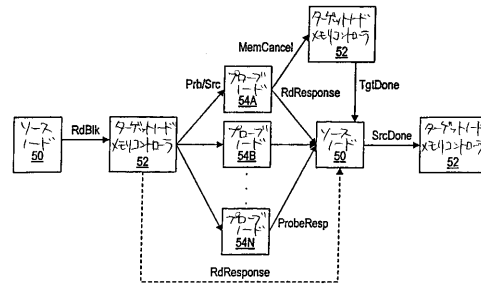
|       |                   |   |              |             |   |   |                   |      |
|-------|-------------------|---|--------------|-------------|---|---|-------------------|------|
| ビット番号 | 7                 | 6 | 5            | 4           | 3 | 2 | 1                 | 0    |
| 1     | RespNode<br>[1:0] |   | CMD[5:0]     |             |   |   |                   |      |
| 2     | SrcTag<br>[1:0]   |   | SrcNode[3:0] |             |   |   | RespNode<br>[3:2] |      |
| 3     | Count             |   |              | SrcTag[6:2] |   |   |                   |      |
| 4     | Rsv               |   |              |             |   |   | Prb               | Type |

44

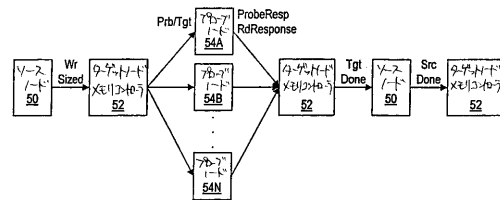
【図 11】



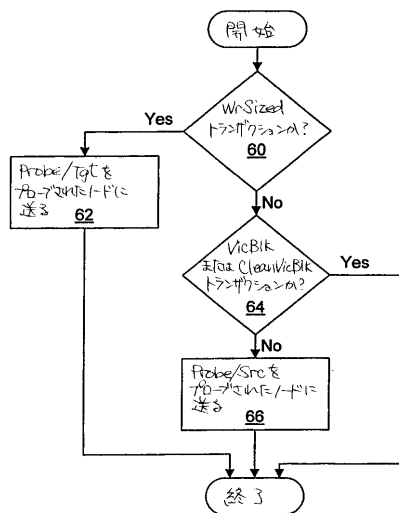
【図 12】



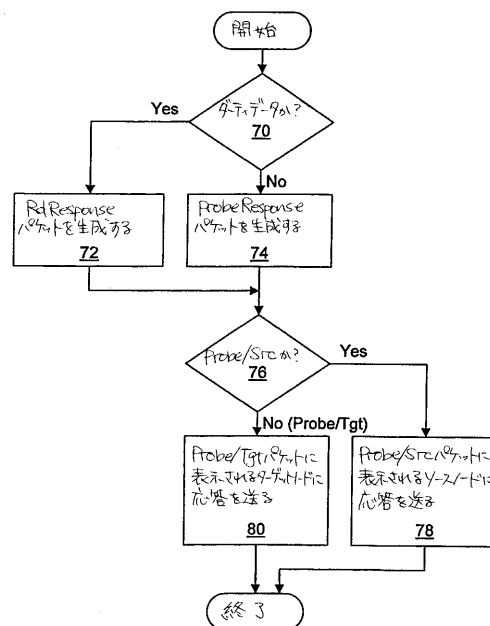
【図 13】



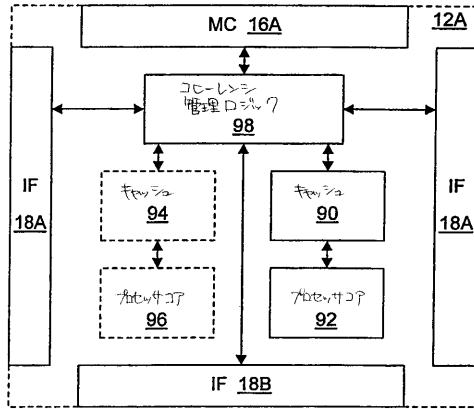
【図 14】



【図 15】



【図 16】



## フロントページの続き

- (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・ビー  
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、アイリス・ウェイ、2 1 0
- (72)発明者 ギューリック, デイル・イー  
アメリカ合衆国、7 8 7 4 8 テキサス州、オースティン、フェスタス・ドライブ、3 1 2 2

審査官 高瀬 勤

- (56)参考文献 特開平 1 0 - 1 0 5 4 6 4 ( J P , A )  
特開平 0 6 - 2 7 4 4 6 1 ( J P , A )  
特開平 1 0 - 3 4 0 2 2 7 ( J P , A )  
特開平 0 8 - 3 2 0 8 2 7 ( J P , A )  
特開平 0 8 - 0 1 6 4 7 4 ( J P , A )  
特表 2 0 0 2 - 5 3 3 8 1 3 ( J P , A )  
Akhilesh Kumar, Phanindra K. Mannava, Laxmi N. Bhuyan, Efficient and Scalable Cache Coherence Schemes for Shared Memory Hypercube Multiprocessors, Proceedings of Supercomputing'94, IEEE, 1994年11月14日, Pages:498-507  
James Laudon, Daniel Lenoski, The SGI Origin: A ccNUMA Highly Scalable Server, Proceedings of the 24th Annual International Symposium on Computer Architecture (ISCA'97), ACM, 1997年 6月, Pages:241-251

- (58)調査した分野(Int.Cl., D B 名)  
G06F 12/08