

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3981192号

(P3981192)

(45) 発行日 平成19年9月26日(2007.9.26)

(24) 登録日 平成19年7月6日(2007.7.6)

(51) Int. Cl.

F I

G 0 6 F 1/12 (2006.01)

G 0 6 F 1/04 3 4 0 A

G 0 6 F 13/00 (2006.01)

G 0 6 F 13/00 3 5 1 C

請求項の数 3 (全 13 頁)

(21) 出願番号	特願平9-264326	(73) 特許権者	398038580
(22) 出願日	平成9年9月29日(1997.9.29)		ヒューレット・パカード・カンパニー
(65) 公開番号	特開平10-161984		HEWLETT-PACKARD COMPANY
(43) 公開日	平成10年6月19日(1998.6.19)		アメリカ合衆国カリフォルニア州パロアルト
審査請求日	平成16年9月27日(2004.9.27)		ハノーバー・ストリート 3000
(31) 優先権主張番号	720332	(74) 代理人	100087642
(32) 優先日	平成8年9月27日(1996.9.27)		弁理士 古谷 聡
(33) 優先権主張国	米国 (US)	(74) 代理人	100063897
			弁理士 古谷 馨
		(74) 代理人	100076680
			弁理士 溝部 孝彦

最終頁に続く

(54) 【発明の名称】 S C I 相互接続を用いた T O C カウンタの同期

(57) 【特許請求の範囲】

【請求項 1】

複数の S C I リングによって相互接続された複数のノードを備えた、同期処理イベントに関する低スキュー時間カウンタを有するマルチプロセッサコンピュータシステムであって、前記ノードの各々が、

時間値を生成する時間カウンタと、

前記ノードと前記 S C I リングとのインタフェースをとる S C I コントローラと、及び前記 S C I コントローラと前記時間カウンタとの間で同期信号を搬送する同期信号配信経路とを備えており、

前記複数のノードのうちの1つがマスタノードとして指定され、該マスタノードが、前記配信経路を介して前記マスタノード上の前記時間カウンタ及び前記 S C I コントローラに配信される同期信号を生成するパルス同期信号生成器を備えており、前記 S C I コントローラが、前記同期信号を受信して利用可能なデータパケットを見つけ、前記同期信号に対応する時間カウンタ同期ビットをデータパケット中にセットし、セットされた該ビットを含むデータパケットが、前記 S C I リングを介して該コンピュータシステムの残りのノード上における S C I コントローラへ配信され、前記残りのノード上における各 S C I コントローラが、受信したデータパケットから前記同期信号を取り出して、その同期信号を各残りのノード上の時間カウンタへ前記配信経路を介して配信して、前記コンピュータシステム内の前記時間カウンタの各々の時間値が同期される、マルチプロセッサコンピュータシステム。

10

20

【請求項 2】

ＳＣＩデータパケットを使用するマルチプロセッサコンピュータシステムであって、そのコンピュータシステムが、複数のＳＣＩリングによって相互接続された、マスタノードを含む複数のノードを備え、前記マスタノードが、時間値を生成する時間カウンタと、前記ノードと前記ＳＣＩリングとのインタフェースをとるＳＣＩコントローラと、前記ＳＣＩコントローラと前記時間カウンタとの間で同期信号を搬送する同期信号配信経路とを備えており、前記ＳＣＩデータパケットが、

時間カウンタ同期ビットとデータ記号とを含むヘッダを備え、

前記マスタノードが、前記配信経路を介して前記マスタノード上の前記時間カウンタ及び前記ＳＣＩコントローラに配信される同期信号を生成するパルス同期信号生成器を含み、前記ＳＣＩコントローラが、前記同期信号を受信して利用可能なデータパケットを見つけ、データパケット中に前記同期信号に対応する前記時間カウンタ同期ビットをセットし、セットされた該ビットを含むデータパケットが、前記ＳＣＩリングを介して該コンピュータシステムの残りのノード上におけるＳＣＩコントローラへ配信され、前記残りのノード上における各ＳＣＩコントローラが、受信したデータパケットから前記同期信号を取り出して、その同期信号を各残りのノード上の時間カウンタへ前記配信経路を介して配信して、前記コンピュータシステムの各時間カウンタの時間値が同期される、ＳＣＩデータパケットを使用するマルチプロセッサコンピュータシステム。

10

【請求項 3】

複数のＳＣＩリングによって相互接続された複数のノードを備え、前記ノードの各々が、時間値を生成する時間カウンタと、前記ノードと前記ＳＣＩリングとのインタフェースをとるＳＣＩコントローラと、前記ＳＣＩコントローラと前記時間カウンタとの間で同期信号を搬送する同期信号配信経路とを備えている、マルチプロセッサコンピュータシステムにおいて、各ノードにおける時間カウンタの時間値の同期を取る方法であって、

20

(a) 複数のノードのうちの１つをマスタノードとして指定し、該マスタノードが、前記配信経路を介して前記マスタノード上の前記時間カウンタ及び前記ＳＣＩコントローラに配信される同期信号を生成するパルス同期信号生成器を備え、

(b) 前記同期信号を前記マスタノード上の前記ＳＣＩコントローラに送り、前記ＳＣＩコントローラが利用可能なデータパケットを見つけて、データパケット中に前記同期信号に対応する時間カウンタ同期ビットをセットし、及び

30

(c) 前記セットされた該ビットを含むデータパケットを、前記ＳＣＩリングを介して前記コンピュータシステムの残りのノード上におけるＳＣＩコントローラへ配信し、前記残りのノード上における各ＳＣＩコントローラは、受信したデータパケットから前記同期信号を取り出して、その同期信号を各残りのノード上の時間カウンタに前記配信経路を介して配信することを含む、方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は一般にマルチプロセッサシステムに関し、特にマルチプロセッサシステムの同期をとる方法およびシステムに関する。

40

【0002】

【従来の技術】

マルチノード・マルチプロセッサシステムの同期をとるには、システム内の各プロセッサのクロックを他のプロセッサのクロックに相対的に同期させなければならない。これを行うには、一定の処理ポイントにおける概略時間が分からなければならない、システム全体にわたって時間がほぼ同じである必要がある。

【0003】

かかるシステムでは、各ノードがクロックカウンタを有し、該ノード上の各プロセッサが該クロックカウンタを読み取る。残念ながら、かかるシステム内の各カウンタクロックひいては各ノードは、僅かに異なるクロック周波数で動作する。このクロック周波数の差は

50

、各カウンタ内のクリスタルが厳密に同一でないことによるものである。この異なるクリスタル周波数のため、カウンタの時間値がドリフトすることになる。クリスタルの物理的な差を制御することはできない。

【 0 0 0 4 】

既知の従来技術は、各ノード間に接続された余剰ワイヤを使用することによりこのドリフトの問題を解決している。かかるワイヤは、同期信号を搬送する別々の信号経路を形成する。ハードウェアにより画定される一定時間の経過後に、該ワイヤを介して同期パケットが配信され、次いで各ノードが該同期パケットを受信して、そのカウンタ時間を適切に変更する。

【 0 0 0 5 】

この従来技術の解決策に関する問題は、余剰信号経路のために、この解決策がコスト/パフォーマンスの点で高価になることにある。この従来技術の解決策はまた、システムの複雑さを増大させるものとなる。これは、回路に（特に接続部間の接地に関して）より多くの問題および誤りを生じさせる恐れのある追加的な接続をワイヤが必要とするからである。

【 0 0 0 6 】

【発明が解決しようとする課題】

したがって、当業界では、互いに同期された異なるノード上の低スキュークロックへのアクセスを提供するためのシステムおよび方法が必要とされている。

【 0 0 0 7 】

また、当業界では、同期中に待ち時間を生じさせることなく異なるノード上のクロックを同期させるためのシステムおよび方法も必要とされている。

【 0 0 0 8 】

更に、当業界では、システム性能を低下させることなく異なるノード上のクロックを同期させるためのシステムおよび方法も必要とされている。

【 0 0 0 9 】

【課題を解決するための手段】

上述その他の必要性は、マルチプロセッサシステムが低スキュークロックにアクセスして処理イベントを同期させるシステムおよび方法によって満たされる。本発明は、S C I又はスケーリング可能なコヒーレント相互接続ネットワーク等の既存のハードウェアを用いて低スキュー信号を配信して、異なるノード上のタイム・オブ・センチュリー・クロック(time of century clock: 百年制クロック)を同期させる。これらのカウンタを、選択されたマスタカウンタからの信号と周期的に同期させることにより、全てのノードがほぼ等しいカウンタ値を維持することになる。送信パケット、エコーパケット、及びアイドルパケットのS C Iヘッダ中の単一のビットが、S C Iリングを介して全てのノードに経路指定される。該ビットが既存のパケットまたはルーチンパケットに挿入されるので、特殊な同期パケットを作成する必要がない。更に、該ビットは、既存の線を介して移動するので、追加的な信号経路または余剰ワイヤが不要となる。

【 0 0 1 0 】

本発明の技術的な利点は、S C Iを使用してシステム上の全てのクロックへ同期パルスを送信することにある。

【 0 0 1 1 】

本発明の別の技術的な利点は、既存のデータパケットを使用して同期パルスを伝送することにある。

【 0 0 1 2 】

本発明の更に別の技術的な利点は、既存のデータパケットのヘッダに同期パルスを配置することにある。

【 0 0 1 3 】

上記説明は、下記の本発明の詳細な説明が一層良好に理解されるように本発明の特徴および技術的な利点をかなり広範に概説したものである。特許請求の範囲に記載の本発明の要

10

20

30

40

50

旨を形成する本発明の他の特徴および利点について以下で説明する。当業者であれば、本発明と同一の目的を達成するために、本開示の概念および特定の実施形態を、修正あるいは別構造の設計のための基礎として容易に利用可能であることが理解されよう。当業者であれば、そのような等価的な構成が、特許請求の範囲に記載の本発明の思想および範囲から逸脱しないものであることもまた理解されよう。

【0014】

【発明の実施の形態】

図1は、システム中の合計112個の考え得るノードのうちの2つのノード、具体的にはノード0およびノード1を概略的に示すブロック図である。また、図2は、単一のノードの構成要素を示すブロック図である。それぞれの異なるノードをクラスタ状に区画してシステムの耐久生存性(survivability)を向上させることができる。これについては、1996年9月27日出願の「ERROR CONTAINMENT CLUSTER OF NODES」と題する同時係属中の米国特許出願第08/720368号に記載されている。なお、本引用をもってその開示内容を本明細書中に包含させたものとし、その詳細な説明は省略する。

【0015】

マルチプロセッサコンピュータシステムは、2つのノードを有することができ、また最大で112個のノードを有することができる。図3に示すように、112ノードシステムにおいて、ノード24は、壁23を形成する7つのX次元リング26×4つのY次元リング27として構成される。かかる4つの壁が4つのZ次元リング28によって相互接続される。Y次元リング27をZ次元リング28に接続するためにブリッジノードが使用される。

【0016】

プロセッサエージェントチップを表すPACと記した1つのプロセッサエージェント11には、最大2つのプロセッサ10を接続することができる。単一のノードは、最大で8つのPAC11を有することができる。多数の同一要素が存在することに留意されたい。説明を明瞭にするため、本明細書では同一要素を単一の符号で示すこととする。なお、2つ以上の同一要素を区別する場合には、異なる要素に新たな符号を付することとする。

【0017】

プロセッサ10は、HEWLETT-PACKARD PA-8000プロセッサであることが好ましい。しかしながら、本発明はプロセッサタイプやアーキテクチャにより制限されるものではない。プロセッサ10は、ランウェイバスを介してPAC11に接続される。PAC11は、入出力(I/O)サブシステムを有し、クロスバー12およびコア論理アクセスバスに接続される。コア論理アクセスバスは主としてシステムブート動作に使用される。このバスは、全てのPACを、消去可能なプログラマブル読み出し専用メモリ(EPROM)、同期ダイナミックランダムアクセスメモリ(SDRAM)、リアルタイムクロック、RS-232インタフェース、及びEthernetインタフェースに結合させる、低帯域幅マルチドロップバスである。更に、プロセッサは、バスを使用してアクセスされる制御状態レジスタ(CSR)に書き込みを行って、クロスバーの初期設定及び構成を行うことができる。

【0018】

PAC11の機能は、要求をプロセッサ10からクロスバー12を介してメモリアクセスシステム14へ送信し、次いでその応答を要求側プロセッサ10へ送り返すことである。各PAC11内には、TOCと記したタイム・オブ・センチュリー・カウンタ13がある。各PACが2つのプロセッサを処理し、1ノード内に最大で8つのPACがあるので、各ノードは最大で16個のプロセッサを有することができる。図2は、4つのクロスバー12を示しているが、各PACは、そのうちの2つのクロスバーと通信する。

【0019】

PACは、4つの一方向データ経路を用いてクロスバー12を介してメモリコントローラ14と通信する。RAC(routing attachment chip: 経路指定接続チップ)と記したクロスバー12は、エージェント11からパケットを受信し、次いで該パケットをMACと記したメモリアクセスコントローラ14に経路指定する経路指定手段である。各RACは16個の32ビット幅の一方向相互接続手段を有しており、その各RACは4つのPACおよび4つのMA

10

20

30

40

50

Cに接続される。クロスバーは、それ自体のCSRを有さず、その代わりに、コアアクセス論理バス上にあるCSRへの書込みによって初期設定される。それらのCSRは、どのポートをアクティブにするかを制御すると共にエラー検出をイネーブルにする。

【0020】

MAC14は、コヒーレントメモリへのアクセスを制御する。メモリアクセスコントローラは、2の倍数で2から8までの番号付けを行うことができ、各MACは、4つのバンクの最大2Gbyte（各バンク29が512Mバイトを有する）までサポートする。したがって、各ノードは最大16Gbyteまでアクセスすることができ、28ノードシステムは最大448Gバイトまでにアクセスすることができる。メモリバンクは、同期DRAMまたはSDRAMからなるSIMMを備える。図2は、図示の簡素化のために2つのメモリバンク29のみを示したものである。該メモリは、ノードローカルメモリ、ネットワークキャッシング、及びメッセージングに使用される。キャッシュのコヒーレンシーを維持する方法については、1996年9月27日出願の「METHOD AND SYSTEM FOR MAINTAINING STRONG ORDERING IN A COHERENT MEMORY SYSTEM」と題する同時係属中の米国特許出願第08/720330号で論じられている。なお、本引用を持ってその開示内容を本明細書中に包含させたものとし、その詳細な説明は省略する。

10

【0021】

プロセッサ10がメモリその他のリソースにアクセスするための要求を生成すると、PAC11は、要求されたアドレスを調べて、該要求の処理に適当なMACを決定し、次いで該要求をRAC12を介して適当なMAC14へと送信する。MAC14は、ノードIDがローカルメモリアドレスに対するものでないと判定した場合には、該要求をTACと記したリングインタフェースコントローラ15へ送る。また、MAC14は、要求されたアドレスがローカルノードについてのものであると判定した場合には、そのMAC14に接続されているメモリ29にアクセスする。

20

【0022】

TAC15は、トロイダルアクセスチップまたはSCIコントローラとしても知られるものである。TACは、ノードからSCIリングへのインタフェースとして働く。TACは、2つの一方向データ経路を使用してMACと通信する。各TACは、2つのSCIリング、即ち、X次元リングおよびY次元リングとのインタフェースをとる。図1は、図示の簡素化のため単一の次元のみを示している。図1はまた、リング16とのインタフェースをとる1つのTAC15と、リング18とのインタフェースをとるもう1つのTAC17とを示している。

30

【0023】

TAC15は、別々のリング16を操作することが可能なものである。最大8つのMAC/TAC対が存在できるので、ノードの各セクションを接続する合計で最大8つのSCIリング（即ち、8つのX次元リングおよび8つのY次元リング）が単一の次元に存在することができる。SCIインタフェースリングについては、「IEEE Standard for Scalable Coherent Interface (SCI)」(IEEE Std.1596-1992 ISBN 1-55937-222-2)で規定されている。尚、本引用をもってその内容を本書中に包含させたものとし、その詳細な説明は省略する。TAC15は、MAC14から非ローカルメモリアクセス要求を受信し、その要求をSCIリング16に送る。図1において、受信側TAC19は、送信側TAC15からの要求を受信し、次いで該要求をそのローカルMAC20へ送る。メモリアクセスが該要求を満たす場合には、その応答が、TAC19、リング16、TAC15、MAC14、RAC12、及びPAC11を介してプロセッサ10に返される。

40

【0024】

各PACプロセッサエージェント内には、TOC13またはタイム・オブ・センチュリー・カウンタと呼ばれる論理構成がある。このカウンタは、ローカルクロック周波数に応じてカウントを行う。PACに取り付けられた各プロセッサは、その2つの異なるプロセッサがほぼ同時にTOCを読み出す場合に各プロセッサがほぼ同じ値に設定され又は少なくとも許容可能な公差限界内に設定されるように、それらのプロセッサ間における比較的等し

50

い待ち時間で該カウンタにアクセスする。各ノードはクリスタルクロックを1つずつ有し、該クロックにより同一ノード上のT O Cが動作する。

【0025】

各ノードが異なるクリスタルを有しているため、異なるノード上で動作するタイム・オブ・センチュリー・カウンタが僅かに異なる周波数で動作するという問題が生じる。T O Cによるカウントの同期を周期的にとり、これにより、異なるノード上のリモートプロセッサがローカルノード上のメモリその他のデバイスに対して読み出しまたはアクセスを行うとき各プロセッサ（ローカルプロセッサおよびリモートプロセッサ）がそれ自体のT O Cに対して読み出しを行う際にほぼ同じ値を読み出すようにする必要がある。

【0026】

各ノードには8つのP A Cがあり、各P A Cは、それ自体のタイム・オブ・センチュリー・カウンタすなわちT O Cを有している。ノード内のこれら8つのP A Cの全てがワイヤ21によって接続される。周期的に、ワイヤ21に沿って同期パルスが送信され、これにより、各P A CがそのT O Cに同期される。同一ノード上の全てのT O Cは同一のクリスタルで動作するので、同一ノード上のT O C間にドリフトは存在しない。

【0027】

ノード内の全てのP A Cの接続を行うワイヤ21はまた、同ノード内の全てのT A Cにも接続される。1つのノード内の1つのT A CはT O Cマスタとして選択される。T O Cマスタのタスクは、同期パルスをS C Iリングを介してそのS C Iリングに接続された全てのノードへ送信することである。同期パルスは、既存のデータパケットのアイドル記号またはヘッダ記号に挿入されるので、該同期パルスを送信するという目的だけのためにデータパケットを作成する場合よりも他のノードに一層高速に到達することができる。更に、同期パルスは、パケットのヘッダ内にあるので、該同期パルスは、パケット内のデータの残りの部分よりも前に作用を受ける。したがって、プロセッサが他のノード上のT O Cの読み出しを行う場合、T O C同期信号が他のパケットよりも高速であるため、それぞれの異なるT O C間のドリフトは知覚されない。

【0028】

図4に示すように、単一のノード30（通常はノード0）がマスタとして指定され、このマスタノードがT O C同期信号を生成し、該T O C同期信号が残りのノードまたはスレーブノード31へ送信される。マスタノード上のP A Cのうちの1つがマスタP A C 11として指定され、該マスタP A C 11がT O C同期信号を生成して該T O C同期信号をマスタノード30上の他の非マスタP A Cへ送信する。これと同時に該T O C同期信号がT A C 15を介してS C Iリング16へ送られ、次いでスレーブノード31へと送られる。T A C 19は、該スレーブノード31上でT O C同期信号を受信し、次いでそのスレーブノード31上の全てのP A Cに接続されたワイヤ22上にそのパルスを送る。したがって、スレーブノード上の全てのP A Cがほぼ同時にT O C同期信号を受信することになる。

【0029】

同期ワイヤ21は、T O C同期信号がマスタノード内の8つのT A C全てに実際に送信されるように全てのP A C及び全てのT A Cを接続する。実際には複数のT A Cのうちの1つだけを使用して他のスレーブノードへ同期信号を送信するが、使用すべきT A Cをソフトウェアが選択することが可能であるため、ハードウェアが故障した場合には、異なるリングを使用するバックアップ用T A Cを選択することが可能であり、したがって故障の修理のために停止させることなく動作を継続させることが可能である。

【0030】

図5は、T O C用のハードウェアを示すものである。T O Cは、システム全体の同期クロックに非常に短い待ち時間でアクセスするための機構を備えている。T O Cを使用して、タイムスタンプ付きトレースデータを後の分析のために生成することができる。各ノードから得たタイムスタンプ付きのトレースデータを後の処理ステップでマージして、5 ~ 10 μ secの範囲のイベントシーケンスを有する正確な大域ピクチャを提供することが可能である。T O Cはまた、送信されたメッセージのタイムスタンプも提供する。受信側は、現

10

20

30

40

50

在時刻からタイムスタンプを減算することによって送信時間を求めることができる。

【 0 0 3 1 】

システム内の各 P A C は T O C 同期パルス生成器 32 を有する。該 T O C 同期パルス生成器 32 は、それがマスタノード上のマスタ P A C でしか使用されない場合であっても設けられる。マスタ T O C 同期パルス生成器は、T O C 同期マスタセクタ 33 によって作動される。全ての残りの P A C 内の T O C 同期セクタは、それらの個々の T O C 同期パルス生成器の選択解除を行い、または T O C 同期パルス生成器をオフにセットする。したがって、1 つの P A C だけがパルスを生じてそのパルスをシステム内の他の全ての P A C に配信することになる。T O C 同期パルス生成器 32 は、ワイヤ 21 を含む配信論理回路 34 へその信号を送信する。T O C 同期信号は、全てのローカル P A C へ送られ、また全てのローカル T A C に送られる。該 T O C 同期信号を S C I リングを介して全てのリモート P A C へ送信するために、1 つの T A C、即ちマスタ T A C が選択される。受信側 T A C 19 は、T O C 同期パルスを受信し、該 T O C 同期パルスをそのノード上の全ての 8 つの P A C に配信する。次いで、各 P A C が該 T O C 同期パルスを受信し、該 T O C 同期パルスを使用してその T O C の再同期を行う。

【 0 0 3 2 】

各ノード上のクリスタルクロック 35 及びクロック生成器 36 は、各 P A C 上の T O C 用の 16 MHz クロックを生成する。P A C は、7 つまたは 8 つの T O C クロック毎にクリスタルクロックをその P A C 自体の T O C に同期させる。マスタ P A C は、256 クロックまたは 16 μ sec 毎に T O C 同期パルスを生成する。

【 0 0 3 3 】

一般に、16 MHz クロック 35 が、プレスケール / シンクロナイザ 37 によりスケールダウンされて、タイムオブセンチュリーレジスタ 38 となる。これは、この特定の P A C 上に配設されたローカルプロセッサにより読み出しが行われるレジスタである。チェック論理回路 39 は、T O C カウンタレジスタが特定の分解能内に同期を維持することを確実にする。配信論理回路 34 は、同期パルス間の時間が、同期周期に対して同期分解能の 1 / 2 を加算又は減算した範囲内であることを確実にするために検査を行う。該分解能は、T O C 同期分解能論理回路 40 により設定される。次の表 1 は、サポートされる幾つかの分解能についての検査範囲を示すものである。

【 0 0 3 4 】

【表 1】

分解能	チェック範囲 (16 μ sec)
1 μ sec	256 \pm 7
2 μ sec	256 \pm 15
4 μ sec	256 \pm 31

【 0 0 3 5 】

チェック論理回路が、進んだパルスまたは遅れたパルスを検出した場合には、その P A C に接続されたプロセッサのうちの 1 つに割り込みが送信される。

【 0 0 3 6 】

プリスケール論理回路 37 は、クロック 35 の 16 による除算を実行し、その結果、1 μ sec の周期信号が得られる。該周期信号は、T O C カウンタレジスタ 38 のインクリメントを可能にするために使用される。レジスタ 38 の同期は、同期パルスの到着時にプリスケール値を切り上げあるいは切り捨てることにより実行される。その丸め量は、T O C 分解能 40 の関数となる。

【0037】

SCIは、パケットベースのプロトコルである。各パケットは、基本的に、ヘッダと、それに続く、パケットのタイプに応じた0～8個のデータ記号とを備える。前記ヘッダは、CLKと記す更なるビットを有する。該ビットはTOC同期ビットである。図6は、ヘッダにCLKビット41を有する典型的なSCIパケットを示すものである。PACは、マスタPACからTOC同期信号を受信すると、修正することができる最初の使用可能なヘッダを見つけて、CLKビットをセットする。リング上の他のあらゆるTACは、そのパケットをCLKビットと共に受信した際に、該CLKビットを取り出し、該ビットをローカルPACへ送り、及びリングを介して次のTACへ送る。最後に、CLKビットは、リングを介して当初のTACまたはマスタTACへと戻され、該マスタTACが、該CLKビ

10

【0038】

CLKは、各パケットに含まれるサイクル冗長コードまたはCRCコードの計算には使用されない。このため、CRCを再計算することなく直ちにCLKビットを変更することが可能となる。CRCはSCIの仕様で規定される。CRCは、本質的には、パケット内の全てのビットの大きなXORであり、最後のパケット内に保存されたものであり、各TACでパケットが受信される際に、新たなCRCが計算されて、送信されたCRCと比較される。2つのCRCが異なる場合にはエラーが発生している。CLKビットはあらゆるヘッダに追加される。このため、TACは、長くとも現在のパケットが終了するのを待った

20

【0039】

各TACは、システム全体にわたり同期パルスを如何に伝搬させるかを制御する制御状況レジスタCSRを有する。CSRは、到来する同期パルスについてのソースを指定する。CSRはまた、同期パルスをSCI X次元リングに伝搬させるかSCI Y次元リングに伝搬させるかを指定する。

【0040】

図7に示すように、TAC TOC構成レジスタは3つのフィールドを有する。ソースフ

ィールド42は、2ビットフィールドであり、イネーブルされた同期パルス出力にどの同期パルス入力（同期信号、またはX到来リンク、またはY到来リンク）を伝搬させるべきかを指定する。該2ビットにより、4つの選択肢、即ち、値0（解決策をとらず、あるいは何も行わない）、値1（PACから信号を取り出して配信する）、値2（二次元リング構造のX入力から信号を取り出して配信する）、及び値3（Y入力から信号を取り出して配信する）が与えられる。最後の2つのフィールド43,44は、ビットを如何に配信するかを指示するものである。XリングビットまたはYリングビットに1がある場合には、そのリング上の最初の使用可能なヘッダにTOC同期信号が配信される。このX-Yレイアウトについては、1996年9月27日出願の「ROUTING METHODS FOR A MULTINODE SCI COMPUTER SYSTEM」と題する同時係属中の米国特許出願第08/720331号に記載されている。

30

40

【0041】

本発明およびその利点について詳細に説明したが、特許請求の範囲に記載の本発明の思想および範囲から逸脱することなく、本開示内容に対して様々な変更、置換、及び修正を加えることが可能であることが理解されよう。

【0042】

以下においては、本発明の種々の構成要件の組み合わせからなる例示的な実施態様を示す。

【0043】

1. 複数のSCIリングによって相互接続された複数のノードを備えた、同期処理イベントに関する低スキュー時間カウンタを有するマルチプロセッサコンピュータシステムであ

50

って、前記ノードの各々が、
時間値を生成する時間カウンタと、
前記ノードと前記ＳＣＩリングとのインタフェースをとるＳＣＩコントローラと、
前記ＳＣＩコントローラと前記時間カウンタとの間で同期信号を搬送する同期信号配信経路とを備えており、
前記複数のノードのうちの１つがマスタノードとして指定され、該マスタノードが同期信号を生成する手段を備えており、該同期信号が、前記同期信号配信経路および前記ＳＣＩリングを介して該コンピュータシステムの残りのノードへと配信されて該コンピュータシステム内の前記時間カウンタの各々の時間値が変更されることを特徴とする、コンピュータシステム。

10

【００４４】

２．前記ノードの各々が、
データを処理する少なくとも１つのプロセッサと、
データを記憶するメモリと、
時間カウンタを有し、ターゲットメモリとのトランザクションを求める前記プロセッサからの要求のディスパッチを行い、及び前記プロセッサに応答を経路指定する、少なくとも１つのプロセッサエージェントと、
該プロセッサエージェントからの要求を受信してターゲットメモリの位置を判定する、前記メモリに対するアクセスを制御する少なくとも１つのメモリエージェントと、
前記要求および前記応答の経路指定を前記プロセッサエージェントと前記メモリエージェントとの間で行う少なくとも１つのクロスバーとを備えており、
前記メモリエージェントが、前記ターゲットメモリが前記メモリであると判定した場合に、前記メモリにアクセスして前記プロセッサに応答し、
前記メモリエージェントが、前記ターゲットメモリがリモートノード上に位置していると判定した場合に、前記要求を前記ＳＣＩコントローラへ転送し、該ＳＣＩコントローラが前記要求を前記ＳＣＩリングを介して前記リモートノードへ送信する、前項１に記載のコンピュータシステム。

20

【００４５】

３．前記ノードの各々が、８つの時間カウンタを有する８つのプロセッサエージェントと、１６個のプロセッサと、８つのＳＣＩコントローラとを備えており、
前記プロセッサエージェントの各々が２つの前記プロセッサに接続されている、前項２に記載のコンピュータシステム。

30

【００４６】

４．前記マスタノード内の前記プロセッサエージェントのうちの１つが、マスタプロセッサエージェントとして指定されており、
前記マスタノード内の前記ＳＣＩコントローラのうちの１つが、マスタＳＣＩコントローラとして指定されており、
同期信号を生成する前記手段が、前記マスタ処理エージェントに存在し、
前記同期信号配信経路が、各ノード上の全ての前記ＳＣＩコントローラと全ての前記時間カウンタとの間で同期信号を搬送し、
前記マスタ処理エージェントが、前記同期信号を前記同期信号配信経路を介して前記マスタノード内の残りの処理エージェントへ配信し、前記同期信号により前記マスタノード内の前記時間カウンタの各々の時間値を変更し、
前記マスタＳＣＩコントローラが、前記同期信号を前記ＳＣＩリングを介して該コンピュータシステムの残りのノードへ配信し、
前記残りのノードの前記ＳＣＩコントローラが、前記同期信号を受信して該同期信号を残りの各ノード内の配信経路を介して該残りのノード内の処理エージェントへ配信し、該同期信号により残りの各ノード内の各時間カウンタの時間値を変更する、前項３に記載のコンピュータシステム。

40

【００４７】

50

５．前記複数のノードが二次元アレイとして構成され、該二次元アレイが、その個々の行における各ノードを接続する少なくとも１つの個々のＳＣＩリングと、前記二次元アレイの個々の列における各ノードを接続する少なくとも１つの個々のＳＣＩリングとを有している、前項４に記載のコンピュータシステム。

【００４８】

６．前記各ノードが８つのＳＣＩコントローラをそれぞれ備えており、前記二次元アレイが、前記各行におけるノードを接続する８つのＳＣＩリングと、前記各列におけるノードを接続する８つのＳＣＩリングとを有するようになっている、前項５に記載のコンピュータシステム。

【００４９】

７．同期信号を生成する前記手段が、
前記同期信号配信経路を介して前記マスタノード上の前記時間カウンタおよび前記ＳＣＩコントローラへ配信される信号を生成するパルス同期信号生成器を更に備えており、前記ＳＣＩコントローラが、前記信号を受信して利用可能なデータパケットを見つけ、データパケット中の時間カウンタ同期ビットをセットし、セットされた該ビットを含むデータパケットが、前記ＳＣＩリングを介して該コンピュータシステムの残りのノードへ配信される、前項１ないし前項６の何れか１つに記載のコンピュータシステム。

【００５０】

８．前記時間カウンタ同期ビットが前記データパケットのヘッダにある、前項１ないし前項６の何れか１つに記載のコンピュータシステム。

【００５１】

９．送信パケット、エコーパケット、またはアイドルパケットからなる群から前記データパケットが選択される、前項８に記載のコンピュータシステム。

【００５２】

１０．前記時間カウンタがタイム・オブ・センチュリー・カウンタである、前項１ないし前項６の何れか１つに記載のコンピュータシステム。

【図面の簡単な説明】

【図１】本発明の同期構成を有するＳＣＩ相互接続を用いたマルチノード・マルチプロセッサシステムの概要を示すブロック図である。

【図２】信号ノードを示す図１のシステムの概要を一層詳細に示すブロック図である。

【図３】 １１２ノードシステムの概要を示す説明図である。

【図４】同期パルス配信構成の概要を示すブロック図である。

【図５】タイム・オブ・センチュリー・クロック・ハードウェアを示すブロック図である。

。

【図６】典型的なＳＣＩデータパケットのレイアウトを示す説明図である。

【図７】ＴＡＣ ＴＯＣ構成レジスタを示す説明図である。

【符号の説明】

- １０ プロセッサ
- １１ プロセッサエージェント
- １２ クロスバー
- １３ タイム・オブ・センチュリー・カウンタ
- １４, ２０ メモリアクセスコントローラ
- １５, １９ トロイダルアクセスチップ
- １６, １８ ＳＣＩリング
- ２１, ２２ ワイヤ
- ２９ メモリバンク

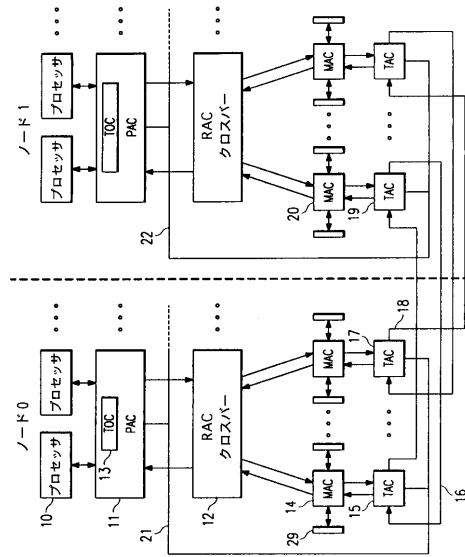
10

20

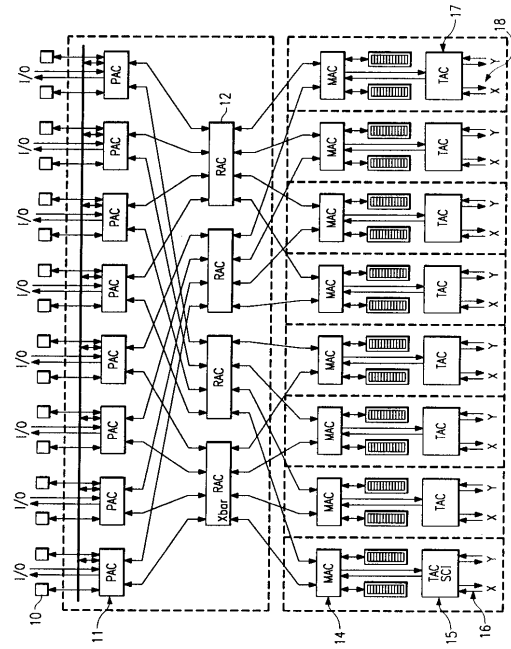
30

40

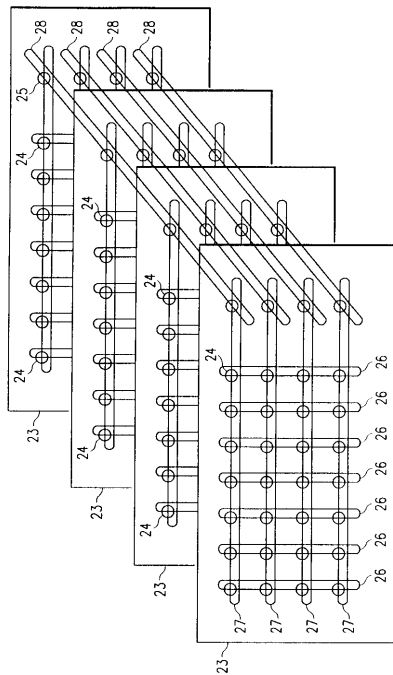
【図 1】



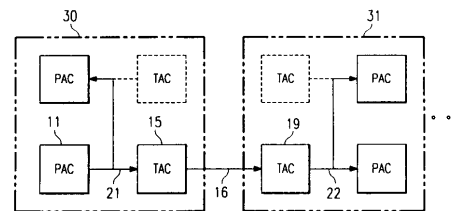
【図 2】



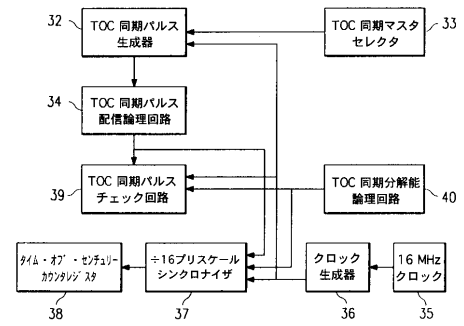
【図 3】



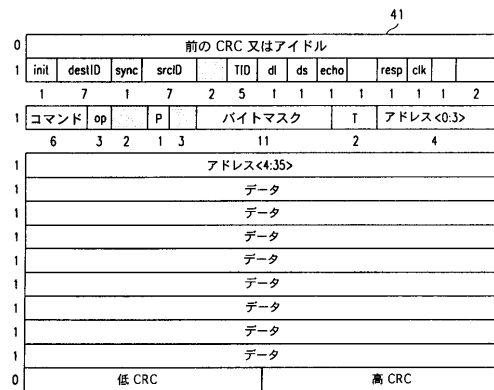
【図 4】



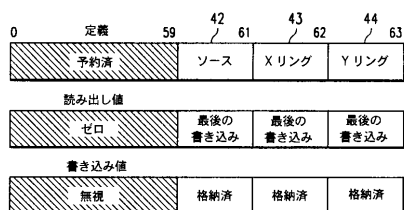
【図 5】



【図 6】



【図 7】



フロントページの続き

(72)発明者 ブライアン・ディー・ホーナンク

アメリカ合衆国テキサス州75075, プラノ, クリークフィールド・ドライブ・1108

(72)発明者 トニー・エム・ブレワー

アメリカ合衆国テキサス州75080, リチャードソン, タム・オシャンター・レーン・3201

審査官 鳥居 稔

(56)参考文献 特開平02-216574(JP, A)

特開平10-097505(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 1/12

G06F 13/00