

(19) 日本国特許庁(JP)

(12) 公表特許公報(A)

(11) 特許出願公表番号

特表2007-526540

(P2007-526540A)

(43) 公表日 平成19年9月13日(2007.9.13)

(51) Int. Cl.	F I	テーマコード (参考)
G06F 15/177 (2006.01)	G06F 15/177 A	5B045
G06F 15/173 (2006.01)	G06F 15/173 640A	

審査請求 未請求 予備審査請求 未請求 (全 20 頁)

(21) 出願番号	特願2006-517396 (P2006-517396)	(71) 出願人	504056174 ニューイシス・インコーポレーテッド NEWISYS INCORPORATED アメリカ合衆国 テキサス州78730, オースティン、ブリッジポイント・パーク ウェイ、6200、ビルディング 4, ス イート 100
(86) (22) 出願日	平成16年6月18日 (2004. 6. 18)	(74) 代理人	110000028 特許業務法人明成国際特許事務所
(85) 翻訳文提出日	平成18年2月23日 (2006. 2. 23)	(72) 発明者	コタ・ラジェッシュ アメリカ合衆国 テキサス州78727 オースティン、アラメダ・トレイス・サー クル, 12320, アpartment #1 107
(86) 国際出願番号	PCT/US2004/019524		
(87) 国際公開番号	W02005/003893		
(87) 国際公開日	平成17年1月13日 (2005. 1. 13)		
(31) 優先権主張番号	10/607, 819		
(32) 優先日	平成15年6月27日 (2003. 6. 27)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 動的マルチクラスタシステムリコンフィギュレーション

(57) 【要約】

【課題】動的なマルチクラスタシステムのコンフィギュレーションの変更を可能にする方法および装置を提供する。

【解決手段】ある例では、マルチクラスタシステムにおけるプロセッサ群は、仮想アドレス空間を共有する。動的にプロセッサ、I/Oリソース、およびクラスタを導入および除去するメカニズムが提供される。このメカニズムは、リセットのあいだ、またはシステムが動作中に実現されえる。リンクは動的にイネーブルまたはディセーブルされえる。

【選択図】 図 8

【特許請求の範囲】

【請求項 1】

コンピュータシステムであって、

第 1 複数プロセッサおよび第 1 相互接続コントローラを含む第 1 クラスタであって、前記第 1 複数プロセッサおよび前記第 1 相互接続コントローラはポイントツーポイントアーキテクチャを用いて通信する、第 1 クラスタ、

第 2 複数プロセッサおよび第 2 相互接続コントローラを含む第 2 クラスタであって、前記第 2 複数プロセッサおよび前記第 2 相互接続コントローラはポイントツーポイントアーキテクチャを用いて通信する、第 2 クラスタ

を備え、

前記第 1 相互接続コントローラから前記第 2 相互接続コントローラへのリンクのためのポリングが、前記第 1 相互接続コントローラをコンフィギュラすることによってイネーブルまたはディセーブルされえる

コンピュータシステム。

10

【請求項 2】

請求項 1 に記載のコンピュータシステムであって、プロセッサの前記第 1 クラスタおよびプロセッサの前記第 2 クラスタは単一の仮想アドレス空間を共有するコンピュータシステム。

【請求項 3】

請求項 1 ~ 2 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、物理レイヤイネーブル指示子を含むコンピュータシステム。

20

【請求項 4】

請求項 1 ~ 3 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、前記第 1 相互接続コントローラおよび前記第 2 相互接続コントローラ間で論理パケットの送信を防ぐようコンフィギュラ可能なフェンス指示子を含むコンピュータシステム。

【請求項 5】

請求項 1 ~ 4 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、前記第 1 相互接続コントローラが前記リンクを再初期化するよう指示するようコンフィギュラ可能な再初期化指示子を含むコンピュータシステム。

30

【請求項 6】

請求項 5 に記載のコンピュータシステムであって、再初期化は、前記第 1 相互接続コントローラに関連付けられる送信機にトレーニングシーケンスを前記第 2 相互接続コントローラへ送信させることを含むコンピュータシステム。

【請求項 7】

請求項 6 に記載のコンピュータシステムであって、前記送信機は、ポリングアクティブ状態が設定されるときに前記トレーニングシーケンスを送るコンピュータシステム。

【請求項 8】

請求項 7 に記載のコンピュータシステムであって、前記送信機は、前記ポリングスリーブ状態が設定されるときにトレーニングシーケンスを送らないコンピュータシステム。

40

【請求項 9】

請求項 5 に記載のコンピュータシステムであって、再初期化は、前記第 1 相互接続コントローラに関連付けられる送信機に初期化シーケンスを前記第 2 相互接続コントローラへ送信させることを含むコンピュータシステム。

【請求項 10】

請求項 1 ~ 9 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、プロセッサのリモートクラスタを特定する値を保持するよう動作可能である複数のクラスタ ID 指示子を含むコンピュータシステム。

【請求項 11】

請求項 1 ~ 10 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続

50

コントローラは、物理レイヤイネーブル、フェンス、再初期化、およびクラスタIDビットを備えるコンフィギュレーションスペースレジスタを含むコンピュータシステム。

【請求項12】

プロセッサのクラスタを導入する方法であって、

第2相互接続コントローラの存在を探してポリングするポイントツーポイントアーキテクチャを用いて通信する第1複数プロセッサを含む第1クラスタ中の第1相互接続コントローラをコンフィギュラすること、

ポイントツーポイントアーキテクチャを用いて通信する第2複数プロセッサを含む第2クラスタ中で第2クラスタ中の第2相互接続コントローラ上のリセット信号をアサートすること、

リンクレイヤプロトコルを前記第1および第2相互接続コントローラ間の接続上で確立すること

を含む方法。

10

【請求項13】

請求項12に記載の方法であって、ポリングは連続的に実行される方法。

【請求項14】

請求項12～13のいずれかに記載の方法であって、前記第1相互接続コントローラは、物理レイヤイネーブル指示子を含む方法。

【請求項15】

請求項12～14のいずれかに記載の方法であって、前記第1相互接続コントローラは、前記第1相互接続コントローラおよび前記第2相互接続コントローラ間で論理パケットの送信を防ぐようコンフィギュ可能なフェンス指示子を含む方法。

20

【請求項16】

請求項12～15のいずれかに記載の方法であって、前記第1相互接続コントローラは、前記第1相互接続コントローラが前記リンクを再初期化するよう指示するようコンフィギュ可能な再初期化指示子を含む方法。

【請求項17】

請求項12～16に記載の方法であって、再初期化は、前記第1相互接続コントローラに関連付けられる送信機にトレーニングシーケンスを前記第2相互接続コントローラへ送信させることを含む方法。

30

【請求項18】

請求項17に記載の方法であって、前記送信機は、ポリングアクティブ状態が設定されるときに前記トレーニングシーケンスを送る方法。

【請求項19】

請求項18に記載の方法であって、前記送信機は、前記ポリングスリープ状態が設定されるときにトレーニングシーケンスを送らない方法。

【請求項20】

請求項16に記載の方法であって、再初期化は、前記第1相互接続コントローラに関連付けられる送信機に初期化シーケンスを前記第2相互接続コントローラへ送信させることを含む方法。

40

【請求項21】

請求項12～20のいずれかに記載の方法であって、前記第1相互接続コントローラは、プロセッサのリモートクラスタを特定する値を保持するよう動作可能である複数のクラスタID指示子を含む方法。

【請求項22】

請求項12～21のいずれかに記載の方法であって、前記第1相互接続コントローラは、物理レイヤイネーブル、フェンス、再初期化、およびクラスタIDビットを備えるコンフィギュレーションスペースレジスタを含む方法。

【請求項23】

コンピュータシステムであって、

50

第 2 相互接続コントローラの存在を探してポリングするポイントツーポイントアーキテクチャを用いて通信する第 1 複数プロセッサを含む第 1 クラスタ中の第 1 相互接続コントローラをコンフィギュアする手段、

ポイントツーポイントアーキテクチャを用いて通信する第 2 複数プロセッサを含む第 2 クラスタ中で第 2 クラスタ中の第 2 相互接続コントローラ上のリセット信号をアサートする手段、

リンクレイヤプロトコルを前記第 1 および第 2 相互接続コントローラ間の接続上で確立する手段

を備えるコンピュータシステム。

【請求項 2 4】

請求項 2 3 に記載のコンピュータシステムであって、ポリングは連続的に実行されるコンピュータシステム。

【請求項 2 5】

請求項 2 3 ~ 2 4 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、物理レイヤイーネブル指示子を含むコンピュータシステム。

【請求項 2 6】

請求項 2 3 ~ 2 5 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、前記第 1 相互接続コントローラおよび前記第 2 相互接続コントローラ間で論理パケットの送信を防ぐようコンフィギュア可能なフェンス指示子を含むコンピュータシステム。

【請求項 2 7】

請求項 2 3 ~ 2 6 のいずれかに記載のコンピュータシステムであって、前記第 1 相互接続コントローラは、前記第 1 相互接続コントローラが前記リンクを再初期化するように指示するようコンフィギュア可能な再初期化指示子を含むコンピュータシステム。

【請求項 2 8】

請求項 2 3 ~ 2 7 に記載のコンピュータシステムであって、再初期化は、前記第 1 相互接続コントローラに関連付けられる送信機にトレーニングシーケンスを前記第 2 相互接続コントローラへ送信させることを含むコンピュータシステム。

【請求項 2 9】

請求項 2 8 に記載のコンピュータシステムであって、前記送信機は、ポリングアクティブ状態が設定されるときに前記トレーニングシーケンスを送るコンピュータシステム。

【請求項 3 0】

請求項 2 9 に記載のコンピュータシステムであって、前記送信機は、前記ポリングスリープ状態が設定されるときにトレーニングシーケンスを送らないコンピュータシステム。

【請求項 3 1】

請求項 2 7 に記載のコンピュータシステムであって、再初期化は、前記第 1 相互接続コントローラに関連付けられる送信機に初期化シーケンスを前記第 2 相互接続コントローラへ送信させることを含むコンピュータシステム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

本発明は、一般にマルチプロセッサシステムに関する。より具体的には、本発明は、マルチクラスタシステム中のリソースを効果的かつ効率的に管理する技術に関する。

【背景技術】

【0 0 0 2】

パフォーマンスの制限のために、単一のメモリ空間でシステム中のプロセッサ群を接続するポイントツーポイントアーキテクチャが開発されるようになってきている。ある例では、個別プロセッサは、互いに複数のポイントツーポイントリンクを通して直接に接続されることによってプロセッサ群のクラスタを形成しえる。プロセッサ群の別個のクラスタも接続されえる。ポイントツーポイントリンクは、コプロセッシングおよびマルチプロセッシ

10

20

30

40

50

ング機能のための帯域幅を大きく増加させる。

【0003】

マルチプロセッサクラスタを接続するのにさまざまなメカニズムが用いられる。しかし、動的にマルチクラスタシステムのコンフィギュレーションを変化させるメカニズムは限られてきた。さまざまな困難のために、マルチクラスタシステムの物理レイヤプロトコルを含む通信プロトコルのその機能性が制限されてきた。

【発明の開示】

【発明が解決しようとする課題】

【0004】

したがって、ポイントツーポイントリンクを用いて接続されるマルチプロセッサのマルチクラスタを有するシステム中のクラスタ群の間の通信のためのメカニズムおよび技術を改良する技術を提供することが望ましい。

【課題を解決するための手段】

【0005】

本発明によれば、動的なマルチクラスタシステムのコンフィギュレーションの変更を可能にする方法および装置が提供される。ある例では、マルチクラスタシステムにおけるプロセッサ群は、仮想アドレス空間を共有する。動的にプロセッサ、I/Oリソース、およびクラスタを導入および除去するメカニズムが提供される。このメカニズムは、リセットのあいだ、またはシステムが動作中に実現されえる。リンクは動的にイネーブルまたはディセーブルされえる。

【0006】

ある実施形態において、コンピュータシステムが提供される。このコンピュータシステムは、第1クラスタおよび第2クラスタを含む。第1クラスタは、第1複数プロセッサおよび第1相互接続コントローラを含む。前記第1複数プロセッサおよび前記第1相互接続コントローラはポイントツーポイントアーキテクチャを用いて通信する。第2クラスタは、第2複数プロセッサおよび第2相互接続コントローラを含む。前記第2複数プロセッサおよび前記第2相互接続コントローラはポイントツーポイントアーキテクチャを用いて通信する。前記第1相互接続コントローラから前記第2相互接続コントローラへのリンクのためのポリングが、前記第1相互接続コントローラをコンフィギュラすることによってイネーブルまたはディセーブルされえる。

【0007】

他の例では、プロセッサのクラスタを導入する方法が提供される。第2相互接続コントローラの存在を探してポリングするよう、ポイントツーポイントアーキテクチャを用いて通信する第1複数プロセッサを含む第1クラスタ中の第1相互接続コントローラがコンフィギュラされる。ポイントツーポイントアーキテクチャを用いて通信する第2複数プロセッサを含む第2クラスタ中で第2クラスタ中の第2相互接続コントローラ上でリセット信号がアサートされる。リンクレイヤプロトコルが前記第1および第2相互接続コントローラ間の接続上で確立される。

【0008】

本発明の性質および優位性のさらなる理解は、明細書および図面の残りの部分を参照することによって実現されよう。

【発明を実施するための最良の形態】

【0009】

本発明は、以下の説明を添付の図面と併せて参照することによって最もよく理解されえ、図面は本発明の具体的な実施形態を図示するものである。

【0010】

本発明を実施するための発明者によって考えられたベストモードを含む本発明のいくつかの具体的な実施形態が詳細にここで参照される。これら具体的な実施形態の例は、添付の図面に示される。本発明は、これら具体的な実施形態について説明されるが、本発明を記載された実施形態に限定するよう意図されてはいないことが理解されよう。むしろ、代

10

20

30

40

50

替物、改変、および等価物を、添付の特許請求の範囲によって規定される本発明の精神および範囲内に含まれるとしてカバーするよう意図される。そのプロセッサ間でポイントツーポイント通信を有するマルチプロセッサアーキテクチャは、本発明の具体的な実施形態を実現するのに適する。以下の記載において、本発明の完全な理解を促すために多くの具体的な詳細が述べられている。本発明は、これら具体的な詳細の一部または全てなしでも実施されえる。本発明の趣旨を不必要にぼかさないために、よく知られたプロセス操作は詳細には記載されていない。さらに本発明が単一のエンティティを参照する場合、本発明の方法および装置は、文脈上明らかにそうでないと表されない限り、1より多いエンティティを用いて実現可能である。

【0011】

図1Aは、本発明の技術を採用しえるマルチクラスタ、マルチプロセッサシステムの一例の概略図である。それぞれの処理クラスタ101、103、105、および107は、複数のプロセッサを含む。処理クラスタ101、103、105、および107は、互いにポイントツーポイントリンク111a~fを通して接続される。ある実施形態において、図1Aに示されるマルチクラスタアーキテクチャ中のマルチプロセッサは、同じメモリ空間を共有する。この例では、ポイントツーポイントリンク111a~fは、マルチクラスタ101、103、105、および107内のマルチプロセッサを接続するために従来のフロントサイドバスに代わって用いられる内部システム接続である。このポイントツーポイントリンクは、任意のポイントツーポイントのコヒーレンスプロトコルをサポートしえる。

10

20

【0012】

図1Bは、本発明の技術を採用しえるマルチクラスタ、マルチプロセッサシステムの他の例の概略図である。それぞれの処理クラスタ121、123、125、および127は、ポイントツーポイントリンク141a~dを通してスイッチ131に結合される。スイッチおよびポイントツーポイントリンクを用いることは、システム内でマルチクラスタを接続するとき、より少ないポイントツーポイントリンクで実現することを可能にすることに注意されたい。スイッチ131は、コヒーレンスプロトコルインタフェースを持つプロセッサを含みえる。さまざまな実現例によれば、図1Aに示されるマルチクラスタシステムは、図1Bに示されるスイッチ131を用いて拡張されえる。

30

【0013】

図2は、例えば図1Aに示されるクラスタ101のようなマルチプロセッサクラスタの概略図である。クラスタ200は、プロセッサ202a~202d、1つ以上の基本I/Oシステム(BIOS)204、メモリバンク206a~206dを備えるメモリサブシステム、ポイントツーポイント通信リンク208a~208e、およびサービスプロセッサ212を含む。ポイントツーポイント通信リンクは、プロセッサ202a~202d、I/Oスイッチ210、および相互接続コントローラ230の間で相互接続を可能にするよう構成される。サービスプロセッサ212は、リンク214a~214fによって図2に表されるJTAGインタフェースを介してプロセッサ202a~202d、I/Oスイッチ210、および相互接続コントローラ230と通信を可能にするよう構成される。他のインタフェースがサポートされることに注意されたい。ある実現例においては、サービスプロセッサは、マルチプロセッサクラスタに含まれないことにも注意されたい。I/Oスイッチ210は、システムの残りをI/Oアダプタ216および220に接続する。

40

【0014】

具体的な実施形態によれば、本発明のサービスプロセッサは、システムリソースを以前に特定されたパーティションスキーマに従って分割する(partition)ためのインテリジェンスを有する。パーティショニングは、システムプロセッサに関連付けられたルーティングテーブルをサービスプロセッサによって直接に操作することによって達成されえ、これはポイントツーポイント通信インフラストラクチャによって可能になる。ルーティングテーブルは、さまざまなシステムリソースを制御および分離するために用いられ、リソース間の接続がテーブル中で定義される。

50

【0015】

プロセッサ202a~dは、相互接続コントローラ230にもポイントツーポイントリンク232a~dを通して結合される。さまざまな実施形態によれば、また後で詳述されるように、相互接続コントローラ230は、システム中の相互接続されたプロセッサの数が、ノードIDスペースおよび複数のプロセッサクラスタのそれぞれに関連付けられたマッピングテーブル制限を越えることができるようにするさまざまな機能を実行する。ある実施形態によれば、相互接続コントローラ230は、クラスタにわたるキャッシュコヒーレンスの維持を含むさまざまな他の機能を実行する。相互接続コントローラ230は、他のマルチプロセッサクラスタ群に関連付けられた同様のコントローラ群に結合されえる。1つのクラスタにこのような相互接続コントローラが1つより多く存在しえることに注意されたい。相互接続コントローラ230は、リモートクラスタと共にプロセッサ202a~dともポイントツーポイントプロトコルを用いて通信する。

10

【0016】

より一般には、図2に示される具体的なアーキテクチャは、単に例示的であって、本発明の実施形態は、異なる構成およびリソース相互結合、および示されたシステムリソースのそれぞれについてさまざまな代替物を有すると想定されることが理解されよう。しかし例示する目的のためにサーバ200の特定の詳細が想定される。例えば、図2に示されるリソースのほとんどは、単一の電子アセンブリ上に常駐すると想定される。加えて、メモリバンク206a~206dは、デュアルインラインメモリモジュール(DIMM)として物理的に提供されるダブルデータレート(DDR)メモリを備えうる。I/Oアダプタ216は、例えば、永久記憶デバイスにアクセスを提供するウルトラダイレクトメモリアクセス(UDMA)コントローラ、またはスモールコンピュータシステムインタフェース(SCSI)コントローラでありえる。I/Oアダプタ220は、例えばローカルエリアネットワーク(LAN)またはインターネットのようなネットワークと通信を提供するよう構成されたイーサネットカード(イーサネットは登録商標である)でありえる。BIOS204は、フラッシュメモリのような任意の持続性メモリでありえる。

20

【0017】

ある実施形態によれば、サービスプロセッサ212は、集積化されたチップセット機能を含むマイクロプロセッサMPC855Tマイクロプロセッサである。相互接続コントローラ230は、ローカルポイントツーポイントコヒーレンスプロトコルをサポートする特定用途向け集積回路(ASIC)である。相互接続コントローラ230は、I/Oデバイスとの通信を可能にするために非コヒーレントプロトコルを扱うようにも構成されえる。ある実施形態において、相互接続コントローラ230は、プログラマブルロジックデバイスまたはフィールドプログラマブルゲートアレイのような特別に構成されたプログラマブルチップである。他の実施形態において、相互接続コントローラ230は、特定用途向け集積回路(ASIC)である。さらに他の実施形態において、相互接続コントローラ230は、相互接続パケットトラフィックにアクセスし処理する機能によって拡張された汎用プロセッサである。

30

【0018】

図3は、相互接続コントローラ230の一例の概略図である。さまざまな実施形態によれば、相互接続コントローラは、マルチプロセッサシステムのさまざまなクラスタ中のプロセッサ群から受け取られたプローブおよびリクエストのようなパケットを扱うよう構成されたプロトコルエンジン305を含む。プロトコルエンジン305の機能は、パフォーマンスを改善するためにいくつかのエンジンにわたって分割されえる。ある例では、パーティショニングは、パケットタイプ(リクエスト、プローブおよびレスポンス)、方向(着信および発信)、またはトランザクションフロー(リクエストフロー、プローブフローなど)に基づいてなされえる。

40

【0019】

プロトコルエンジン305は、最近のリクエストおよびプローブのようなトランザクションを相互接続コントローラが追跡し、これらトランザクションを特定のプロセッサ群と

50

関連付けることを可能にするペンディングバッファ309へのアクセスを有する。ペンディングバッファ309で維持されるトランザクション情報は、トランザクション宛先ノード、後で起こるコリジョン検出およびプロトコル最適化のためのリクエストのアドレス、レスポンス情報、タグ、および状態情報を含みえる。ヒストリバッファ311は、他のクラスタ中の相互接続コントローラによって成功して受け取られなかったパケットの効率的な再送信を可能にするためにも提供されえる。

【0020】

相互接続コントローラは、外部プロセッサクラスタと共に、クラスタ中の他のプロセッサと相互接続コントローラが通信することを可能にするコヒーレントプロトコルインタフェース307を有する。相互接続コントローラは、I/Oデバイス（例えば図2においてリンク208cおよび208dによって表されるような）と通信するための非コヒーレントプロトコルインタフェース311のような他のインタフェースも含みえる。さまざまな実施形態によれば、それぞれのインタフェース307および311は、フルクロスバーとして、またはマルチプレクサおよびバッファのような要素を用いた別個の受信および送信ユニットとしてのいずれかで実現される。相互接続コントローラ230は、コヒーレントおよび非コヒーレントインタフェースの両方を提供する必要は必ずしもないことに注意されたい。またあるクラスタ内の相互接続コントローラ230は、他のクラスタ内の相互接続コントローラ230と通信しえることにも注意されたい。

10

【0021】

本発明のさまざまな実施形態によれば、プロセッサ202a~202dは、実質的に同一である。図4は、複数のポート404a~404cを有するインタフェース402およびそれに関連付けられたルーティングテーブル406a~406cを含むそのようなプロセッサ202の簡略化されたブロック図である。それぞれのポート404は、例えば図2のリンク208a~208eのような関連するリンクを介してコンピュータシステム中の他のリソース、例えばプロセッサまたはI/Oデバイスと通信できる。

20

【0022】

図4に示されるインタフェースは、さまざまなトポロジー、例えばリング、メッシュなどの任意のものによるシステムプロセッサを相互接続する複数のセグメントを備えるポイントツーポイント、分散化ルーティングメカニズムとして一般化されえる。これらセグメントのそれぞれのエンドポイントのそれぞれは、ユニークなノードIDおよびそれが「所有する」複数の関連付けられたリソース、例えばメモリおよびそれが接続されるI/Oを有する、接続されたプロセッサに関連付けられる。

30

【0023】

分散ルーティングメカニズム中のノードのそれぞれに関連付けられたルーティングテーブルは、コンピュータシステムリソース内の相互接続の現在の状態を全体として表現する。任意の与えられたノード（例えばプロセッサ）によって所有されるリソース（例えば特定のメモリ範囲またはI/Oデバイス）のそれぞれは、ルーティングテーブル（群）中で、アドレスとしてそのノードと関連付けられる。リクエストがノードに到着すると、リクエストされたアドレスが、適切なノードおよびリンク、すなわちアドレスのある範囲内の与えられた特定のアドレスを特定するノードのルーティングテーブル中の2つのレベルのエントリと比較され、これらはノードxに行け（go to node x）およびノードxについてリンクyを使え（for node x use link y）である。

40

【0024】

図4に示されるように、プロセッサ202は、関連付けられたルーティングテーブル中の情報によって3つの他のプロセッサ群とポイントツーポイント通信を行いえる。具体的な実施形態によれば、ルーティングテーブル406a~406cは、2レベルテーブルを備え、第1レベルは、システムリソース（例えばメモリバンク）のユニークなアドレスを対応するノード（例えばプロセッサのうちの一つ）に関連付け、第2レベルは、現在のノードからそのノードへ行くのに用いられるように、それぞれのノードをリンク（例えば208a~208e）に関連付ける。

50

【0025】

プロセッサ202は、JTAGハンドシェイクレジスタ408のセットも有し、これはとりわけサービスプロセッサ（例えば図2のサービスプロセッサ212）およびプロセッサ202の間の通信を促進する。すなわち、サービスプロセッサは、ルーティングテーブル406a~406cでの最終的な記憶のために、ルーティングテーブルエントリをハンドシェイクレジスタ408に書き込める。図4に図示されたプロセッサアーキテクチャは、本発明の具体的な実施形態を記載する目的のための単なる例示であることが理解されよう。例えば、本発明の他の実施形態を実現するために、より少ないまたはより多い数のポートおよび/またはルーティングテーブルが用いられる。

【0026】

上述のように、本発明の具体的な実施形態が基づく基本プロトコルは、特定の實現例によれば、3ビット空間である制限されたノードIDスペースを提供し、したがって8ノードのみユニークな特定が可能になる。すなわち、もしこの基本プロトコルが、本発明によって示される技術革新なしで採用されるなら、単一のクラスタ中ではポイントツーポイントインフラストラクチャを介して8ノードしか相互接続できない。この制限を回避するために、本発明は、特定のクラスタ内で単一レイヤ識別スキームを保持しつつ、一方で他の同様に配置されるクラスタおよび処理ノードとの相互接続およびそれら間の通信を可能にする階層メカニズムを導入する。

【0027】

具体的な実施形態によれば、それぞれのマルチプロセッサクラスタ内のノード群のうちの一つは、相互接続コントローラ、例えば図2の相互接続コントローラ230であり、これが情報の階層マッピングを管理し、それによってマルチクラスタが単一のアドレス空間を共有することを可能にし、一方で同時にそのクラスタ内のプロセッサが、それら自身のクラスタの外のものの「知識」なしに、任意のクラスタ内の任意のプロセッサと動作し、かつ対話することを可能にする。相互接続コントローラは、その関連付けられたプロセッサには、クラスタ内のプロセッサまたはノードのうちの単なる別の一つであるように見える。

【0028】

基本プロトコルにおいて、クラスタ内の特定のプロセッサがリクエストを生成するとき、アドレスマッピングテーブルのセットは、そのリクエストをクラスタ内の他のノードのうちの一つにマッピングするために採用される。すなわち、クラスタ内のそれぞれのノードは、それが関連付けられる共有されたメモリ空間の一部を有する。異なるタイプのアドレスマッピングテーブルがメインメモリ、メモリマッピングされたI/O、異なるタイプのI/O空間などのために存在する。これらアドレスマッピングテーブルは、クラスタ中の特定のノードへのリクエスト中で特定されたアドレスをマッピングする。

【0029】

それからルーティングテーブルのセットは、リクエストするノードからアドレスマッピングテーブルから特定されたノードへどのようにたどり着くかを決定するために採用される。すなわち、上述のように、それぞれのプロセッサ（すなわちクラスタノード）は、リクエストを現在のノードからアドレスマッピングテーブルから特定されたノードへ送信するために用いられる、ポイントツーポイントインフラストラクチャ内の特定のリンクを特定するルーティングテーブルを関連付ける。一般に、ノードは、1つまたは複数のリソース（例えばプロセッサを含む）に対応しえるが、ノードおよびプロセッサという語はここではしばしば交換できるように用いられる。特定の實現例によれば、ノードは、複数のサブユニット、例えばCPU、メモリコントローラ、I/Oブリッジなどを備え、これらのそれぞれはユニットIDを有する。

【0030】

加えて、個別のトランザクションは、非連続のパケット（non-consecutive packets）中にセグメント化されるので、それぞれのパケットは、パケットがトランザクションを開始するノードに対して関連付けられるトランザクションを特定するために、ユニークな

10

20

30

40

50

トランザクションタグを含む。具体的な実現例によれば、トランザクションタグは、ソースノード（3ビットフィールド）、ソースノードユニット（2ビットフィールド）、およびトランザクションID（5ビットフィールド）を特定する。

【0031】

よって、トランザクションが特定のノードにおいて開始されるとき、アドレスマッピングテーブルは、それからパケットにアペンドされ、ルーティングテーブルによってそのパケットをルーティングする適切なリンク（群）を特定するよう用いられる宛先ノード（およびユニット）を特定するために採用される。ソース情報は、リクエストに適切に応答するためにリクエストでプロンプされる宛先ノードおよび任意の他のノードによって用いられる。

10

【0032】

具体的な実施形態によれば、かつ上述のように、それぞれのクラスタ中の相互接続コントローラは、そのクラスタ中の他のプロセッサには、クラスタ中の他のプロセッサであるように単に見える。しかし、相互接続コントローラに関連付けられる共有されたメモリ空間の一部は、実際にはグローバルに共有されたメモリ空間の残り、すなわちシステム中の他のクラスタに関連付けられるメモリを含む。すなわち、特定のクラスタ中のローカルプロセッサの視点からは、システム中の他のマルチプロセッサクラスタの全てに関連付けられるメモリ空間は、それら自身のクラスタ中の相互接続コントローラ（群）によって表される。

【0033】

図5を参照して説明されるさらにより具体的な実施形態によれば、それぞれのクラスタは、4つのプロセッサ202a~dおよび相互接続コントローラ230を含む5つのノード（例えば図2示されるように）を有し、それらのそれぞれは、クラスタ内でユニークな3ビットノードIDによって表される。上述のように、それぞれのプロセッサ（すなわちクラスタノード）は、例えばCPU、メモリコントローラなどを含む多くのサブユニットを表しえる。

20

【0034】

本発明によって設計され、そのようなクラスタコンフィギュレーションを仮定する例示的地址マッピングスキームの図示は、図5に示される。図示された例において、グローバルメモリ空間は、ここではクワッドとも呼ばれる（4つのローカルプロセッサをそれぞれ含むので）4つのそのようなクラスタによって共有されるとも仮定される。理解されるように、それぞれのクラスタ内のクラスタおよびノードの個数は、異なる実施形態によって変わりえる。

30

【0035】

アドレスマッピング機能を単一のクラスタを越えて拡張するために、それぞれのクラスタは、そのローカルメモリ空間、すなわちそのクラスタ中のプロセッサに関連付けられるグローバルメモリ空間の一部を、連続した領域内にマッピングし、一方、この領域の上および下のグローバルメモリ空間の残りの部分は、ローカル相互接続コントローラ（群）にマッピングされる。それぞれのクラスタ中の相互接続コントローラは、2つのマッピングテーブルを維持する。すなわち、グローバルマップおよびローカルマップである。ローカルマップは、リモートクラスタからローカルクラスタ内の特定のノードへの着信リクエストをマッピングする。

40

【0036】

図5をここで参照して、それぞれのローカルクラスタは、ローカルメモリマップ（501~504）を有し、これらはローカルメモリ空間（すなわちローカルプロセッサに関連付けられるグローバルメモリ空間の連続部分）をそれぞれのノードにマッピングし、全てのリモートメモリ空間（すなわちグローバルメモリ空間の残り）を、ローカル相互接続コントローラ（群）、例えばクワッド3のノード4に関連付けられる1つまたは2つのマップエントリにマッピングする。ローカルクラスタのそれぞれのノードは、ローカルマップのコピーを有する。それぞれのクラスタ中の相互接続コントローラは、システム中の他の

50

クラスタのそれぞれと共にこれらリモートメモリ空間に関するグローバルマップ(505~508)も保持する。それぞれの相互接続コントローラは、リモートクラスタからそのクラスタ内の個々のノードへ受け取られるリクエストをマッピングするためにローカルマップ(509~511)のそのコピーを用いる。

【0037】

図5を参照して説明される例としてのトランザクションは、例示的でありえる。この例では、クワッド3のノード2は、ローカル相互接続コントローラ(すなわちノード4)へマッピングする(マップ501を介して)リクエストを生成する。相互接続コントローラがこのリクエストを受け取るとき、そのグローバルマップ505は、アドレスをクワッド2にマッピングする。それから相互接続コントローラは、そのリクエストをクワッド2に転送する。クワッド2における相互接続コントローラは、そのローカルメモリマップを用いて、リクエストに対するターゲットへの適切なノード、この例ではノード1、を決定する。

10

【0038】

具体的な実現例において、それぞれのプロセッサまたはクラスタノードは、8つのメモリマップレジスタに限定される。図5を参照して説明されるスキームは、ローカルメモリ空間について4つのエントリを、リモート空間について多くとも2つのレジスタを必要とする。したがって、具体的な実施形態によれば、2つの残りのエントリが領域を再分割するために用いられる。8つのマッピングレジスタの制限は、クワッドにローカルな全てのメモリが連続なブロック内で割り当てられることを要求する。そのような実施形態における相互接続コントローラのローカルメモリも8つのエントリである。しかし相互接続コントローラのグローバルマップのサイズは、システム中のクラスタの個数によって決定される。さまざまな実施形態によれば、メモリマッピングされたI/O空間は、マッピングレジスタの同一のセットによってマッピングされる。

20

【0039】

上述のように、ローカルクラスタレベルで、アドレスマッピングテーブルからの情報は、クラスタ内の宛先ノードへ情報を送信する適切なリンクを特定するために用いられる。上述のグローバルマッピングを用いてクラスタ間の送信を行うためには、同様のメカニズムが必要とされる。したがって、さまざまな実施形態によれば、クラスタ中のそれぞれのノードに関連付けられるローカルルーティングテーブルに加えて、相互接続コントローラは、システム中の他のクラスタを、クラスタを相互接続するさまざまなポイントツーポイント送信リンク(例えば図1Aのリンク111)にマッピングするグローバルルーティング情報を保持する。

30

【0040】

本発明の具体的な実施形態によれば、2つのタイプのローカルルーティングテーブルが採用される。一つはディレクテッドパケットのためであり、一つはブロードキャストパケットのためである。それぞれのテーブル(例えば図4のテーブル406)は、ターゲットノードおよびリンク間のマッピングを保持する。ディレクテッドパケットについて、別個のテーブルがリクエストについておよび応答について用いられる。これは、応答がリクエストにリクエストと同じパスを通過して元に戻ることを可能にする。同じルートを維持することは、デバッグを容易にし、正確さのためには要求されない。ブロードキャストパケットについては、対応するテーブルは、どのリンクにブロードキャストパケットが転送されるかを示す。ブロードキャストパケットは、よって複数のリンクへとルーティングされる。

40

【0041】

本発明の相互接続コントローラの特定の實現例において、そのローカルテーブルは、ローカル宛先ノードをディレクテッドパケットについては4つのリンクのうちの1つにマッピングし、ブロードキャストパケットについては任意の個数のリンクにマッピングする。相互接続コントローラはまた、リモート宛先クラスタを特定のリモートリンクにマッピングするグローバルルーティングテーブルを維持する。特定の實施形態によれば、相互接続

50

コントローラは、グローバルルーティングレベルにおいてパケットのマルチキャストもサポートする。

【0042】

本発明によって設計されたルーティングメカニズムの具体的な実施形態が図6Aおよび6Bを参照してこれから説明される。図6Aのシステム600は、ノード N_0 および N_1 を含む複数のローカルノードをそれぞれ有する4つのクラスタを含む。図6Bのテーブルは、例示的目的のためにシステムの全てのローカルおよびグローバルルーティングテーブルを結合する。

【0043】

例示的トランザクションの一部として、クラスタ0中のノード N_0 におけるCPU602は、クラスタ3中のノード N_0 におけるCPU604に向けられたパケットを生成する。このパケットは、例えば、そのノードにおけるメモリコントローラにマッピングするメモリリクエストでありえる。CPU602が、そのクラスタの外の何ものについても知識を有しないので、それは宛先としてクラスタ0中のノード N_1 （すなわちローカル相互接続コントローラ606）にターゲットするパケットを生成する。上述のように、これは、ノード N_0 によって所有されるローカルメモリマップ（図6Bのテーブルの関連部分を参照）が、ノード N_1 をリモートクラスタによって所有された全てのメモリに対応するものとして特定するという事実による。相互接続コントローラ606は、パケットを受け取り、そのグローバルアドレスマップを用いて（例えば上述のもの）、パケットの最終宛先がクラスタ3であることを決定し、クラスタ3をターゲットするリモートパケットを生成する。それから、そのグローバルルーティングテーブル（図6Bのテーブルの関連部分を参照）を用いて、相互接続コントローラ606は、このパケットがリンク L_1 上に送られなければならないことを決定する。上述のローカルルーティングメカニズムと同様に、ソースおよび宛先クラスタを特定する情報がパケットにアペンドされる。

【0044】

クラスタ1における相互接続コントローラ608がパケットを受け取ると、それはパケットがクラスタ3に宛先が設定されていることも決定し、そのグローバルルーティングテーブル（図6B）からリンク L_2 がそのパケットを送るために用いられなければならないことを決定する。クラスタ3における相互接続コントローラ610は、パケットを受け取り、パケットがローカルクラスタをターゲットとしていることを決定して、そのローカルルーティングテーブル（図6B）を用いて、そのパケットをその宛先に送るためにはローカルリンク L_0 が用いられなければならないことを決定する。それからノード N_0 におけるCPU604は、パケットをリンク L_0 を介して受け取る。ノードID空間が3ビットID空間である具体的な実施形態によれば、このマルチレベルルーティングメカニズムは、クラスタの数に具体的な制限がない8つのローカルノードに拡張されえる。

【0045】

マルチレベルルーティングメカニズムを有することは、システムが比較的束縛されない個数のプロセッサおよび処理クラスタを有することを可能にする。しかし、システムサイズおよび複雑さが増すと共に、システムパフォーマンスを低下させることなくプロセッサおよび処理クラスタを動的に追加および削除することが急激に重要になる。より多くの個数のプロセッサおよび処理クラスタを持つさらなるコンフィギュレーションも可能であり、シームレスにコンフィギュレーション変更を可能にすることが有利になる。コンフィギュレーション変更の一つのタイプは、プロセッサ、処理クラスタ、リンクなどのような欠陥のある要素の除去または動的置換である。

【0046】

従来のメカニズムによってリンクを開いたり閉じたりすることは制限がある。既存のメカニズムは、例えば、リセット中に見つからないリンクを動的に開いたり、リセット中に見つかったリンクを閉じたりすることを提供できない。したがって、ユーザが動的にリンクをイネーブルまたはディセーブルすることを可能にし、マルチプロセッサ、マルチクラスタシステムのコンフィギュレーションを変更する改良点を提供することが望ましい。あ

10

20

30

40

50

る例では、ユーザは、既存のシステムにプロセッサおよび処理クラスタのようなリソースを動的に追加できる。他の例では、ユーザは、クラスタ間の相互接続コントローラまたはハイパートランスポートケーブルのような欠陥のある要素を置換できる。

【 0 0 4 7 】

本発明のさまざまな実施形態によれば、相互接続コントローラがリンクの動的なイネープリングおよびディセープリングを可能にするメカニズムが提供される。ある実施形態において、コンフィギュレーションスペースレジスタは、相互接続コントローラと関連付けられる。コンフィギュレーションスペースレジスタは、物理レイヤおよびリンクレイヤ通信をイネープリングする情報、リンクを再初期化する指示子、および相互接続コントローラに結合された処理クラスタをトラッキングする識別子を含む。

10

【 0 0 4 8 】

図 7 は、コンフィギュレーションスペースレジスタを示す概略図である。コンフィギュレーションスペースレジスタ 7 0 1 は、プロセッサおよびクラスタ間の物理およびリンクレイヤ通信をイネーブルおよびディセーブルするために提供される。コンフィギュレーションスペースレジスタ 7 0 1 は、相互クラスタ通信を管理する相互接続コントローラ中に含まれる。コンフィギュレーションスペースレジスタ 7 0 1 は、他のリソース中に含まれることに注意されたい。物理レイヤおよびリンクレイヤトグル指示子と共に複数のプロセッサクラスタ ID を維持するメカニズムは、ここではコンフィギュレーションスペースレジスタと呼ばれる。ある例では、コンフィギュレーションスペースレジスタ 7 0 1 は、物理レイヤ通信をイネーブルおよびディセーブルするための物理レイヤ指示子 7 1 1 を含む。コンフィギュレーションスペースレジスタへワイヤリングすることによって物理レイヤ通信をディセーブルすることは、プロセッサおよびクラスタ間のリンクをシャットダウンする。ある例では、物理マクロがオフにされ、ディセーブルされた物理マクロから発する電気的アクティビティは存在しなくなる。ポリングは行われず、関連付けられる相互接続コントローラは、ポリングに応答せず、リンクを維持しない。さまざまな実施形態によれば、コンフィギュレーションスペースレジスタ 7 0 1 は、リンクレイヤ指示子 7 1 3 も含む。ある例では、リンクレイヤ指示子 7 1 3 は、フェンスビットである。コンフィギュレーションスペースレジスタへのワイヤリングによってリンクレイヤをディセーブルすることは、リンクが物理レイヤ通信以外の任意の通信のために用いられることを防ぐ。

20

【 0 0 4 9 】

しかし、相互接続コントローラは、ポリングに依然として応答し、既存の接続を維持する。物理レイヤ通信以外のものは何も許されないので、ハイパートランスポートコマンドまたはデータパケットは、ディセーブルされたリンクレイヤを持つリンクにわたっては送られない。ポリング状態は、コンフィギュレーションスペースレジスタを用いて維持される。ある例では、システムがリセットから出てくるときにもしリンクが見つからないなら、物理レイヤプロトコルは、何かのデバイスが取り付けられていないかを決定するためにリンクを連続的にポリングする能力を有する。このポリングの特徴は、コンフィギュレーションスペースレジスタ中の適切なトグルを設定することによってイネーブルまたはディセーブルされる。ポリングすることは、自動認識および通知によってリソースの既存のクラスタへのホットプラグングを可能にする。ある例では、ポリングは、新しいマルチ

30

40

【 0 0 5 0 】

システムコンフィギュレーションを維持するプロセッサは、新しいクラスタが取り付けられるときにリンクをイネーブルし、システムを拡張するか、または単にユーザに最近なされた接続がシステムによって認識されたことを通知するよう選択しえる。さまざまな実施形態によれば、コンフィギュレーションスペースレジスタは、新しい初期化シーケンスが実行されるべきかを示すための再初期化指示子 7 1 5 も含む。コンフィギュレーションスペースレジスタ 7 0 1 が接続されたクラスタの ID を維持することを可能にするために、クラスタ識別子 7 1 7、7 1 9、および 7 2 1 を保持するレジスタも提供される。コンフィギュレーションスペースレジスタ 7 0 1 中の全ての値が必要なわけではないことに注

50

意されたい。例えば、2つのクラスタしかサポートしないシステムにおいては、クラスタIDを維持するためには単一のレジスタしか提供されない。同じように、コンフィギュレーションスペースレジスタ701は、特定のポリング状態情報のような他の情報も維持しえる。

【0051】

図8は、システムのリセット中にコンフィギュレーションスペースレジスタを用いるある技術を示すフロープロセス図である。801において、物理レイヤ通信がリセットからイネーブルされる。物理レイヤをイネーブルすると同時に標準物理レイヤ初期化シーケンスが実行されえる。ポリング状態も維持されえる。ある例では、ポリング状態がアクティブに設定され、相互接続コントローラは、他のリソースの存在のためにリンクをモニタするだけでなく、ビットのシーケンスをアクティブに送る。他の例では、ポリング状態がパッシブに設定され、相互接続コントローラは単に他のリソースの存在だけをモニタする。803において、フェンスビットは、リンクレイヤ通信をディセーブルするように設定される。リンクレイヤ通信をディセーブルすることは、データローカル送信からの干渉なしに、物理レイヤ通信が確立されることを可能にする。しかし、リンクレイヤ通信は必ずしもディセーブルされる必要はないことを注意されたい。805において、初期化シーケンスまたはトレーニングシーケンスが送信される。

【0052】

送信は、シリアライゼーションおよびデシリアライゼーション(SERDES)メカニズムおよび8b/10bエンコーディングを用いて典型的には実行される。8b/10bエンコーディングは、1または0の連続したシーケンスを混合されたシーケンスで置換し、さまざまな理由のために用いられる。8b/10bは、正しいPLL動作を確実にするために送信されたデータ中に十分な信号遷移を発生するのを助ける。この8b/10bエンコーディングスキームなしでは、データ中の1または0のストリングが、クロックがドリフトまたは同期を失い、データの損失をきたすことを起こしえる。8b/10bは、信号がDC的にバランスがとれていることも確実にし、これはDCオフセットが時間と共にリンク内で大きくなっていかないことを意味する。8b/10bエンコードされた文字は、多くの信号の誤りをすぐに検出する特定の規則に従う。

【0053】

用いられえる初期化またはトレーニングシーケンスの例は、Infiniband Trade Associationから入手可能なInfiniband Architecture Specification Volumes I and IIにおいて記載され、その全体が全ての目的のために参照によって援用される。Infinibandは必ずしも用いられなければならないのではないことに注意されたい。物理レイヤ通信の初期化のあいだ、情報は相互接続コントローラ群間で交換される。811において、リンク幅パラメータが交換される。813において、リンクスピードパラメータが交換される。815において、オプションの誤り訂正情報が交換される。817において、フェンスビットをトグルオフすることによってリンクレイヤがそれからイネーブルされる。リンクレイヤがイネーブルされると、データ通信が進められる。819において、クラスタID情報が交換される。821において、クラスタIDがコンフィギュレーションスペースレジスタにおいて設定される。823において、接続されたリソースを反映するように相互接続コントローラに関連付けられるルーティングテーブルが必要に応じてアップデートされる。

【0054】

他の例において、複数のクラスタにわたってのリセットの同期化については仮定はなされていない。システムリセット信号の任意の時間的順序のアサーションおよびデアサーションが有効である。ある例では、それぞれのクラスタ中のそれぞれのリモートリンクレイヤは、リセットのデアサーションの後、フェンスビットがクリアにされて立ち上がる。それぞれのリモートリンクは、ディセーブルされて立ち上がり、リモート物理レイヤは、接続の他端上に他のクラスタがないかポリングする。リモート接続の両端上のクラスタがリセットから出た後、リモート物理レイヤはポーリングと共に互いを検出する。さまざまな実施形態によれば、リモート物理レイヤは、物理リンク接続の信頼性のある確立のために

10

20

30

40

50

必要とされる任意の他の情報と共に、リンク幅パラメータ、リンクスピードパラメータ、誤り訂正情報を交換することを伴う初期化/トレーニングシーケンスをここで実行しえる。

【 0 0 5 5 】

初期化/トレーニングシーケンスの後、クラスタIDフィールドが交換されえる。クラスタはリセットから出てきたばかりなので、クラスタIDは、所定の固定値に、例えば全てゼロに初期化されえる。それからリモートクラスタIDは、CSR中に維持される。物理レイヤはイネーブルされ、リンクレイヤはまだディセーブルされている。この時点で、サービスプロセッサは、主クラスタ中のリモートリンクのCSRを読み、クラスタIDを割り当て、かつリモート接続の物理レイヤを再初期化するためにCSRに書き込むことによって、主クラスタから始まってシステム全体の探索をしている。

10

【 0 0 5 6 】

さまざまな実施形態によれば、主クラスタ中のサービスプロセッサは、どのクラスタがどの他のクラスタにそれぞれのクラスタ上のリモートリンクを通して接続されているかについての情報をここで得る。この情報は、ルーティングテーブルを生成するのに用いられる。この情報は、BIOSへも渡される。BIOSは、クラスタのうちの一つの中のCSR中の「リンクレイヤを初期化」ビットをセットする。これは、リンクレイヤスタートパケットで片方のリンクレイヤがリンクレイヤ初期化を開始するようにする。他端のリンクレイヤがリンクレイヤスタートパケットを受け取るとき、それはリンクレイヤスタートパケットで応答し、リンクレイヤが確立される。

20

【 0 0 5 7 】

リンクレイヤが立ち上がると、リモートクラスタ中のプロセッサおよびコントローラをプログラム/初期化するためにパケットがリモートリンクをわたって主クラスタから送られえる。これでシステムが動作している。ある例では、もし任意のリモート接続が用いられないなら、サービスプロセッサまたはBIOSは、これらリモート接続についてリンクレイヤを分離し、物理レイヤをディセーブルすることを選びえる。

【 0 0 5 8 】

図9は、ホットプラグのあいだのスペースレジスタの使用を示すフロープロセス図である。システムがまだ動作しているあいだのリソースのシステムへの任意の動的付加は、ここではホットプラグと呼ばれる。901において、ポリング状態が相互接続コントローラにおいて維持される。903において、そのリンクレイヤをディセーブルするためにフェンスビットが維持される。しかしフェンスビットは必ずしもディセーブルされる必要はないことに注意されたい。新しい処理クラスタが検出されるとき、再初期化が905においてトリガされる。ある例では、再初期化は、コンフィギュレーションスペースレジスタの再初期化指示子であるビットに書き込むことによってトリガされる。907において、初期化シーケンスまたはトレーニングシーケンスが送られる。911において、リンク幅およびリンクスピードのような情報が交換される。913においてそれからリンクレイヤがイネーブルされる。915において、プロセッサの新しく追加されたグループのクラスタIDがコンフィギュレーションスペースレジスタに書き込まれる。917において上述のルーティングテーブルは、複数のクラスタの通信を可能にするために同様にアップデートされる。

30

40

【 0 0 5 9 】

他の例において、1つ以上のクラスタが動作し、新しいクラスタがシステムに追加される必要がある。さまざまな実施形態によれば、追加される必要がある新しいクラスタ上でリセットがアサートおよびディアサートされる。もし新しいクラスタへのリモート接続についての物理レイヤがディセーブルされるなら、物理レイヤはポリング状態に置かれる。リモート接続は分離されない(unfenced)。クラスタIDが交換され、物理レイヤはイネーブルされる。サービスプロセッサは、上述のように同様のやり方で新しいクラスタおよび主クラスタをもプログラムするために必要とされる情報を生成する。新しいリモート接続のリンクレイヤがここで初期化される。

50

【0060】

図10は、システムからリソースを動的に除去する技術を示すフロープロセス図である。システムがアクティブであるときに、プロセッサまたは処理クラスタのようなリソースを動的にアンプラグすることは一般的に難しいが、これはプロセッサおよびそれらに関連付けられたキャッシュを除去することがシステム動作を非常に大きく左右するからである。ある例では、特定のアプリケーションによって処理されているデータは、システムから除去されるためにプロセッサキャッシュ中において中間状態で保持されえる。同様に、オペレーティングシステムプロセスは、除去が設定された1つ以上のプロセッサ上で走っているかもしれない。さまざまな実施形態によれば、クラスタの動的除去をサポートするオペレーティングシステムが用いられる。ある例では、アプリケーションは、マルチプロセッサクラスタの動的除去をサポートするオペレーティングシステム環境中で1001において終了される。クラスタが1003においてディセーブルされるが、まだアンプラグされない。1005において、除去されるべきそのクラスタに関連付けられるキャッシュがフラッシュされる。1007において、処理クラスタの除去を反映するようにルーティングテーブルが変更される。1009において、フェンスビットがコンフィギュレーションスペースレジスタ中に書き込まれる。フェンスビットは、さまざまなエンティティによって書き込まれえる。ある例では、フェンスビットは、サービスプロセッサによって、またはプロセッサに関連付けられるJTAGインタフェースによって書き込まれる。1011において、マルチプロセッサクラスタが物理に除去されえる。1013において、プロセッサの新しいまたは置換クラスタが導入されることを可能にするために、まだシステム中に存在する相互接続コントローラにおける物理レイヤが維持される。

10

20

【0061】

本発明は、その具体的な実施形態を参照して示され記述されてきたが、記載された実施形態の形態および詳細の変更が本発明の精神または範囲から逸脱することなくなされることが当業者には理解されよう。例えば本発明の実施形態は、ポイントツーポイント、スイッチ、またはバスアーキテクチャを通して接続されたマルチプロセッサクラスタと共に採用されえる。他の例では、プロセッサのマルチクラスタは、単一の相互接続コントローラを共有しえ、または複数の相互接続コントローラが単一のクラスタ中で用いられえる。従って本発明の範囲は添付の特許請求の範囲を参照して決定されるべきである。

30

【図面の簡単な説明】

【0062】

【図1A】複数のクラスタを有するシステムを示す概略図である。

【図1B】複数のクラスタを有するシステムを示す概略図である。

【図2】複数のプロセッサを有する例示的クラスタを示す概略図である。

【図3】本発明のさまざまな実施形態を利用する例示的相互接続コントローラの概略図である。

【図4】ローカルプロセッサの概略図である。

【図5】本発明の特定の実施形態によるメモリマッピングスキームの概略図である。

【図6A】本発明の特定の実施形態を示す4つのクラスタシステムの簡略化されたブロック図である。

40

【図6B】図6Aの4つのクラスタシステムのためのルーティング情報を含む結合されたルーティングテーブルの図である。

【図7】相互接続コントローラに関連付けられるコンフィギュレーションスペースレジスタを示す概略図である。

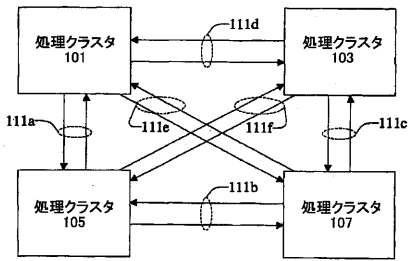
【図8】システムのリスタートを示すフロープロセス図である。

【図9】新しいリンクの動的導入を示すフロープロセス図である。

【図10】既存リンクの動的除去を示すフロープロセス図である。

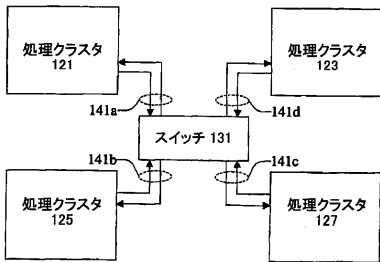
【 図 1 A 】

Fig. 1A



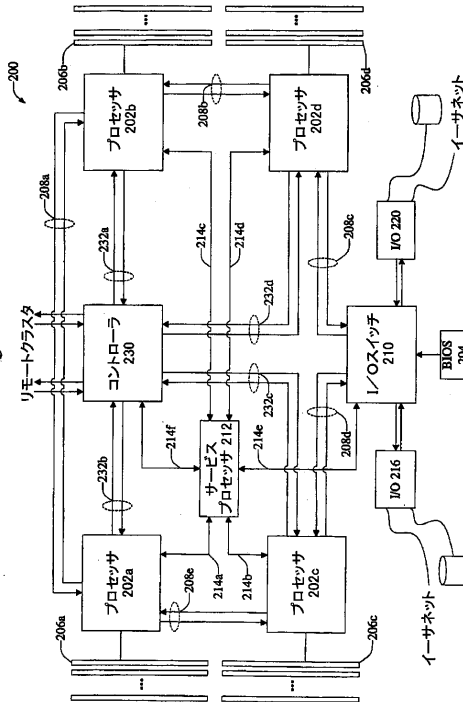
【 図 1 B 】

Fig. 1B



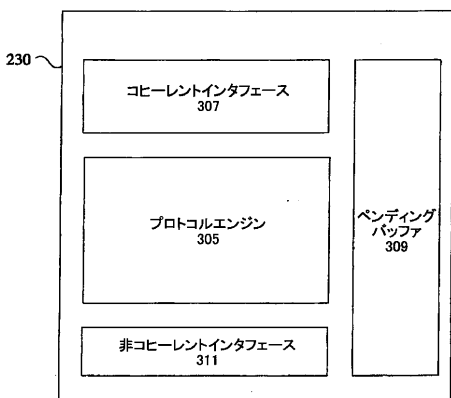
【 図 2 】

Fig. 2



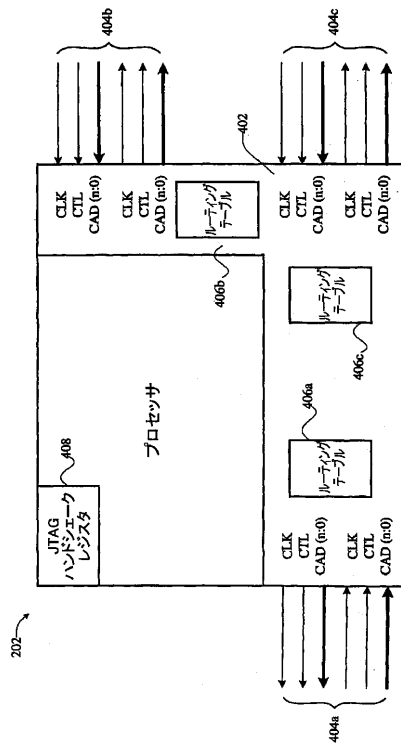
【 図 3 】

Fig. 3



【 図 4 】

Fig. 4



【 図 5 】

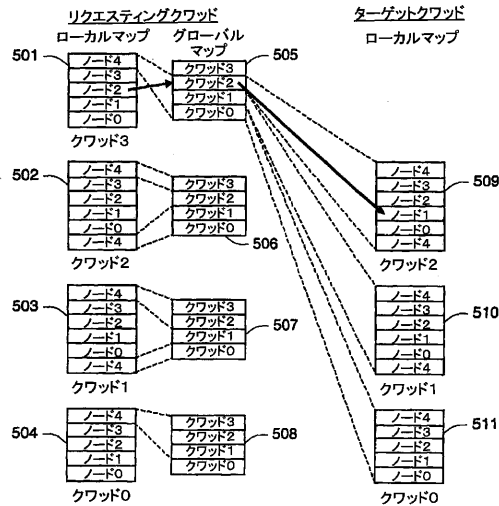


Fig. 5

【 図 6 A 】

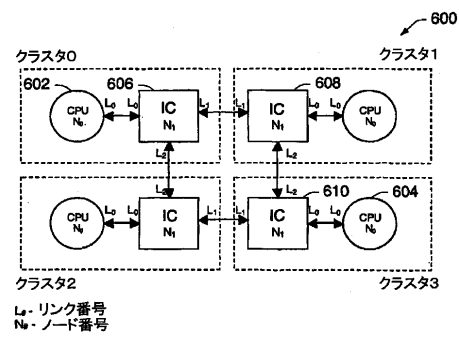


Fig. 6A

【 図 6 B 】

	ローカルテーブル		グローバルテーブル			
	宛先ノード		宛先クラスタ			
ソース	N ₀	N ₁	C ₀	C ₁	C ₂	C ₃
クラスタ0 ノード0	X	L ₀	NA	NA	NA	NA
ノード1	L ₂	X	X	L ₁	L ₂	L ₁
クラスタ1 ノード0	X	L ₂	NA	NA	NA	NA
ノード1	L ₂	X	L ₁	X	L ₂	L ₂
クラスタ2 ノード0	X	L ₂	NA	NA	NA	NA
ノード1	L ₂	X	L ₂	L ₂	X	L ₁
クラスタ3 ノード0	X	L ₀	NA	NA	NA	NA
ノード1	L ₀	X	L ₂	L ₂	L ₁	X

Fig. 6B

【 図 7 】

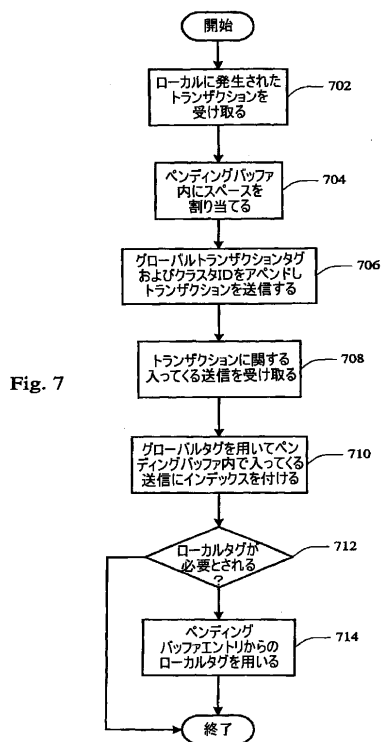


Fig. 7

【 図 8 】

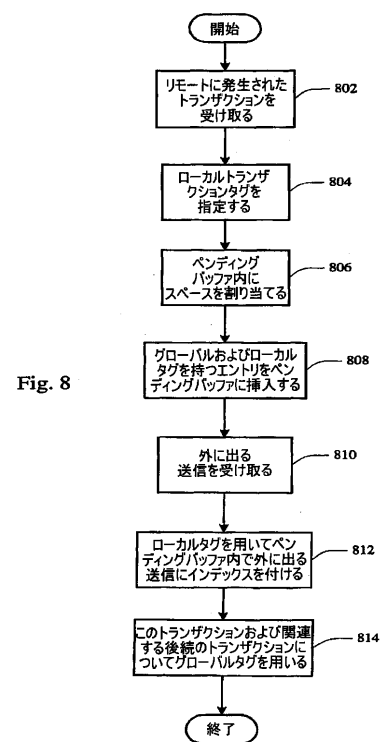
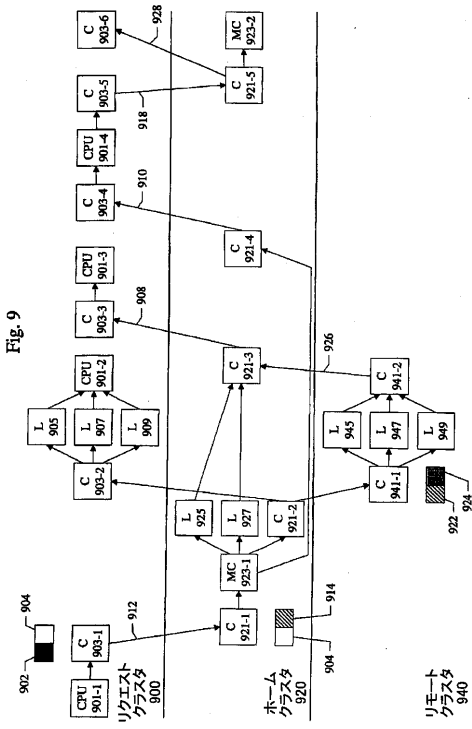


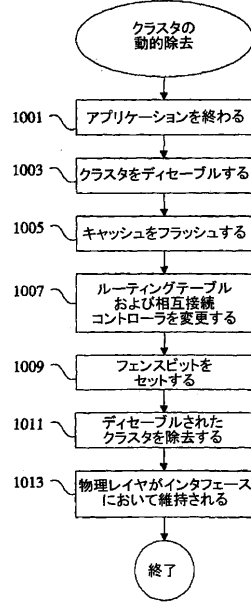
Fig. 8

【 図 9 】



【 図 10 】

Figure 10



フロントページの続き

(81) 指定国 AP(BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), EP(AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PL, PT, RO, SE, SI, SK, TR), OA(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG), AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW

(72) 発明者 ネワワーカー・シャシャンク

アメリカ合衆国 テキサス州 7 8 7 1 7 オースティン, コルデロ・ドライブ, 1 5 0 0 8

(72) 発明者 ブラサドゥ・グル

アメリカ合衆国 テキサス州 7 8 7 5 0 オースティン, カラニッシュ・パーク・ドライブ, 1 1 1 1 1

(72) 発明者 ツァイトラー・カール

アメリカ合衆国 テキサス州 7 7 3 7 5 トムボール, ブラッシュ・キャニオン・ドライブ, 1 1 8 3 5

(72) 発明者 グラスコ・デイビッド・ビー

アメリカ合衆国 テキサス州 7 8 7 2 6 オースティン, エンバー・グレン・ドライブ, 1 0 3 3 7

Fターム(参考) 5B045 BB02 BB15 BB28 DD01 HH05 HH06 JJ38