

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5084197号
(P5084197)

(45) 発行日 平成24年11月28日(2012.11.28)

(24) 登録日 平成24年9月14日(2012.9.14)

(51) Int.Cl.

F I

G 0 6 F 15/173 (2006.01)

G 0 6 F 15/173 6 4 0 M

請求項の数 12 (全 22 頁)

(21) 出願番号	特願2006-217953 (P2006-217953)	(73) 特許権者	310021766
(22) 出願日	平成18年8月10日 (2006. 8. 10)		株式会社ソニー・コンピュータエンタテインメント
(65) 公開番号	特開2008-41027 (P2008-41027A)		東京都港区港南1丁目7番1号
(43) 公開日	平成20年2月21日 (2008. 2. 21)	(74) 代理人	100105924
審査請求日	平成21年7月31日 (2009. 7. 31)		弁理士 森下 賢樹
		(74) 代理人	100109047
			弁理士 村田 雄祐
		(74) 代理人	100109081
			弁理士 三木 友由
		(74) 代理人	100134256
			弁理士 青木 武司

最終頁に続く

(54) 【発明の名称】 プロセッサノードシステムおよびプロセッサノードクラスタシステム

(57) 【特許請求の範囲】

【請求項1】

プロセッサと、前記プロセッサの入出力バスと周辺デバイスが接続されるPCIエクスプレスとの間でデータを中継するブリッジとが搭載されたプロセッサ基板を複数含み、

前記ブリッジのポートは当該プロセッサがホストとなるルートコンプレックスモードまたは当該プロセッサが周辺デバイスとなるエンドポイントモードに設定可能に構成され、

一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートを、別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに接続することにより、前記複数のプロセッサ基板間を相互結合してなることを特徴とするプロセッサノードシステム。

【請求項2】

前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられるPCIエクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられるPCIエクスプレスコネクタとがフレキシブル基板により配線接続されてなることを特徴とする請求項1に記載のプロセッサノードシステム。

【請求項3】

前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられるPCIエクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられるPCIエクスプレスコネクタとを相互

接続するための一枚のバックプレーン基板をさらに設けたことを特徴とする請求項 1 に記載のプロセッサノードシステム。

【請求項 4】

プロセッサと、前記プロセッサの入出力バスと周辺デバイスが接続される P C I エクスプレスとの間でデータを中継するブリッジとが 2 組搭載されたプロセッサ基板を 4 枚含み、

各ブリッジは、当該プロセッサがホストとなるルートコンプレックスモードに設定されたポートと当該プロセッサが周辺デバイスとなるエンドポイントモードに設定されたポートを有し、

一のプロセッサ基板のルートコンプレックスモードに設定されたポートは、別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに接続されることを条件として、各プロセッサ基板の合計 4 個のポートの内、3 個のポートを用いて、前記 4 枚のプロセッサ基板の内、任意の 2 枚のプロセッサ基板間を相互結合してなることを特徴とするプロセッサノードシステム。

10

【請求項 5】

各プロセッサ基板の前記 3 個のポートを第 1 ポート、第 2 ポート、第 3 ポートとし、当該プロセッサ基板以外の 3 枚のプロセッサ基板を、第 1 プロセッサ基板、第 2 プロセッサ基板、第 3 プロセッサ基板とした場合、前記第 1 ポートは前記第 1 プロセッサ基板に接続され、前記第 2 ポートは前記第 2 プロセッサ基板に接続され、前記第 3 ポートは前記第 3 プロセッサ基板に接続されることを特徴とする請求項 4 に記載のプロセッサノードシステム。

20

【請求項 6】

前記プロセッサ基板内の 2 つのプロセッサは入出力バスを介して接続されており、前記プロセッサ基板間のポートの接続により、4 枚の前記プロセッサ基板の合計 8 個のプロセッサの内、任意の 2 個のプロセッサが互いに通信可能に相互接続されることを特徴とする請求項 5 に記載のプロセッサノードシステム。

【請求項 7】

前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられる P C I エクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられる P C I エクスプレスコネクタとがフレキシブル基板により配線接続されてなることを特徴とする請求項 4 から 6 のいずれかに記載のプロセッサノードシステム。

30

【請求項 8】

筐体内に前記 4 枚のプロセッサ基板面を互いに平行に設置し、筐体背面に各プロセッサ基板のブリッジのポートに設けられる P C I エクスプレスコネクタが配置されるように構成し、筐体背面に配置された各プロセッサ基板の P C I エクスプレスコネクタ間をフレキシブル基板により接続してなることを特徴とする請求項 7 に記載のプロセッサノードシステム。

【請求項 9】

筐体内に前記 4 枚のプロセッサ基板面を互いに平行に設置し、筐体背面に各プロセッサ基板のブリッジのポートに設けられる P C I エクスプレスコネクタが配置されるように構成し、前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられる P C I エクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられる P C I エクスプレスコネクタとを相互接続するための一枚のバックプレーン基板をさらに設けたことを特徴とする請求項 4 から 6 のいずれかに記載のプロセッサノードシステム。

40

【請求項 10】

各プロセッサのメモリ空間に、相互接続された他のプロセッサの共有領域が I / O アドレス空間としてメモリマッピングされることにより、各プロセッサは前記他のプロセッサの共有領域にアクセス可能に構成されることを特徴とする請求項 4 から 9 のいずれかに記

50

載のプロセッサノードシステム。

【請求項 1 1】

請求項 4 から 1 0 のいずれかに記載のプロセッサノードシステムを複数含み、各プロセッサノードシステムをクラスタとして相互接続した、クラスタの接続形態を有するプロセッサノードクラスタシステムであって、

隣接する 2 つのプロセッサノードシステム間で当該プロセッサノードシステム内のプロセッサ基板間の接続に使用されていない空きポートを互いに接続することにより、前記複数のプロセッサノードシステム間を相互結合してなることを特徴とするプロセッサノードクラスタシステム。

【請求項 1 2】

前記空きポートに P C I エクスプレスコネクタが設けられ、前記空きポートの P C I エクスプレスコネクタ間がフレキシブル基板により配線接続されてなることを特徴とする請求項 1 1 に記載のプロセッサノードクラスタシステム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

この発明は、複数のプロセッサを相互接続したプロセッサノードシステムおよびプロセッサノードクラスタシステムに関する。

【背景技術】

【0 0 0 2】

パーソナルコンピュータやサーバには、P C I (Peripheral Component Interconnect) バスを介して各種の周辺デバイスが接続され、情報処理システムが構成される。プロセッサの入出力バスと、周辺デバイスの入出力バスである P C I バスとは規格が異なるため、通常、ブリッジを介してプロセッサと周辺デバイスとが接続される。

【0 0 0 3】

情報処理システムの機能拡張や性能強化を図るために、グラフィックプロセッサや高速なメモリデバイスを P C I デバイスとして接続することがあり、より多くの周辺デバイスを P C I バスで接続できるようにすることが要請されている。そのため、P C I エクスプレス (PCI Express) (商標または登録商標) スイッチを用いて、一つのプロセッサに対して複数のデバイスを接続することが行われている。また、複数のプロセッサノードを相互接続したり、プロセッサノードとデバイスを相互接続するために、I n f i n i b a n d と呼ばれる超高速インタフェース技術が用いられることがある。

【発明の開示】

【発明が解決しようとする課題】

【0 0 0 4】

1 0 ギガビットイーサネット (商標または登録商標) や I n f i n i b a n d 技術を用いて複数のプロセッサノードを相互接続したクラスタシステムでは、プロセッサ間の高速な通信を実現することができるという利点があるが、スイッチが未だ高価であるため、クラスタシステムを低価格で提供することは難しく、クラスタ内のプロセッサノード数を増やしていくには限界がある。さらに、イーサネット (商標または登録商標) や I n f i n i b a n d では、パケットの生成、プロトコル処理などソフトウェアのオーバーヘッドが大きいというデメリットがある。

【0 0 0 5】

本発明はこうした課題に鑑みてなされたものであり、その目的は、複数のプロセッサを安価な手段により結合して、高速なプロセッサ間通信を実現する技術およびその技術を利用したプロセッサノードシステムやプロセッサノードクラスタシステムを提供することにある。

【課題を解決するための手段】

【0 0 0 6】

上記課題を解決するために、本発明のある態様のプロセッサノードシステムは、プロセ

10

20

30

40

50

ッサと、前記プロセッサの入出力バスと周辺デバイスが接続されるPCIエクスプレスとの間でデータを中継するブリッジとが搭載されたプロセッサ基板を複数含む。前記ブリッジのポートは当該プロセッサがホストとなるルートコンプレックスモードまたは当該プロセッサが周辺デバイスとなるエンドポイントモードに設定可能に構成され、一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートを、別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに接続することにより、前記複数のプロセッサ基板間が相互結合される。

【0007】

前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられるPCIエクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられるPCIエクスプレスコネクタとがフレキシブル基板により配線接続されてもよい。

10

【0008】

前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられるPCIエクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられるPCIエクスプレスコネクタとを相互接続するための一枚のバックプレーン基板をさらに設けてもよい。

【0009】

本発明の別の態様もまた、プロセッサノードシステムである。このプロセッサノードシステムは、プロセッサと、前記プロセッサの入出力バスと周辺デバイスが接続されるPCIエクスプレスとの間でデータを中継するブリッジのセットが2組搭載されたプロセッサ基板を4枚含む。各ブリッジは、当該プロセッサがホストとなるルートコンプレックスモードに設定されたポートと当該プロセッサが周辺デバイスとなるエンドポイントモードに設定されたポートを有する。一のプロセッサ基板のルートコンプレックスモードに設定されたポートは、別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに接続されることを条件として、各プロセッサ基板の合計4個のポートの内、3個のポートを用いて、前記4枚のプロセッサ基板の内、任意の2枚のプロセッサ基板間を相互結合される。

20

【0010】

筐体内に前記4枚のプロセッサ基板面を互いに平行に設置し、筐体背面に各プロセッサ基板のブリッジのポートに設けられるPCIエクスプレスコネクタが配置されるように構成し、筐体背面に配置された各プロセッサ基板のPCIエクスプレスコネクタ間をフレキシブル基板により接続してもよい。

30

【0011】

筐体内に前記4枚のプロセッサ基板面を互いに平行に設置し、筐体背面に各プロセッサ基板のブリッジのポートに設けられるPCIエクスプレスコネクタが配置されるように構成し、前記一のプロセッサ基板のブリッジのルートコンプレックスモードに設定されたポートに設けられるPCIエクスプレスコネクタと、前記別のプロセッサ基板のブリッジのエンドポイントモードに設定されたポートに設けられるPCIエクスプレスコネクタとを相互接続するための一枚のバックプレーン基板をさらに設けてもよい。

40

【0012】

本発明のさらに別の態様は、プロセッサノードクラスタシステムである。このプロセッサノードクラスタシステムは、プロセッサノードシステムを複数含む。隣接する2つのプロセッサノードシステム間で当該プロセッサノードシステム内のプロセッサ基板間の接続に使用されていない空きポートを互いに接続することにより、前記複数のプロセッサノードシステム間が相互結合される。

【0013】

なお、以上の構成要素の任意の組合せ、本発明の表現を方法、装置、システム、コンピュータプログラム、データ構造、記録媒体などの間で変換したものもまた、本発明の態様として有効である。

50

【発明の効果】

【0014】

本発明によれば、複数のプロセッサを相互接続して安価で高性能なシステムを構成することができる。

【発明を実施するための最良の形態】

【0015】

実施の形態に係るクラスタシステムは、プロセッサが搭載された基板（ボード）をフレキシブル基板で密結合することにより構成される。図1を参照して、各プロセッサ基板の構成を説明し、図2を参照して、4つのプロセッサ基板をフレキシブル基板により密結合したノードの構成を説明する。図3を参照して、複数のノード間をフレキシブル基板により連結することにより構成されるクラスタシステムを説明する。また、図4～図8を参照して、プロセッサ基板間をフレキシブル基板により接続する形態について説明する。

10

【0016】

図1は、プロセッサ基板50の構成図である。プロセッサ基板50には、2つのマルチコアプロセッサ（Multicore Processor）（以下、「MCP」と呼ぶ）20、21が搭載されている。各MCP20、21は、複数のプロセッサコアを1つのパッケージに集積したものであり、プロセッサコアとして、1つのプロセッシングエレメント（PE）と、複数のサブプロセッシングエレメント（SPE）を含む。PEは、キャッシュメモリを有し、DRAM10から読み込んだデータをキャッシュしながら、情報処理を行う。また、PEは、各MCP20、21全体を統括的に制御する。各SPEはローカルメモリを内部にもち、ローカルメモリに対してデータを読み書きしながら、情報処理を行う。複数のSPEは非同期で動作する。

20

【0017】

2つのMCP20、21は入出力インタフェース（以下、「IOIF」と呼ぶ）64を介して相互に接続されており、高速なデータ通信が可能である。さらに、各MCP20、21は、IOIF62、63を介してブリッジ30、31の上流（アップストリーム）ポートに接続されている。ブリッジ30、31の下流（ダウンストリーム）ポートには、PCIエクスプレス66、67を介して各種の周辺（ペリフェラル）デバイスや他のプロセッサ基板が接続される。

【0018】

ここで、PCIエクスプレス66、67は、PCIエクスプレス（PCI Express）（商標または登録商標）の仕様にしたがうものであるが、現行のPCIエクスプレス規格に限定する趣旨ではなく、現行のPCIエクスプレス規格に準拠するものや、現行のPCIエクスプレス規格をさらに拡張したり、発展させた規格によるものであってもかまわない。PCIエクスプレス66、67で接続された周辺デバイスや他のプロセッサ基板を以下、「PCIデバイス」という。

30

【0019】

IOIF62、63、64は、上りと下りの2つのチャンネルをもち、メモリバスに匹敵する高い帯域幅、たとえば、数十ギガバイト/秒を実現している。各MCP20、21の所定のメモリ領域は、IOIF62、63、64を介して参照可能なI/Oアドレス空間にメモリマッピングされる。各MCP20、21は、IOIF62、63、64を介してI/Oアドレス空間にマッピングされた他のMCPのメモリ領域にアクセスすることが可能であり、高速なプロセッサ間通信が実現される。

40

【0020】

各ブリッジ30、31は、IOIF62、63とPCIエクスプレス66、67とを「橋渡し」することで、MCP20、21とPCIデバイスとを相互接続する。IOIF62、63と、PCIエクスプレス66、67とは、バスの規格が異なるため、ブリッジ30、31は、2つのバスの間でプロトコルの変換を行い、MCP20、21とPCIデバイスとがやりとりするデータのフォーマットを各バスの仕様に合わせる。

【0021】

50

PCIエクスプレス66、67に接続されたPCIデバイスの先にさらにPCIエクスプレスを介してPCIデバイスを接続していくと、MCP20、21をルート(根)とし、リーフ(葉)にはPCIデバイスが接続されたPCIデバイスのツリー(木)構造が形成される。以下、このPCIデバイスのツリー構造を「PCIツリー」という。

【0022】

各ブリッジ30、31の下流ポートは2つ設けられており、一方は、ルートコンプレックス(RC; Root Complex)として、他方は、エンドポイント(EP; Endpoint)としてコンフィギュレーションして用いることができる。1つのポートがルートコンプレックスモードとエンドポイントモードを切り替えられるように構成されてもよい。ブリッジ30、31の下流ポートをルートコンプレックスとして用いると、MCP20、21は、PCIツリーのルートとなって、PCIデバイスを接続するホストとして機能する。ブリッジ30、31の下流ポートをエンドポイントとして用いると、MCP20、21は、ホストに接続されるPCIデバイスとして機能する。

10

【0023】

ブリッジ30、31の下流ポートには、ルートコンプレックス用のコネクタ40R、41Rと、エンドポイント用のコネクタ40E、41Eとが設けられる。本実施の形態では、プロセッサ基板50に設けられた合計4個のコネクタ40R、40E、41R、41Eの内、3個のコネクタを同一ノード内の他の3つのプロセッサ基板との接続に用い、残りの1個のコネクタを他のノードのプロセッサ基板との接続に用いる。

【0024】

図2は、4つのプロセッサ基板が密結合されたノード100の構成図である。ノード100は、第1プロセッサ基板50、第2プロセッサ基板51、第3プロセッサ基板52、および第4プロセッサ基板53をフレキシブル基板で相互に接続したものである。各プロセッサ基板50～53の構成は、図1で説明した通りである。

20

【0025】

以下、第1プロセッサ基板50、第2プロセッサ基板51、第3プロセッサ基板52、第4プロセッサ基板53をそれぞれ「プロセッサ基板0」、「プロセッサ基板1」、「プロセッサ基板2」、「プロセッサ基板3」と呼ぶ。

【0026】

プロセッサ基板0に搭載された2つのMCP20、21をそれぞれ「MCP0」、「MCP1」と呼び、MCP0、MCP1に接続されたブリッジ30、31をそれぞれ「ブリッジ0」、「ブリッジ1」と呼ぶ。同様に、プロセッサ基板1に搭載された2つのMCP22、23をそれぞれ「MCP2」、「MCP3」と呼び、MCP2、MCP3に接続されたブリッジ32、33をそれぞれ「ブリッジ2」、「ブリッジ3」と呼ぶ。プロセッサ基板2に搭載された2つのMCP24、25をそれぞれ「MCP4」、「MCP5」と呼び、MCP4、MCP5に接続されたブリッジ34、35をそれぞれ「ブリッジ4」、「ブリッジ5」と呼ぶ。プロセッサ基板3に搭載された2つのMCP26、27をそれぞれ「MCP6」、「MCP7」と呼び、MCP6、MCP7に接続されたブリッジ36、37をそれぞれ「ブリッジ6」、「ブリッジ7」と呼ぶ。

30

【0027】

各プロセッサ基板内の2つのMCPを相互接続するIOIFを「IOIF0」と呼び、MCPとブリッジの上流ポート間のIOIFを「IOIF1」と呼ぶ。

40

【0028】

ブリッジ0のRC用コネクタ、EP用コネクタをそれぞれ「コネクタRC0」、「コネクタEP0」と呼ぶ。同様に、ブリッジ1～ブリッジ7のRC用コネクタをそれぞれ「コネクタRC1」～「コネクタRC7」と呼び、ブリッジ1～ブリッジ7のEP用コネクタをそれぞれ「コネクタEP1」～「コネクタEP7」と呼ぶ。

【0029】

プロセッサ基板0のMCP0側のブリッジ0のコネクタRC0は、プロセッサ基板1のMCP3側のブリッジ3のコネクタEP3と接続される。この接続により、プロセッサ基

50

板 0 の M C P 0 から見た場合、M C P 0 はルートコンプレックスとして機能し、プロセッサ基板 1 の M C P 3 はエンドポイントとして機能する。すなわち、プロセッサ基板 0 の M C P 0 はホストであり、プロセッサ基板 1 を P C I デバイスとして接続した形態となり、M C P 0 をルートとして M C P 3 をつないだ P C I ツリーが形成される。

【 0 0 3 0 】

プロセッサ基板 0 の M C P 1 側のブリッジ 1 のコネクタ R C 1 は、プロセッサ基板 2 の M C P 4 側のブリッジ 4 のコネクタ E P 4 と接続される。ルートコンプレックスであるプロセッサ基板 0 の M C P 1 から見た場合、プロセッサ基板 0 の M C P 1 をホスト、プロセッサ基板 2 をデバイスとする P C I ツリーが形成される。

【 0 0 3 1 】

プロセッサ基板 0 の M C P 1 側のブリッジ 1 のコネクタ E P 1 は、プロセッサ基板 3 の M C P 6 側のブリッジ 6 のコネクタ R C 6 と接続される。ルートコンプレックスであるプロセッサ基板 3 の M C P 6 から見た場合、プロセッサ基板 3 の M C P 6 をホスト、プロセッサ基板 0 をデバイスとする P C I ツリーが形成される。

【 0 0 3 2 】

同様に、プロセッサ基板 1 の M C P 2 側のブリッジ 2 のコネクタ R C 2 は、プロセッサ基板 2 の M C P 5 側のブリッジ 5 のコネクタ E P 5 と接続され、ブリッジ 2 のコネクタ E P 2 は、プロセッサ基板 3 の M C P 7 側のブリッジ 7 のコネクタ R C 7 と接続される。プロセッサ基板 2 の M C P 4 側のブリッジ 4 のコネクタ R C 4 は、プロセッサ基板 3 の M C P 7 側のブリッジ 7 のコネクタ E P 7 と接続される。

【 0 0 3 3 】

ノード 1 0 0 内のプロセッサ基板間の接続に用いられないブリッジのコネクタ、すなわち、プロセッサ基板 0 のブリッジ 0 のコネクタ E P 0、プロセッサ基板 1 のブリッジ 3 のコネクタ R C 3、プロセッサ基板 2 のブリッジ 5 のコネクタ R C 5、およびプロセッサ基板 3 のブリッジ 6 のコネクタ E P 6 は、空きスロットとして、他のノードのプロセッサ基板との接続に利用される。

【 0 0 3 4 】

図 3 は、複数のノードを連結したクラスタシステム 2 0 0 の構成図である。クラスタシステム 2 0 0 は、図 2 で説明した構成のノード 1 0 0 ~ 1 0 2、1 1 0 ~ 1 1 2、1 2 0 ~ 1 2 0 を上下左右に連結したものである。たとえば、ノード 1 0 0 の右にはノード 1 0 1 が接続され、ノード 1 0 1 のさらに右にはノード 1 0 2 が接続される。ノード 1 0 0 の下にはノード 1 1 0 が接続され、ノード 1 1 0 のさらに下にはノード 1 2 0 が接続される。

【 0 0 3 5 】

同図に示すように、左右に並ぶ 2 つのノードは、左側のノードのプロセッサ基板 3 のコネクタ E P 6 と、右側のノードのプロセッサ基板 1 のコネクタ R C 3 とを接続することにより、結合される。上下に並ぶ 2 つのノードは、上側のノードのプロセッサ基板 2 のコネクタ R C 5 と、下側のノードのプロセッサ基板 0 のコネクタ E P 0 とを接続することにより、結合される。

【 0 0 3 6 】

クラスタシステム 2 0 0 において、端部に位置するノードの隣接ノードが存在しない側のコネクタは空きスロットになるが、この空きスロットには各種の周辺デバイスを接続したり、さらにノードを接続することにより、システムを拡張することができる。

【 0 0 3 7 】

このように、クラスタシステム 2 0 0 では、ノードを上下左右に結合する平面上の配置により、ノード数を自由自在に増やしていくことができるという利点がある。

【 0 0 3 8 】

クラスタシステム 2 0 0 において、各ノード内の 4 枚のプロセッサ基板間の接続、およびノード間の接続には、フレキシブル基板が用いられる。以下、図 4 ~ 図 8 を参照して、フレキシブル基板を用いた接続形態を説明する。

10

20

30

40

50

【 0 0 3 9 】

図4は、プロセッサ基板50の裏面の配線の模式図である。同図において、MCP0と複数のDRAM10の間の配線、MCP0とブリッジ0の間の配線、ブリッジ0とコネクタRC0、EP0の間の配線が示されている。また、MCP1と複数のDRAM11の間の配線、MCP1とブリッジ1の間の配線、ブリッジ1とコネクタRC1、EP1の間の配線が示されている。各コネクタRC0、EP0、RC1、EP1はPCI-Express x16コネクタであり、フレキシブル基板を接続することができる。

【 0 0 4 0 】

図5は、ノード100内の4枚のプロセッサ基板50～53間をフレキシブル基板によって接続した構成を示す図である。フレキシブル基板は、プリント配線基板の一種であり、FPC(Flexible Printed Circuit)とも呼ばれ、薄くて屈曲性がある。

10

【 0 0 4 1 】

図2で説明したプロセッサ基板50～53(プロセッサ基板0～3)を、フレキシブル基板による接続がしやすいように、プロセッサ基板1(符号51)、プロセッサ基板0(符号50)、プロセッサ基板2(符号52)、プロセッサ基板3(符号53)の順に、基板面を互いに平行にして配置する。

【 0 0 4 2 】

プロセッサ基板1のコネクタRC2は、フレキシブル基板201によりプロセッサ基板2のコネクタEP5と接続される。プロセッサ基板1のコネクタEP2は、フレキシブル基板202によりプロセッサ基板3のコネクタRC7と接続される。プロセッサ基板1のコネクタEP3は、フレキシブル基板203によりプロセッサ基板0のコネクタRC0と接続される。

20

【 0 0 4 3 】

プロセッサ基板0のコネクタRC1は、プロセッサ基板2のコネクタEP4とフレキシブル基板204によって接続される。プロセッサ基板0のコネクタEP1は、プロセッサ基板3のコネクタRC6とフレキシブル基板205によって接続される。プロセッサ基板2のコネクタRC4は、プロセッサ基板3のコネクタEP7とフレキシブル基板206によって接続される。

【 0 0 4 4 】

プロセッサ基板1のコネクタEP2とプロセッサ基板3のコネクタRC7をつなぐフレキシブル基板202は、プロセッサ基板0のコネクタRC1とプロセッサ基板2のコネクタEP4をつなぐフレキシブル基板204の上側をまたいでいる。このようにフレキシブル基板を用いれば、配線の上に別の配線が通るような接続形態も可能であり、4枚のプロセッサ基板を平行に並べて相互に密結合させ、省スペース化を図ることができる。

30

【 0 0 4 5 】

また、汎用品のPCIエクスプレスコネクタとフレキシブル基板を用いてプロセッサ基板間を接続する構成であるため、PCIエクスプレススイッチなどでプロセッサ基板間を相互接続した構成に比べて、はるかに安価であり、製造コストを削減することができる。

【 0 0 4 6 】

さらに、プロセッサ基板の部品実装密度を高くし、プロセッサ基板を小型化することによって、より短いフレキシブル基板でプロセッサ基板を相互接続することができ、高速信号を扱うことが可能になる。PCI-Expressは高速通信を前提としており、ケーブル接続では信号の伝搬が遅く、ケーブル接続によってプロセッサ基板間の密結合を実現することは困難である。本実施の形態では、フレキシブル基板でプロセッサ基板間を配線するため、高速信号の伝搬が可能である。

40

【 0 0 4 7 】

図6は、クラスタシステム200内の複数のノード間をフレキシブル基板によって接続した構成を示す図である。同図では、図3の4つの隣接するノード100、101、110、111の接続形態が示されている。各ノード100、101、110、111内の4枚のプロセッサ基板間は、図5で説明したようにフレキシブル基板で接続されている。た

50

だし、ノード110については、ノード間の接続形態を把握しやすくするため、ノード内のプロセッサ基板間を接続するフレキシブル基板を図示していない。

【0048】

ノード100のプロセッサ基板3のコネクタEP6は、ノード101のプロセッサ基板1のコネクタRC3とフレキシブル基板211によって接続される。これにより2つのノード100、101が左右方向に結合する。同様にノード110のプロセッサ基板3のコネクタEP6は、ノード111のプロセッサ基板1のコネクタRC3とフレキシブル基板221によって接続され、2つのノード110、111が左右方向に結合する。

【0049】

ノード100のプロセッサ基板2のコネクタRC5は、ノード110のフレキシブル基板0のコネクタEP0とフレキシブル基板214によって接続される。これにより2つのノード100、110が上下方向に結合する。同様に、ノード101のプロセッサ基板2のコネクタRC5は、ノード111のフレキシブル基板0のコネクタEP0とフレキシブル基板224によって接続され、2つのノード101、111が上下方向に結合する。

【0050】

クラスタシステム200では、ノード間の接続にもフレキシブル基板が用いられ、省スペース化とコストダウンを図ることができる。クラスタシステム200は、複数のノードを平面上で上下左右に配置して接続する形態であるため、隣り合うノード間の距離を短くすることができ、ノード間接続に用いるフレキシブル基板の長さを十分に短くすることができ、PCI-Expressの高速信号を扱うことができる。

【0051】

図7は、ノード100の筐体を説明する図である。ノード100の筐体には、4枚のプロセッサ基板50から53が収納されており、背面のコネクタ間は図5で説明したように6個のフレキシブル基板201~206で接続されている。さらに、図6で説明したように、プロセッサ基板0には、上方向に隣接するノードのプロセッサ基板2と接続するためのフレキシブル基板213が設けられ、プロセッサ基板2には、下方向に隣接するノードのプロセッサ基板0と接続するためのフレキシブル基板214が設けられる。一方、プロセッサ基板1には、左方向に隣接するノードのプロセッサ基板3と接続するためのフレキシブル基板212が設けられ、プロセッサ基板3には、右方向に隣接するノードのプロセッサ基板1と接続するためのフレキシブル基板211が設けられる。

【0052】

図8は、クラスタシステム200の筐体を説明する図である。図7のノード100の筐体を上下左右に並べ、図7で説明したフレキシブル基板211、212によって左右方向にノード間を接続し、フレキシブル基板213、214によって上下方向にノード間を接続する。このように、クラスタシステム200は、ノードの筐体を平面に配置してフレキシブル基板で接続することで容易に構成することができる。また、ノードの追加がしやすく、スケラビリティがあり、多数のノードを結合したノードクラスタを省スペースで安価に提供することができる。

【0053】

図4~図8では、プロセッサ基板にフレキシブル基板用コネクタが設けられ、フレキシブル基板用コネクタ間をフレキシブル基板で接続する形態を説明した。このように汎用PCIエクスプレスコネクタをフレキシブル基板で接続する形態は、接続形態の一例に過ぎず、これ以外の接続形態も考えられる。別の接続形態として、プロセッサ基板のカードエッジを差し込むための汎用のPCIエクスプレスコネクタを搭載したバックプレーン基板を一枚用意して、4枚のプロセッサ基板をバックプレーン基板に差し込むことで図5で説明したPCIエクスプレスコネクタ間の接続をバックプレーン基板上で実現してもよい。また、さらに別の接続形態として、プロセッサ基板に差動信号用コネクタペアであるZDコネクタを設け、バックプレーン基板上でZDコネクタを接続するように構成してもよい。このようなバックプレーン基板を用いた接続形態もフレキシブル基板を用いた接続形態と同様、安価な高速通信を実現することができ、また、省スペース化を図ることができる

10

20

30

40

50

【 0 0 5 4 】

図 9 A ~ 図 9 D を参照して、図 2 で説明したノード 1 0 0 内の 4 枚のプロセッサ基板のフルメッシュ型の結合により形成される P C I ツリーを説明する。P C I ツリーは、ノード内の各 M C P が P C I エクスプレスで接続された P C I デバイスを検索することにより得られる。

【 0 0 5 5 】

図 9 A は、M C P 0 または M C P 1 を中心に置いた場合の P C I ツリーを説明する図である。同図では P C I ツリー構造において同じ階層にある M C P を水平に配置し、ルートに近い方を上に、リーフに近い方を下に配置している。

10

【 0 0 5 6 】

ルートコンプレックスである M C P 0 のすぐ下の階層には、M C P 3 がエンドポイントとして接続されている。M C P 2 が M C P 3 と同階層にあって、M C P 3 に接続されている。これにより、M C P 0 をルートとする第 1 の P C I ツリーが形成される。M C P 1 は M C P 0 と同階層にあって、M C P 0 に接続されている。ルートコンプレックスである M C P 1 のすぐ下の階層には、M C P 4 がエンドポイントとして接続されている。M C P 5 が M C P 4 と同階層にあって、M C P 4 に接続されている。これにより、M C P 1 をルートとする第 2 の P C I ツリーが形成される。ルートコンプレックスである M C P 6 は、M C P 1 のすぐ上の階層にあって、M C P 1 をエンドポイントとして接続している。M C P 7 が M C P 6 と同階層にあって、M C P 6 に接続されている。これにより、M C P 6 をルートとする第 3 の P C I ツリーが形成される。

20

【 0 0 5 7 】

図 9 B は、M C P 2 または M C P 3 を中心に置いた場合の P C I ツリーを説明する図である。ルートコンプレックスである M C P 2 のすぐ下の階層には、M C P 5 がエンドポイントとして接続されている。M C P 4 が M C P 5 と同階層にあって、M C P 5 に接続されている。これにより、M C P 2 をルートとする第 1 の P C I ツリーが形成される。ルートコンプレックスである M C P 7 は、M C P 2 のすぐ上の階層にあって、M C P 2 をエンドポイントとして接続している。M C P 6 が M C P 7 と同階層にあって、M C P 7 に接続されている。これにより、M C P 7 をルートとする第 2 の P C I ツリーが形成される。M C P 3 は M C P 2 と同階層にあって、M C P 2 に接続されている。ルートコンプレックスである M C P 3 のすぐ下の階層には、他のノードのエンドポイントが接続される。ルートコンプレックスである M C P 0 は、M C P 3 のすぐ上の階層にあって、M C P 3 をエンドポイントとして接続している。M C P 1 が M C P 0 と同階層にあって、M C P 0 に接続されている。これにより、M C P 0 をルートとする第 3 の P C I ツリーが形成される。

30

【 0 0 5 8 】

図 9 C は、M C P 4 または M C P 5 を中心に置いた場合の P C I ツリーを説明する図である。ルートコンプレックスである M C P 4 のすぐ下の階層には、M C P 7 がエンドポイントとして接続されている。M C P 6 が M C P 7 と同階層にあって、M C P 7 に接続されている。これにより、M C P 4 をルートとする第 1 の P C I ツリーが形成される。ルートコンプレックスである M C P 1 は、M C P 4 のすぐ上の階層にあって、M C P 4 をエンドポイントとして接続している。M C P 0 が M C P 1 と同階層にあって、M C P 1 に接続されている。これにより、M C P 1 をルートとする第 2 の P C I ツリーが形成される。M C P 5 は M C P 4 と同階層にあって、M C P 4 に接続されている。ルートコンプレックスである M C P 5 のすぐ下の階層には、他のノードのエンドポイントが接続される。ルートコンプレックスである M C P 2 は、M C P 5 のすぐ上の階層にあって、M C P 5 をエンドポイントとして接続している。M C P 3 が M C P 2 と同階層にあって、M C P 2 に接続されている。これにより、M C P 2 をルートとする第 3 の P C I ツリーが形成される。

40

【 0 0 5 9 】

図 9 D は、M C P 6 または M C P 7 を中心に置いた場合の P C I ツリーを説明する図である。ルートコンプレックスである M C P 6 のすぐ下の階層には、M C P 1 がエンドポイ

50

ントとして接続されている。MCP0がMCP1と同階層にあって、MCP1に接続されている。これにより、MCP6をルートとする第1のPCIツリーが形成される。MCP7はMCP6と同階層にあって、MCP6に接続されている。ルートコンプレックスであるMCP7のすぐ下の階層には、MCP2がエンドポイントとして接続されている。MCP3がMCP2と同階層にあって、MCP2に接続されている。これにより、MCP7をルートとする第2のPCIツリーが形成される。ルートコンプレックスであるMCP4は、MCP7のすぐ上の階層にあって、MCP7をエンドポイントとして接続している。MCP5がMCP4と同階層にあって、MCP4に接続されている。これにより、MCP4をルートとする第3のPCIツリーが形成される。

【0060】

このように、ノード100内の4つのプロセッサ基板間で図2で説明したようにRCコネクタとEPコネクタを接続することにより、あるMCPをルートとするPCIツリーが複数形成される。ノード100内のMCP0～MCP7はそれぞれ、自己をルートとするPCIツリー内で、もしくは異なるPCIツリーをまたぐことで他のMCPとの間でデータ通信を行うことができる。ノード100の空きスロットのコネクタと接続された隣接ノードのプロセッサ基板上的MCPは、同一PCIツリー内にあるため、ノードをまたいでデータ通信が可能である。しかし、ノード100の空きスロットのコネクタと接続されていない他のノードのMCPとデータ通信をする場合は、同一PCIツリー内にないため、ルーティングが必要となる。このため、ノード100内の各MCPは、ソフトウェアでルーティングを実行して、他のPCIツリー内のMCPとの通信を可能にする。

【0061】

ノード100内の各MCPが、自己のPCIツリー内で、もしくは異なるPCIツリーをまたぐことで他のMCPの所定の共有領域にアクセスできるように、各MCPのメモリ空間には他のMCPの所定の共有領域がメモリマッピングされる。図10～図17を参照して、このメモリマッピングを説明する。

【0062】

図10は、ノード100内の各MCPのメモリ空間300を説明する図である。メモリ空間300には、コヒーレントなローカルメモリ領域351とノンコヒーレントな共有メモリ領域352がある。コヒーレントなローカルメモリ領域351は、メモリアクセスのアトミック性が保証され、同期制御がなされる領域であり、他のMCPからはアクセスすることはできない。ノンコヒーレントな共有メモリ領域352は、他のMCPのメモリ空間にマッピングされ、他のMCPからアクセスされる。メモリ空間300には、さらに各MCPのSPEおよびPEのレジスタやSPEのローカルストアがマッピングされたノンコヒーレント領域353がある。このノンコヒーレント領域353の少なくとも一部は、他のMCPのメモリ空間にマッピングされ、他のMCPからアクセスされる。

【0063】

メモリ空間300には、IOIF0を介してアクセス可能なI/Oアドレス空間がIOIF0領域360としてメモリマッピングされる。また、IOIF1を介してアクセス可能なI/Oアドレス空間がIOIF1領域370としてメモリマッピングされる。

【0064】

各MCPは、自分のメモリ空間300内のノンコヒーレント領域353に含まれるSPE/PEのレジスタやSPEのローカルストア、およびノンコヒーレントな共有メモリ領域352を共有領域(shared area)として、IOIF0を介して他のMCPに開放してアクセスを許可する。各MCPは、他のMCPにアクセスを許可する共有領域の情報をIOIF0用のI/Oページテーブル(以下、「IOPT」という)310に格納する。他のMCPは、このIOPT310を参照して、共有領域を自分のメモリ空間にマッピングしてアクセス可能にする。

【0065】

図11は、MCP0の共有領域がMCP1のメモリ空間301にマッピングされ、MCP1の共有領域がMCP0のメモリ空間300にマッピングされる様子を説明する図であ

10

20

30

40

50

る。MCP0のIOIF0用のIOPT310(「IOPT0」)がIOIF0経由でMCP1に提示されると、IOPT0により指定されたMCP0の共有領域321がMCP1のメモリ空間301のIOIF0領域361にマッピングされる。MCP0の共有領域321には、MCP0のSPE/PEのレジスタ、MCP0のSPEのローカルストア、およびMCP0の共有メモリが含まれる。

【0066】

一方、MCP1のIOIF0用のIOPT311(「IOPT1」)がIOIF0経由でMCP0に提示されると、IOPT1により指定されたMCP1の共有領域320がMCP0のメモリ空間300のIOIF0領域360にマッピングされる。MCP1の共有領域320には、MCP1のSPE/PEのレジスタ、MCP1のSPEのローカルストア、およびMCP1の共有メモリが含まれる。

10

【0067】

このように、IOIF0を介して接続されたMCP0とMCP1は、互いに相手の共有領域が自分のメモリ空間300、301にメモリマッピングされているため、相手の共有領域にアクセスすることができる。

【0068】

図12は、IOIF0で相互接続されたMCP0およびMCP1のそれぞれの共有領域がIOIF1経由で接続された他のMCPのメモリ空間にマッピングされる様子を説明する図である。

【0069】

図11で説明したように、MCP1のメモリ空間のIOIF0領域361には、MCP0の共有領域321がマッピングされている。MCP1は、自分の共有領域とともにMCP0の共有領域321をIOIF1経由で接続された他のMCPに開放してアクセスを許可する。MCP1は、IOIF1用のIOPT331に、自分の共有領域、すなわちMCP1のSPE/PEのレジスタ、MCP1のローカルストア、およびMCP1の共有メモリの情報を格納する。さらにMCP1は、IOIF1のIOPT331に、MCP0の共有領域321、すなわちMCP0のSPE/PEのレジスタ、MCP0のローカルストア、およびMCP0の共有メモリの情報を格納する。

20

【0070】

図9Aで説明したように、MCP1とMCP6の接続関係は、MCP6がルートコンプレックス、MCP1がエンドポイントの関係であるから、エンドポイントであるMCP1が自分の共有領域の情報をルートコンプレックスであるMCP6に提示する。MCP1は、IOIF1用のIOPT331をIOIF1経由で接続されたMCP6に提示する。MCP6は、IOIF1用のIOPT331で指定されたMCP0とMCP1の両方の共有領域342を自分のメモリ空間306のIOIF1領域376にマッピングする。

30

【0071】

図13(a)、(b)は、MCP0がIOIF1経由で接続された他のMCPからIOIF1用のIOPTの提示を受けた場合に、MCP0のメモリ空間300に他のMCPの共有領域がマッピングされる様子を説明する図である。

【0072】

図13(a)に示すように、ルートコンプレックスであるMCP0は、IOIF1を経由してエンドポイントであるMCP3に接続されている。MCP3はIOIF0によりMCP2と相互接続されるから、MCP3のメモリ空間のIOIF0領域にはMCP2の共有領域がマッピングされる。図12で説明したMCP1からMCP6へのIOIF1用のIOPTの提示と同様に、エンドポイントであるMCP3は、自分の共有領域とMCP2の共有領域の情報をIOIF1用のIOPTに格納してルートコンプレックスであるMCP0に提示する。

40

【0073】

MCP0は、MCP3からIOIF1用のIOPTの提示を受けて、図13(b)に示すように、メモリ空間300のIOIF1領域370にMCP2およびMCP3の共有領

50

域 340 をマッピングする。

【0074】

図14(a)、(b)は、MCP1がIOIF1経由で接続された他のMCPからIOIF1用のIOPTの提示を受けた場合に、MCP1のメモリ空間301に他のMCPの共有領域がマッピングされる様子を説明する図である。

【0075】

図14(a)に示すように、ルートコンプレックスであるMCP1は、IOIF1を経由してエンドポイントであるMCP4に接続されている。MCP4はIOIF0によりMCP5と相互接続されるから、MCP4のメモリ空間のIOIF0領域にはMCP5の共有領域がマッピングされる。エンドポイントであるMCP4は、自分の共有領域とMCP5の共有領域の情報をIOIF1用のIOPTに格納してルートコンプレックスであるMCP1に提示する。MCP1は、MCP4からIOIF1用のIOPTの提示を受けて、図14(b)に示すように、メモリ空間301のIOIF1領域371にMCP4およびMCP5の共有領域346をマッピングする。

10

【0076】

同様に、MCP6は、自分の共有領域とMCP7の共有領域の情報をIOIF1用のIOPTに格納してMCP1に提示し、MCP1は、MCP6からIOIF1用のIOPTの提示を受けて、図14(b)に示すように、メモリ空間301のIOIF1領域371にMCP6およびMCP7の共有領域348をマッピングする。

【0077】

次に、MCP0とMCP1は、図13(b)、図14(b)のメモリ空間300、301のIOIF1領域370、371にマッピングされた、IOIF1経由で接続された他のMCPの共有領域の情報を互いに交換する。

20

【0078】

図15は、MCP0とMCP1間でメモリマッピングされた共有領域の情報をやりとりする様子を説明する図である。MCP0は、メモリ空間300のIOIF1領域370にマッピングされたMCP2およびMCP3の共有領域340の情報をIOIF0を介してMCP1に与える。MCP1は、MCP0から与えられた情報にもとづき、MCP2およびMCP3の共有領域を自分のメモリ空間301のIOIF0領域361にマッピングする。

30

【0079】

一方、MCP1は、メモリ空間301のIOIF1領域371にマッピングされたMCP4およびMCP5の共有領域346の情報と、MCP6およびMCP7の共有領域348の情報とをIOIF0を介してMCP0に与える。MCP0は、MCP1から与えられた情報にもとづき、MCP4およびMCP5の共有領域とMCP6およびMCP7の共有領域を自分のメモリ空間300のIOIF0領域360にマッピングする。

【0080】

図11、図13(b)、図14(b)、および図15で説明した手順でメモリ空間に他のMCPの共有領域がメモリマッピングされることにより、MCP0は、図9Aで説明した第1～第3PCIツリー内にあるMCP1～MCP7の共有領域にアクセスすることができるようになる。なぜなら第1～第3PCIツリーをまたがって一つのアドレスマップが構築されているからである。同様に、MCP1は、図9Aで説明した第1～第3PCIツリー内にあるMCP0、MCP2～MCP7の共有領域にアクセスすることができるようになる。

40

【0081】

このように、ノード100内の各MCPは、第1～第3PCIツリー内の他のMCPの共有領域を自分のメモリ空間にメモリマッピングしており、第1～第3PCIツリー内の他のMCPの共有領域にアクセスしたり、第1～第3PCIツリー内の他のMCPと共有領域を介したデータ通信や同期制御を実行することができる。ノード100内のプロセッサ基板間はフレキシブル基板で接続され、高速なPCI-Expressによる通信が可

50

能なハードウェア構成が採用されている。したがって、ノード100内の各MCPは、メモリマッピングされた共有領域を高速にアクセスすることができ、他のMCPとデータのやりとりを効率良く行うことができる。

【0082】

図16は、接続ノードのMCPとの接続も含めたPCIツリーを説明する図である。MCP0のブリッジ0のコネクタEP0は、隣接ノードのブリッジ5'のコネクタRC5と接続され、MCP5'がMCP0に対してルートコンプレックスとなる。MCP4'はMCP5'と同階層にあって、MCP5'と接続されている。MCP0は、隣接ノードのMCP5'からIOIF1用IOPTの提示を受けて、MCP5'およびMCP4'の共有領域349をメモリ空間300のIOIF1領域370にマッピングする。

10

【0083】

図17は、図16のPCIツリーの場合におけるMCP0のメモリ空間300を説明する図である。図17に示すように、IOIF1領域370には、MCP2およびMCP3の共有領域340の他、MCP4'およびMCP5'の共有領域340がメモリマッピングされる。また、IOIF0領域360には、MCP4およびMCP5の共有領域326、MCP6およびMCP7の共有領域328、およびMCP1の共有領域320がメモリマッピングされる。

【0084】

まとめると、PCIのメモリマップは、PCIツリーのルートにあるホストプロセッサが、デバイスやスイッチのベースアドレスを設定することで構成される。エンドポイントであるデバイスは、自分が要求するアドレス領域のサイズをホストプロセッサに通知し、ホストプロセッサは、デバイスが要求したサイズにしたがってメモリマップを構築する。具体的には、要求するアドレスレンジのサイズは、コンフィグレーションレジスタのBARフィールドに実装するビット数で指定される。

20

【0085】

本実施の形態のブリッジデバイスは、エンドポイントとして動作する場合、外部からアクセス可能なコンフィグレーションレジスタと内部からアクセス可能なコンフィグレーションレジスタをそれぞれ別々のレジスタとして実装し、それぞれのレジスタについて要求されるアドレスレンジのサイズ、すなわちBARの実装ビット数を設定することが可能である。これにより、システム初期化時に設定したサイズのアドレスレンジにより、PCIのアドレスマップがそれぞれのホストプロセッサによって構築される。ここで、それぞれのホストプロセッサとは、ルートコンプレックスとなるプロセッサと、エンドポイントとして動作するプロセッサのことである。

30

【0086】

一方、IOIFのメモリマップは、IOPTに共有領域の情報を格納して他のMCPに提示することにより設定される。この作業は外部からメモリアクセスがあった場合、その先にマッピングされる領域を設定するものである。この設定作業は、自分がルートコンプレックスとして動作する場合でも、自分がエンドポイントとして動作する場合でも、PCIからトランザクションを受け、それをメモリアクセスとして許可する場合は必要となる。

40

【0087】

実際の運用としては、PCIで構築するアドレスサイズは余裕をもたせたサイズにしてPCIメモリマップを構築し、その中で実際にメモリをマップする範囲は、IOPTによって設定することになる。また、PCIメモリマップのアドレスレンジに関しては、PCIエクスプレスの規格にしたがい、エンドポイントが通知し、ルートコンプレックスがアドレス構築するということになるが、その中で、どの範囲がメモリやローカルストアにマッピングされているかについての情報は、図15で説明したように、共有メモリを介したオリジナルプロトコルでやりとりする必要がある。

【0088】

以上説明したように、本実施の形態によれば、プロセッサ基板の汎用のPCIエクスブ

50

レスコネクタ間を安価なフレキシブル基板やバックプレーン基板で直接接続することにより、P C I エクスプレススイッチを必要としない、安価でかつ高性能なクラスタシステムを構築することができる。

【 0 0 8 9 】

以上、本発明を実施の形態をもとに説明した。実施の形態は例示であり、それらの各構成要素や各処理プロセスの組合せにいろいろな変形例が可能なこと、またそうした変形例も本発明の範囲にあることは当業者に理解されるところである。

【 0 0 9 0 】

上記の実施の形態では、プロセッサ基板にマルチコアプロセッサが搭載された場合を説明したが、これはシングルプロセッサであってもよい。また、実施の形態では、プロセッサ基板に2つのマルチコアプロセッサが搭載され、4枚のプロセッサ基板で1つのノードを構成する例を説明したが、プロセッサ基板に搭載されるプロセッサの数、1つのノード内のプロセッサ基板の数、ブリッジのコネクタ数などは、設計の自由度がある。ノード内の複数のプロセッサ基板をフレキシブル基板によって密結合し、ノード間をさらにフレキシブル基板で連結してノードクラスタを構成することができる限り、ノード内のプロセッサ基板の数と配置、ノードクラスタ内のノードの配置にはいろいろなパターンがありうる。いずれにしても安価、省スペース、高速通信の各要求を満足する接続形態が好ましい。

【 0 0 9 1 】

上記の実施の形態では、プロセッサ基板のブリッジのポートに他のプロセッサ基板のポートを接続したが、プロセッサ基板のブリッジのポートに周辺デバイスを接続し、プロセッサと各種周辺デバイスを相互結合したシステムを構成してもよい。また、ブリッジはプロセッサの入出力バスをP C I エクスプレスに接続したが、他のプロセッサ基板や周辺デバイスが接続される外部インタフェースとしてP C I エクスプレス以外のインタフェースが用いられてもよい。

【 図面の簡単な説明 】

【 0 0 9 2 】

【 図 1 】 プロセッサ基板の構成図である。

【 図 2 】 4つのプロセッサ基板が密結合されたノードの構成図である。

【 図 3 】 複数のノードを連結したクラスタシステムの構成図である。

【 図 4 】 プロセッサ基板の裏面の配線の模式図である。

【 図 5 】 ノード内の4枚のプロセッサ基板間をフレキシブル基板によって接続した構成を示す図である。

【 図 6 】 クラスタシステム内の複数のノード間をフレキシブル基板によって接続した構成を示す図である。

【 図 7 】 ノードの筐体を説明する図である。

【 図 8 】 クラスタシステムの筐体を説明する図である。

【 図 9 A 】 図 2 のノード内で形成されるP C I ツリーを説明する図である。

【 図 9 B 】 図 2 のノード内で形成されるP C I ツリーを説明する図である。

【 図 9 C 】 図 2 のノード内で形成されるP C I ツリーを説明する図である。

【 図 9 D 】 図 2 のノード内で形成されるP C I ツリーを説明する図である。

【 図 1 0 】 ノード内の各M C Pのメモリ空間を説明する図である。

【 図 1 1 】 あるM C Pのメモリ空間に他のM C Pの共有領域がマッピングされる様子を説明する図である。

【 図 1 2 】 あるM C Pのメモリ空間に他のM C Pの共有領域がマッピングされる様子を説明する図である。

【 図 1 3 】 あるM C Pのメモリ空間に他のM C Pの共有領域がマッピングされる様子を説明する図である。

【 図 1 4 】 あるM C Pのメモリ空間に他のM C Pの共有領域がマッピングされる様子を説明する図である。

【 図 1 5 】 2つのM C P間でメモリマッピングされた共有領域の情報をやりとりする様子

10

20

30

40

50

を説明する図である。

【図16】 接続ノードのMCPとの接続も含めたPCIツリーを説明する図である。

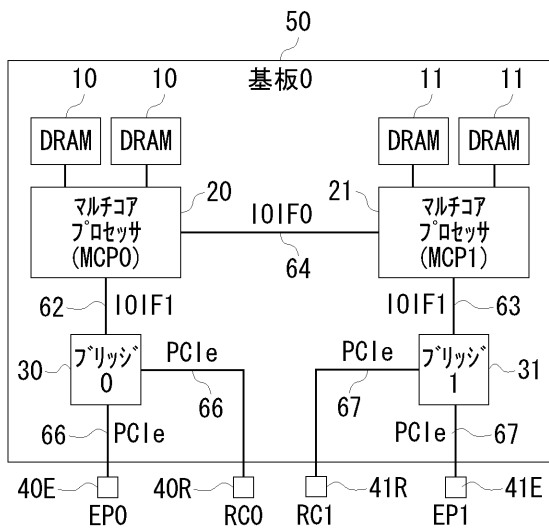
【図17】 図16のPCIツリーの場合におけるMCPのメモリ空間を説明する図である。

【符号の説明】

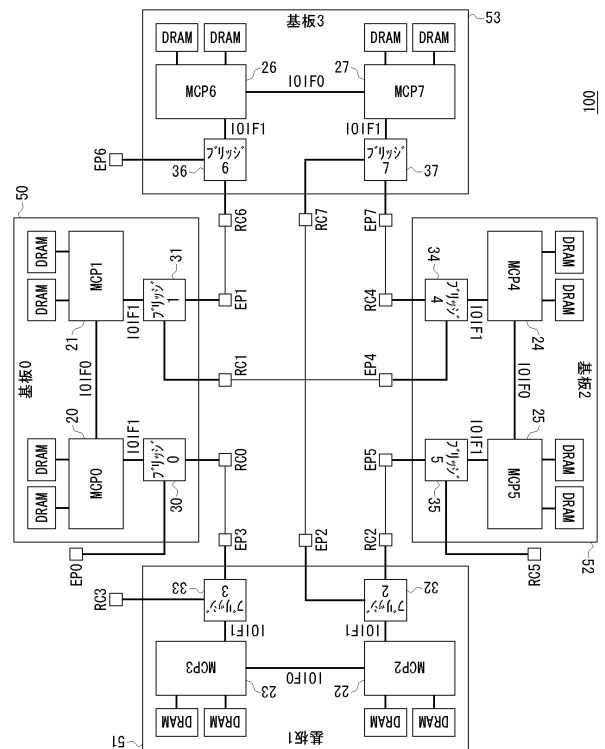
【0093】

10 DRAM、 20 マルチコアプロセッサ、 30 ブリッジ、 50 プロセッサ基板、 100 ノード、 200 クラスタシステム、 201~206、 211~214 フレキシブル基板、 300 メモリ空間、 310 IOPT。

【図1】

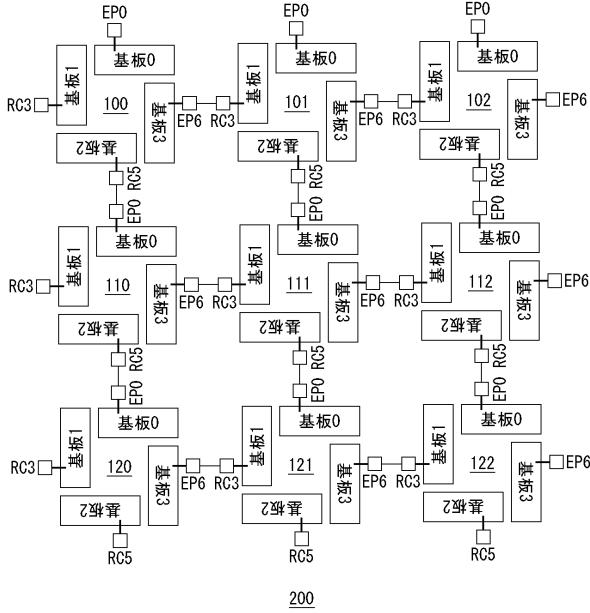


【図2】

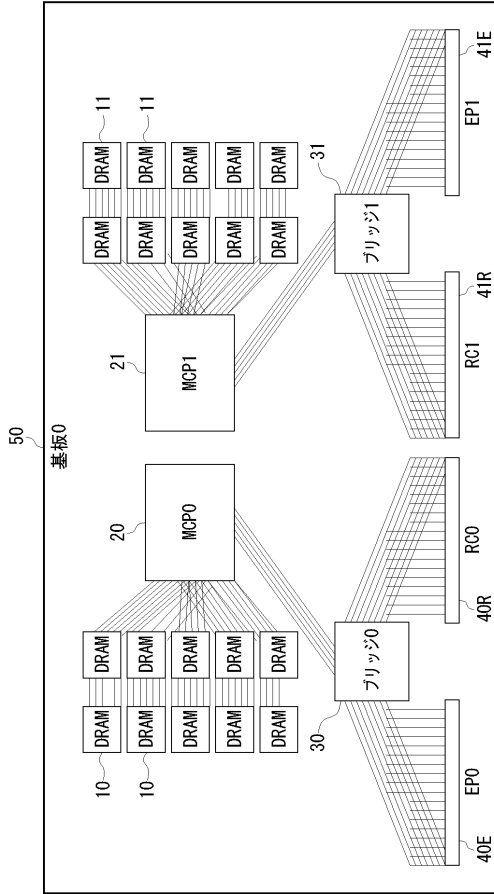


001

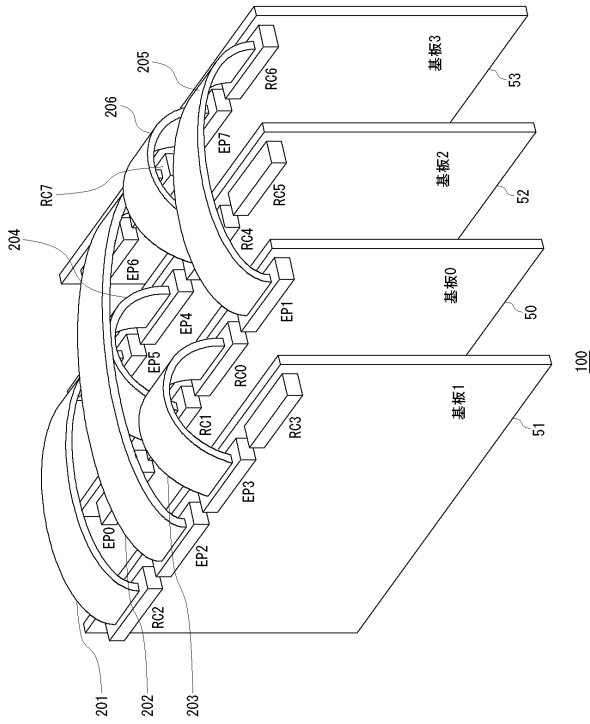
【図3】



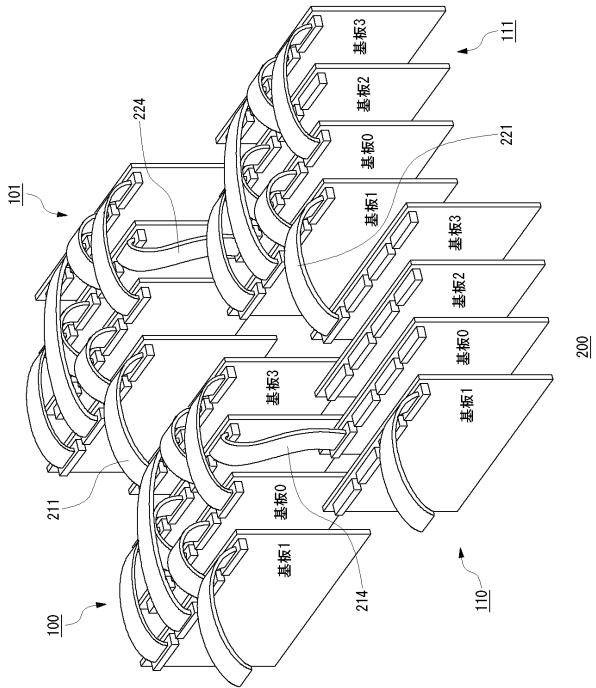
【図4】



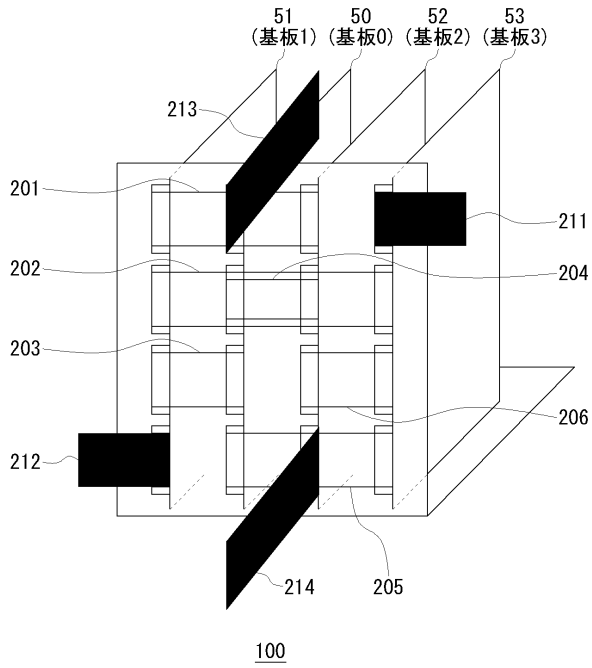
【図5】



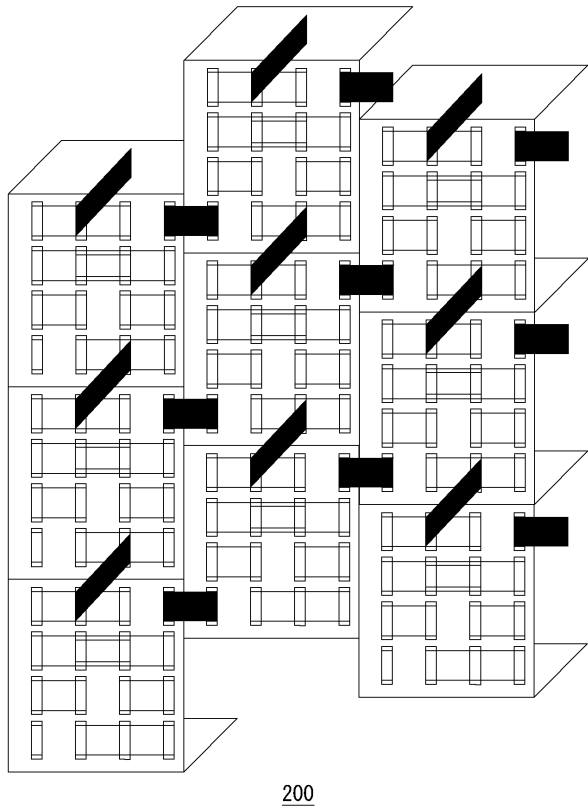
【図6】



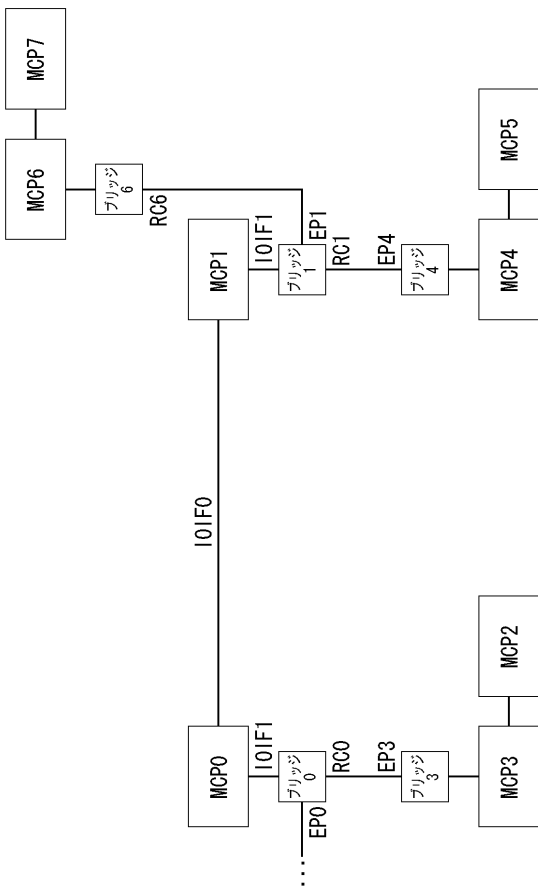
【図7】



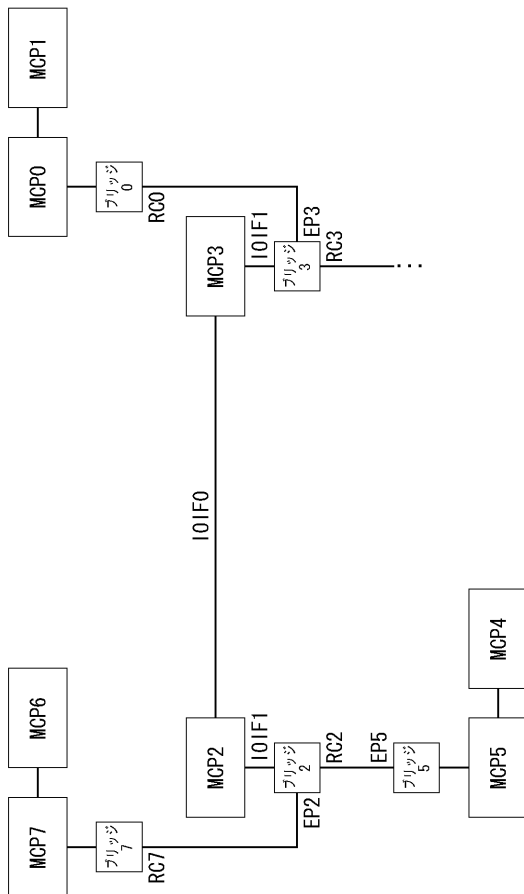
【図8】



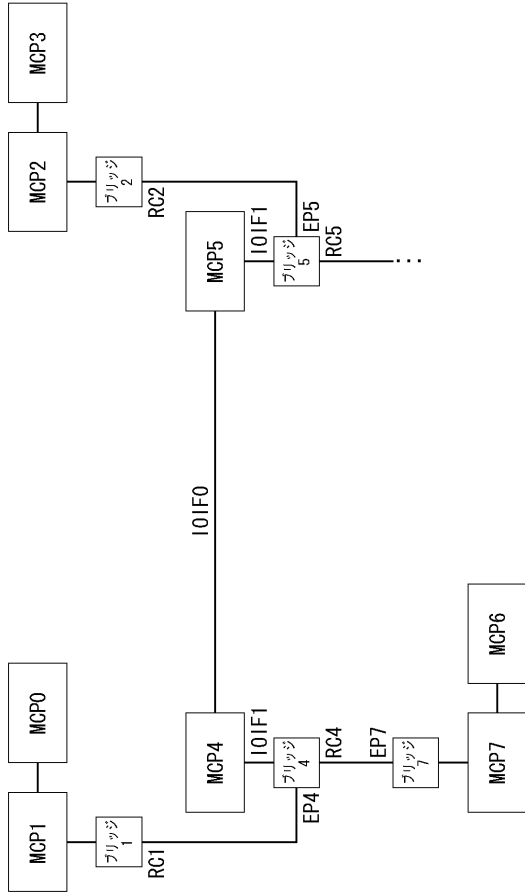
【図9A】



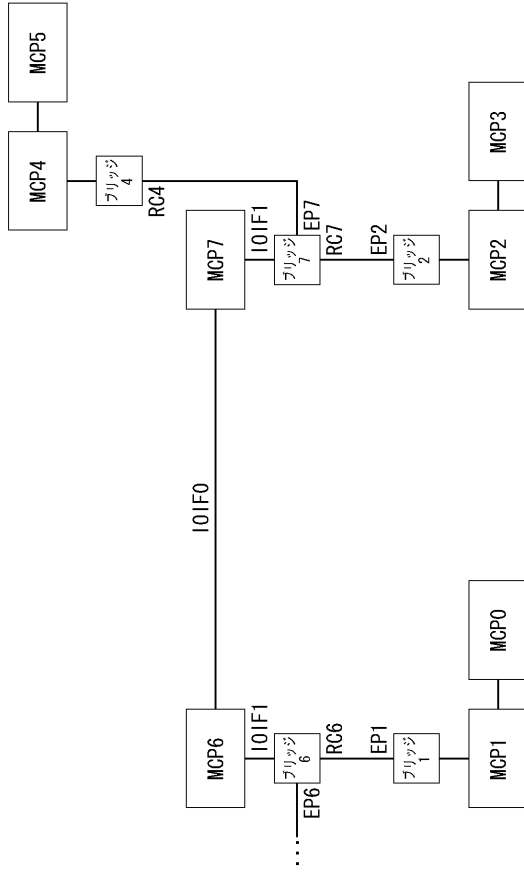
【図9B】



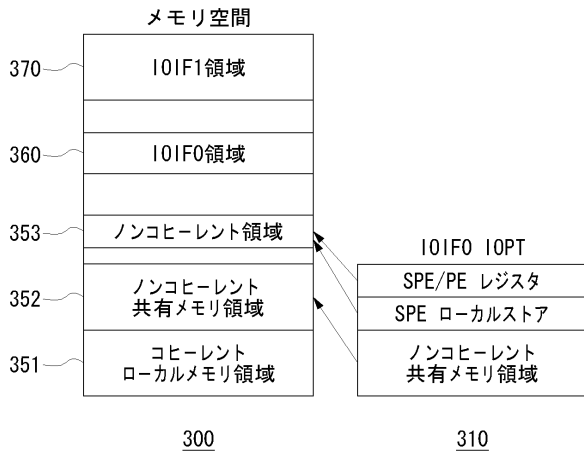
【図 9 C】



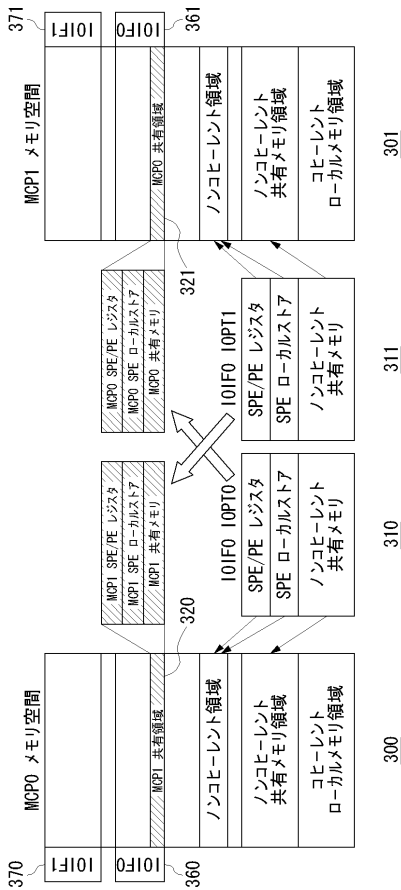
【図 9 D】



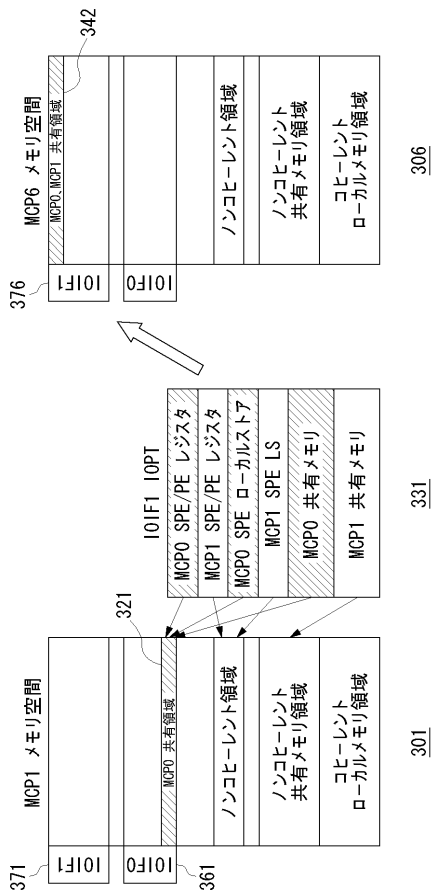
【図 1 0】



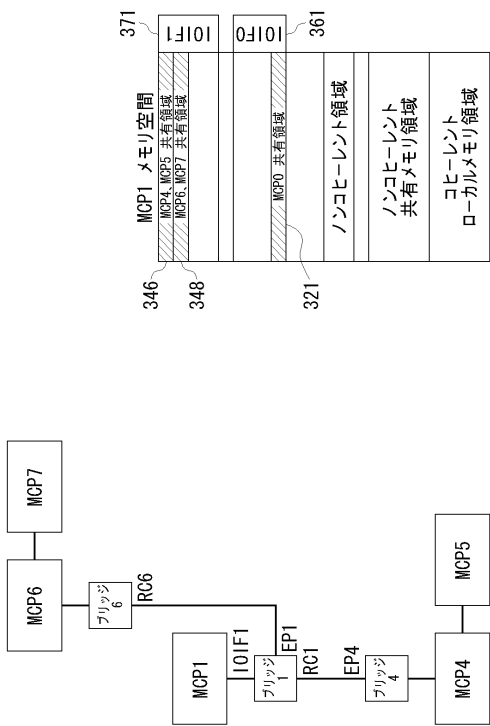
【図 1 1】



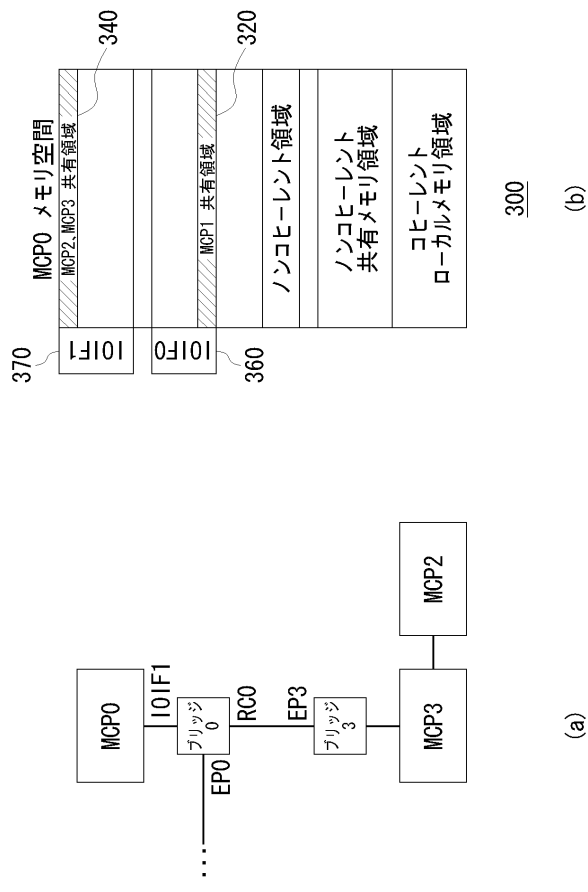
【図 1 2】



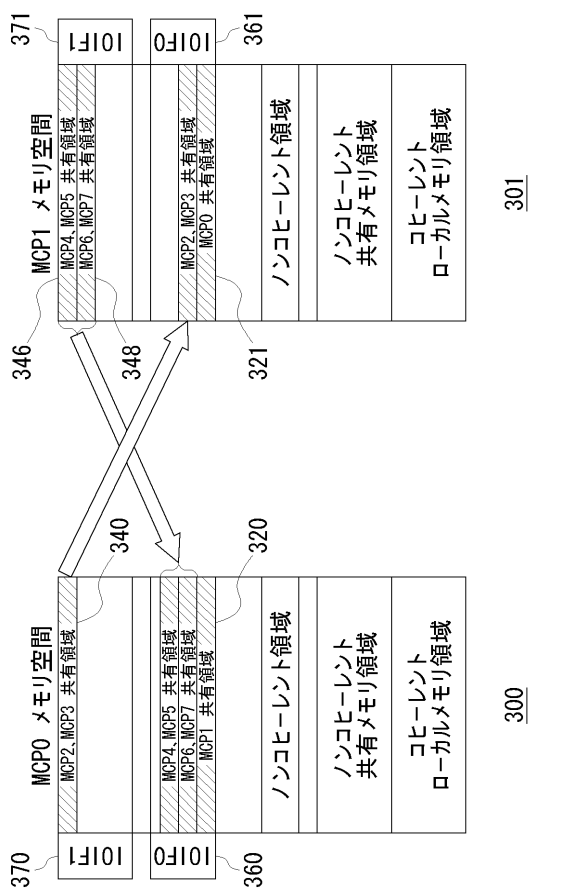
【図 1 4】



【図 1 3】



【図 1 5】



300

(b)

(a)

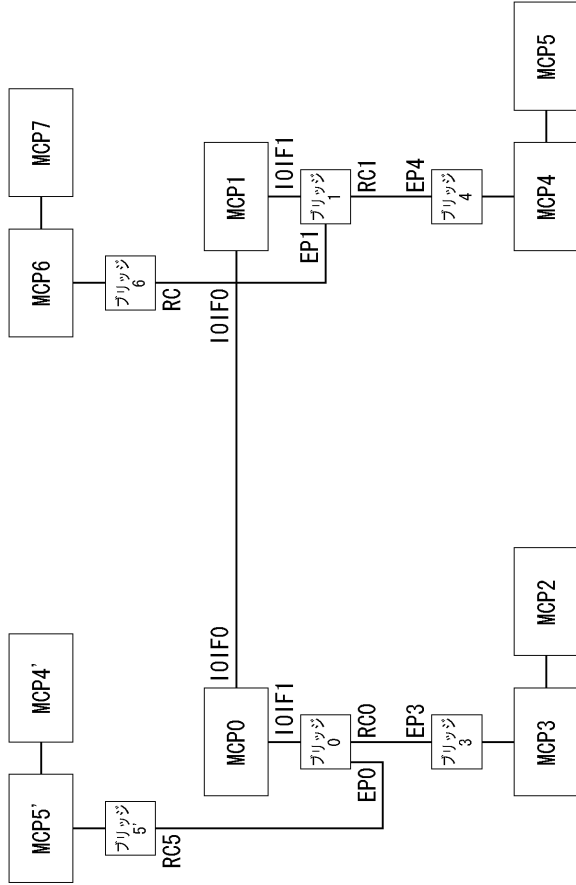
300

(b)

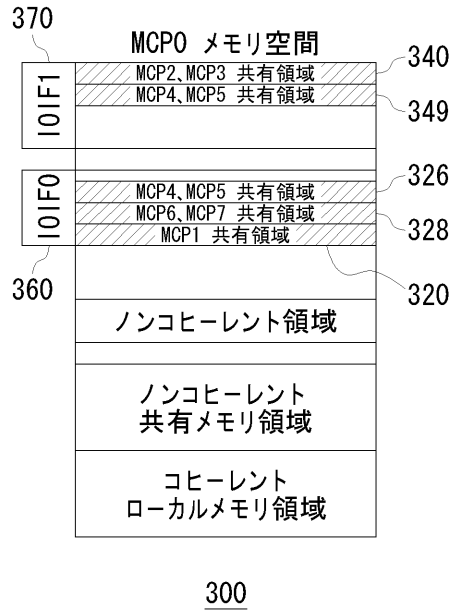
300

(a)

【図16】



【図17】



300

フロントページの続き

(72)発明者 斎藤 英幸

東京都港区南青山2丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内

(72)発明者 堀江 和由

東京都港区南青山2丁目6番21号 株式会社ソニー・コンピュータエンタテインメント内

審査官 三坂 敏夫

(56)参考文献 特開平07-093236(JP,A)

特表2000-506645(JP,A)

特開2006-209456(JP,A)

特開2004-070954(JP,A)

特開平09-146895(JP,A)

特開平08-044460(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 15/173