

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5573829号  
(P5573829)

(45) 発行日 平成26年8月20日 (2014. 8. 20)

(24) 登録日 平成26年7月11日 (2014. 7. 11)

(51) Int. Cl. F I  
**G 0 6 F 12/06 (2006.01)** G O 6 F 12/06 5 3 0 A  
**G 0 6 F 12/08 (2006.01)** G O 6 F 12/08 5 3 1 E

請求項の数 8 (全 37 頁)

(21) 出願番号	特願2011-279022 (P2011-279022)	(73) 特許権者	000005223
(22) 出願日	平成23年12月20日 (2011. 12. 20)		富士通株式会社
(65) 公開番号	特開2013-130976 (P2013-130976A)		神奈川県川崎市中原区上小田中4丁目1番1号
(43) 公開日	平成25年7月4日 (2013. 7. 4)	(74) 代理人	100089118
審査請求日	平成24年7月20日 (2012. 7. 20)		弁理士 酒井 宏明
前置審査		(72) 発明者	鯉沼 秀之
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	岡田 誠之
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内
		(72) 発明者	杉崎 剛
			神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

最終頁に続く

(54) 【発明の名称】 情報処理装置およびメモリアクセス方法

(57) 【特許請求の範囲】

【請求項1】

それぞれがプロセッサ及び記憶装置を備える複数のノードと、前記複数のノード間を接続するインターコネクットとを有する情報処理装置であって、

前記ノードの各々は、

各ノードに備えられたプロセッサを識別するプロセッサ識別情報と、当該プロセッサ識別情報が示すプロセッサを備えるノードの記憶装置に割当てられた物理アドレスとを対応付けて記憶する記憶部と、

論理アドレスと物理アドレスとの変換を行う第1変換部と、

物理アドレスを、前記記憶部に当該物理アドレスと対応付けて記憶されたプロセッサ識別情報に変換する第2変換部と、

前記物理アドレス及び前記プロセッサ識別情報を含み、当該プロセッサ識別情報が識別するプロセッサを備えるノードの記憶装置からデータの読出しを指示する読出要求、または、前記物理アドレス及び前記プロセッサ識別情報を含み、当該プロセッサ識別情報が識別するプロセッサを備えるノードの記憶装置へのデータの書込みを指示する書込要求を送信する送信部と、

他のノードから前記インターコネクットを介して送信された読出要求または書込要求を受信する受信部と、

前記受信部により受信された前記読出要求または前記書込要求に含まれる物理アドレスに基づいて、自ノードの記憶装置のデータ格納領域のうち、自ノード内のアクセスに用い

10

20

られるローカル領域と、他のノードからもアクセス可能な共有領域とのいずれの領域へのアクセスであるかを判定するローカル判定部と、

各ノードへの前記プロセッサまたは前記記憶装置の追加や削除に応じて、前記記憶部に記憶されたプロセッサ識別情報と物理アドレスとの対応付けを書換える制御部とを有することを特徴とする情報処理装置。

【請求項 2】

前記ノードの各々は、

前記ノードの各々が備える記憶装置の物理アドレスのうち、所定の位置のビットが同一の値である物理アドレスを前記共有領域に割当てるとともに、前記所定の位置のビットが前記共有領域に割当てた物理アドレスとは異なる値である物理アドレスを前記ローカル領域に割当し、

前記ローカル判定部は、前記読出要求または前記書込要求に含まれる物理アドレスのうち、前記所定の位置のビットの値に応じて、前記ローカル領域と前記共有領域とのいずれの領域へのアクセスであるかを判定する

ことを特徴とする請求項 1 に記載の情報処理装置。

【請求項 3】

前記ノードの各々は、

前記ノードの各々が備える記憶装置の全物理アドレスを前記ローカル領域と前記共有領域とに分けて割り当て、

前記ローカル判定部は、前記読出要求または前記書込要求に含まれる物理アドレスのうち、最上位のビットの値に応じて、前記ローカル領域と前記共有領域とのいずれの領域へのアクセスであるかを判定する

ことを特徴とする請求項 1 または 2 に記載の情報処理装置。

【請求項 4】

前記ノードの各々は、

前記ローカル判定部が前記ローカル領域へのアクセスであると判定した場合には、前記読出要求または前記書込要求の送信元となるノードに対してアクセスを許可しない旨の否定応答を送信することを特徴とする請求項 1 - 3 のいずれか 1 つに記載の情報処理装置。

【請求項 5】

前記ノードの各々は、前記他のノードから前記読出要求または前記書込要求を受信した場合には、当該読出要求または当該書込要求の対象となるデータにつき、自ノードが備える記憶装置のデータをキャッシュした他のノードを示すディレクトリを用いて、自ノードが備える記憶装置上のデータと他のノードがキャッシュしたデータとの同一性を保持するディレクトリ制御部を有することを特徴とする請求項 1 - 4 のいずれか 1 つに記載の情報処理装置。

【請求項 6】

前記ノードの各々は、

前記記憶装置からデータをキャッシュするキャッシュメモリと、

キャッシュミスが発生した場合は、キャッシュミスした物理アドレスが他のノードが有する記憶装置の物理アドレスであるか否かを判別する判別部とをさらに有し、

前記第 2 変換部は、前記キャッシュミスした物理アドレスが他のノードが有する記憶装置の物理アドレスであると前記ローカル判定部が判定した場合は、当該物理アドレスをプロセッサ識別情報に変換することを特徴とする請求項 1 - 5 のいずれか 1 つに記載の情報処理装置。

【請求項 7】

前記ノードの各々が備えるプロセッサにより実行される各 OS は、アプリケーションから前記共有領域の獲得が要求された場合には、当該アプリケーションが使用する論理アドレスと、前記共有領域に割当てられる物理アドレスとの変換を行うように前記第 1 変換部を設定することを特徴とする請求項 1 - 6 のいずれか 1 つに記載の情報処理装置。

【請求項 8】

10

20

30

40

50

それぞれがプロセッサと、記憶装置と、各ノードに備えられたプロセッサを識別するプロセッサ識別情報と当該プロセッサ識別情報が示すプロセッサを備えるノードの記憶装置に割当てられた物理アドレスとを対応付けて記憶する記憶部とを備え、インターコネクトを介して他ノードと接続されるノードが実行するメモリアクセス方法であって、

アクセス対象の論理アドレスと物理アドレスとの変換を行い、

前記物理アドレスを前記記憶部に当該物理アドレスと対応付けて記憶されたプロセッサ識別情報に変換し、

前記物理アドレス及び前記プロセッサ識別情報を含み、当該プロセッサ識別情報が識別するプロセッサを備えるノードの記憶装置からデータの読出しを指示する読出要求、または、前記物理アドレス及び前記プロセッサ識別情報を含み、当該プロセッサ識別情報が識別するプロセッサを備えるノードの記憶装置へのデータの書込みを指示する書込要求を送信し、

10

他のノードから前記インターコネクトを介して送信された読出要求または書込要求を受信した場合は、当該読出要求または書込要求に含まれる物理アドレスに基づいて、自ノードの記憶装置のデータ格納領域のうち、自ノード内のアクセスに用いられるローカル領域と、他ノードからもアクセス可能な共有領域とのいずれの領域へのアクセスであるかを判定し、

各ノードへの前記プロセッサまたは前記記憶装置の追加や削除に応じて、前記記憶部に記憶されたプロセッサ識別情報と物理アドレスとの対応付けを書換える

処理を実行することを特徴とするメモリアクセス方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置およびメモリアクセス方法に関する。

【背景技術】

【0002】

従来、複数の演算処理装置が主記憶装置を共有するSMP(Symmetric Multiprocessor)の技術が知られている。このようなSMPの技術が適用された情報処理システムの一例として、演算処理装置と主記憶装置とを有する複数のノードを同一のバスで接続し、バスを介して、各演算処理装置が各主記憶装置を共有する情報処理システムが知られている。

30

【0003】

このような情報処理システムでは、例えばスヌープ方式を用いて、各ノードの演算処理装置がキャッシュしたデータの coherence を保持する。しかし、スヌープ方式では、各演算処理装置がキャッシュしたデータの更新状況をバスを介してやり取りするので、ノード数が増加するに従って、バスがボトルネックとなり、メモリアクセスの性能が悪化する。

【0004】

このようなバスのボトルネックを回避するため、インターコネクトを用いて、複数のノードを接続し、各ノードの演算処理装置が各ノードの主記憶装置を共有するNUMA(Non Uniform Memory Access)の技術が知られている。

40

【0005】

このようなNUMAの技術が適用された情報処理システムでは、各ノードの主記憶装置の記憶領域が共通の物理アドレス空間に対して一意にマップされる。そして、各ノードの演算装置は、アクセス対象の物理アドレスが示す記憶領域が存在するノードを識別し、インターコネクトを介して、識別したノードの主記憶装置にアクセスする。

【先行技術文献】

【特許文献】

【0006】

【特許文献1】特開2000-235558号公報

50

【非特許文献】

【0007】

【非特許文献1】Computer Architecture: A Quantitative Approach, Second Edition, John L. Hennessy, David A. Patterson, §8.4

【発明の概要】

【発明が解決しようとする課題】

【0008】

ここで、上述したNUMAの技術では、各ノードの演算処理装置がキャッシュしたデータの coherence を保つことができない。そこで、各ノードの演算処理装置がキャッシュしたデータの coherence を保つ機構を備える ccNUMA (Cache Coherent NUMA) を採用することが考えられる。

10

【0009】

しかしながら、ccNUMAを適用した情報処理システムでは、各ノードがアクセス対象となる記憶領域が存在するノードを識別するので、アドレス変換を効率良く行う必要がある。また、各ノードは、主記憶装置を自ノードのみが利用する記憶領域と、他ノードと共用する記憶領域とに分割する場合がある。このような場合には、各ノードは、アクセス対象記憶領域が、他のノードと共用する記憶領域であるかを効率良く行う必要がある。

【0010】

本発明は、1つの側面では、各演算処理装置によるメモリアクセスを効率的に行うことができる。

20

【課題を解決するための手段】

【0011】

1つの側面では、それぞれがプロセッサ及び記憶装置を備える複数のノードと、各ノード間を接続するインターコネクトとを有する情報処理装置である。このような情報処理装置が有する各ノードは、各ノードに備えられたプロセッサを識別するプロセッサ識別情報と、当該プロセッサ識別情報が示すプロセッサを備えるノードの記憶装置に割当てられた物理アドレスとを対応付けて記憶する記憶部を有する。また、各ノードは、論理アドレスと物理アドレスとの変換を行う。また、各ノードは、物理アドレスを、前記記憶部に当該物理アドレスと対応付けて記憶されたプロセッサ識別情報に変換する。また、各ノードは、第1変換部によって論理アドレスから変換された物理アドレスと、前記第2変換部により物理アドレスから変換されたプロセッサ識別情報を含み、当該プロセッサ識別情報が識別するプロセッサを備えるノードの記憶装置からデータの読出しを指示する読出要求、または、前記物理アドレス及び前記プロセッサ識別情報を含み、当該プロセッサ識別情報が識別するプロセッサを備えるノードの記憶装置へのデータの書込みを指示する書込要求を送信する。また、各ノードは、他のノードからインターコネクトを介して送信された読出要求または書込要求を受信すると、読出要求または書込要求に含まれる物理アドレスに基づいて、自ノード内のアクセスに用いられるローカル領域と、他ノードからもアクセス可能な共有領域とのいずれの領域へのアクセスであるかを判定する。また、前記ノードは、各ノードへの前記プロセッサまたは前記記憶装置の追加や削除に応じて、前記記憶部に記憶されたプロセッサ識別情報と物理アドレスとの対応付けを書換える制御部を有する。

30

40

【発明の効果】

【0012】

1実施形態によれば、各演算処理装置によるメモリアクセスを効率的に行うことができる。

【図面の簡単な説明】

【0013】

【図1】図1は、実施例1に係る情報処理システムの一例を説明するための図である。

【図2】図2は、実施例1に係るビルディングブロックの機能構成を説明するための図である。

【図3】図3は、実施例1に係るビルディングブロックのメモリに振り分けられる物理ア

50

ドレスの範囲を説明するための図である。

【図 4】図 4 は、実施例 1 に係る情報処理システムが各メモリに振り分ける物理アドレスを説明するための図である。

【図 5】図 5 は、物理アドレスの振り分けのバリエーションを説明するための第 1 の図である。

【図 6】図 6 は、物理アドレスの振り分けのバリエーションを説明するための第 2 の図である。

【図 7】図 7 は、実施例 1 に係る CPU の機能構成を説明するための図である。

【図 8】図 8 は、実施例 1 に係るノードマップが記憶する情報の一例を説明するための図である。

10

【図 9】図 9 は、ノードマップが記憶する情報のバリエーションの一例を説明するための第 1 の図である。

【図 10】図 10 は、ノードマップが記憶する情報のバリエーションの一例を説明するための第 2 の図である。

【図 11 A】図 11 A は、キャッシュタグの一例を説明するための図である。

【図 11 B】図 11 B は、実施例 1 に係る CPU が送信するパケットを説明するための図である。

【図 12】図 12 は、実施例 1 に係る CPU がリクエストを送信する処理の一例を説明するための図である。

【図 13】図 13 は、実施例 1 に係る CPU がパケットを受信した際に実行する処理の一例を説明するための図である。

20

【図 14】図 14 は、ノードマップを設定する処理の流れを説明するためのフローチャートである。

【図 15】図 15 は、共有領域を制御する処理の流れを説明するためのフローチャートである。

【図 16】図 16 は、共有メモリの割当処理を説明するためのフローチャートである。

【図 17】図 17 は、共有メモリアタッチ処理を説明するためのフローチャートである。

【図 18】図 18 は、アプリケーションが共有メモリを使用する処理を説明するためのフローチャートである。

【図 19】図 19 は、ノード間の共有メモリデタッチ処理を説明するためのフローチャートである。

30

【図 20】図 20 は、ノード間共有メモリの解放処理を説明するためのフローチャートである。

【図 21】図 21 は、リクエストを発行する処理の流れを説明するためのフローチャートである。

【図 22】図 22 は、リクエストを受信した際に実行する処理の流れを説明するためのフローチャートである。

【図 23】図 23 は、CPU が応答を受信した際に実行する処理の流れを説明するためのフローチャートである。

【図 24】図 24 は、実施例 2 に係る情報処理システムを説明するための図である。

40

【図 25】図 25 は、パーティションの一例を説明するための図である。

【図 26 A】図 26 A は、パーティション # A の CPU が記憶するノードマップの一例を説明するための図である。

【図 26 B】図 26 B は、パーティション # A を示すノードマップの一例を説明するための図である。

【図 26 C】図 26 C は、パーティション # B を示すノードマップの一例を説明するための図である。

【発明を実施するための形態】

【0014】

以下に添付図面を参照して本願に係る情報処理装置及びメモリアクセス方法について説

50

明する。

【実施例 1】

【0015】

まず、本願に係る実施例の説明の前に、従来の情報処理システムが有する問題の具体例について説明する。例えば、従来の情報処理システムは、CPU (Central Processing Unit) が共有メモリ領域にアクセスするために出力した論理アドレスを共有メモリ空間アドレスに変換する。そして、情報処理システムは、共有メモリ空間アドレスを物理アドレスに変換することで、CPUのアクセス対象となる記憶領域を識別する。

【0016】

しかし、このように論理アドレスを共有メモリ空間アドレスに変換し、変換後の共有メモリ空間アドレスを物理アドレスに変換する手法では、アドレス変換に要するハードウェアの物量が多くなってしまふ。また、論理アドレスを共有メモリ空間アドレスに変換し、変換後の共有メモリ空間アドレスを物理アドレスに変換する手法では、アドレス変換に要する時間が増加してしまふ。

【0017】

また、従来の情報処理システムは、CPUが共有メモリ空間のデータをキャッシュする際に、全てのCPUに対してキャッシュ情報を送信することで、コヒーレンスを保持する。しかし、このように、全てのCPUに対してキャッシュ情報を送信する手法では、ボトルネックが発生し、メモリアクセスの性能が悪化してしまふ。また、従来の情報処理システムは、CPUの増設を行った場合は、CPUの個数の増加に比例してバストラフィックが増加するので、ボトルネックが発生し、メモリアクセスの性能が悪化してしまふ。

【0018】

また、例えば、ノードは、自ノードのみがアクセスするローカル領域にカーネルデータやユーザデータを格納する。このため、各ノードは、ローカル領域に格納したデータのセキュリティを確保し、ソフトウェアバグに対する耐性を高くするため、アクセス対象となる記憶領域が他ノードからアクセス可能である共有メモリ域かアクセス不能なローカルメモリ域かの判断を要する。

【0019】

このため、従来の情報処理システムにおいては、ローカル領域に記憶されたデータはキャッシュ可能とし、共有領域に記憶されたデータはキャッシュ不能とする。しかし、このように共有領域に記憶されたデータをキャッシュ不能とする手法では、メモリアクセスにおけるレイテンシが増大してしまふ。なお、他のノードからメモリアクセスが行われる度にアクセス対象が共有領域であるかローカル領域であるかを判断する場合には、判断を行うための回路規模が増大するとともに、アクセスにおけるレイテンシが増大してしまふ。

【0020】

また、従来の情報処理システムにおいては、ノードは、他のノードが有するメモリにアクセスする度に、特殊なチャンネル装置やDMA (Direct Memory Access) エンジンのプログラムの実行を要するため、メモリアクセスの性能が劣化してしまふ。また、従来の情報処理システムでは、メモリが有する記憶領域のどの領域を共有領域とするかを固定的に設定する。このため、例えば、従来の情報処理システムでは、システムを停止させることなく、ノードを追加することで共有領域を追加することができない。

【0021】

また、従来の情報処理システムにおいては、チャンネルやDMA経路を介したメモリアクセスを行うためのハードウェアが追加される。このため、従来の情報処理システムにおいては、ノード間でメモリを共有しないシステムと比較して、設置されるハードウェアが大きく異なる。この結果、従来の情報処理システムにおいては、ノード間でメモリを共有する場合には、OS (Operating System) 等のプログラムを大幅に変更しなければならない。

【0022】

以下の説明では、実施例 1 として、上述した問題を解決する情報処理システムの一例に

10

20

30

40

50

について説明する。まず、図1を用いて、情報処理システムの構成例について説明する。図1は、実施例1に係る情報処理システムの一例を説明するための図である。図1に示す例では、情報処理システム1は、XB(クロスバススイッチ)2と複数のビルディングブロック10~10eとを有する。また、各ビルディングブロック10~10eは、管理用ネットワークを介して管理端末3と接続されている。また、XB2は、サービスプロセッサ2bを有する。

【0023】

ビルディングブロック10は、複数のCPU21~21cと複数のメモリ22~22cとサービスプロセッサ24とを有する。また、他のビルディングブロック10~10eも、ビルディングブロック10と同様の構成を有するものとし、以下の説明を省略する。なお、図1に示す例では、CPU21b、21bおよびメモリ22b、22cについては、記載を省略した。

10

【0024】

XB2は、各ビルディングブロック10~10eを相互に接続するクロスバススイッチである。また、XB2が有するサービスプロセッサ2bは、各ビルディングブロック10~10eが有するサービスプロセッサを管理するサービスプロセッサ、すなわち、マスタとなるサービスプロセッサである。また、管理端末3は、管理用ネットワークを介して、各ビルディングブロック10~10eが有するサービスプロセッサの設定や制御を行う端末である。なお、少数のノードが接続される小規模構成の場合、XB2を介さずにビルディングブロック同士を直接接続してもよい。

20

【0025】

各ビルディングブロック10~10eは、それぞれ独立してOSを動作させる。すなわち、各ビルディングブロック10~10eが実行するOSは、ビルディングブロック毎に異なるパーティションで動作する。ここで、パーティションとは、同一のOSが動作し、動作しているOSから見て1つのシステムとして動作するビルディングブロックの群を示す。

【0026】

例えば、ビルディングブロック10~10aがパーティション#Aとして動作し、ビルディングブロック10b~10dがパーティション#Bとして動作する。このような場合には、ビルディングブロック10が動作させるOSは、ビルディングブロック10、10aが1つのシステムとして動作していると識別し、ビルディングブロック10bが動作させるOSは、ビルディングブロック10b~10dが1つのシステムとして動作していると識別する。

30

【0027】

次に、図2を用いて、ビルディングブロックの構成例について説明する。図2は、実施例1に係るビルディングブロックの機能構成を説明するための図である。図2に示す例では、ビルディングブロック10は、ノード20、サービスプロセッサ24、XB接続部27、27a、PCIE(Peripheral Component Interconnect Express)接続部28を有する。

【0028】

また、ノード20は、複数のCPU21~21cと複数のメモリ22~22cと通信部23とを有する。また、サービスプロセッサ24は、制御部25と通信部26とを有する。また、図2に示す例では、各CPU21~21cは、相互に直接接続されるとともに、通信部23と接続されている。また、各メモリ22~22cは、各CPU21~21cと接続されている。

40

【0029】

また、各CPU21~21cは、XB接続部27またはXB接続部27aと接続されている。なお、XB接続部27、27aは、同一のXB接続部であってもよい。また、各CPU21~21cは、PCIE接続部28と接続されている。また、通信部23は、サービスプロセッサ24が有する通信部26と接続されている。なお、制御部25、通信部2

50

6、通信部23、各CPU21~21cは、例えば、JTAG (Joint Test Action Group) やI2C (Inter-Integrated Circuit) で接続されている。

【0030】

例えば、図2に示す例では、CPU21~21cは、演算処理を実行する演算処理装置である。また、各CPU21~21cには、それぞれ独立したメモリ22~22cが接続されている。また、各CPU21~21cは、メモリ22~22cや、他のビルディングブロック10a~10eが有するメモリを共有メモリとして利用する。また、各CPU21~21cは、後述するように、物理アドレスと、物理アドレスが割り振られたメモリと接続されたCPUの識別子であるCPUID (identification) とを対応付けたノードマップを有する。

10

【0031】

そして、例えば、CPU21は、アクセス対象となる物理アドレスと対応付けられたCPUIDが、ノード20とは異なるノードが有するCPUを示す場合には、XB接続部27およびXB2を介して他のノードにメモリアクセスのリクエストを送信する。また、CPU21は、アクセス対象となる物理アドレスと対応付けられたCPUIDが、CPU21a~21cを示す場合には、CPU間の直接接続を介して、メモリアクセスのリクエストを送信する。すなわち、CPU21は、アクセス対象となる物理アドレスと対応付けられたCPUIDが、CPU21以外のCPUであって、自身と同じノード20に存在するCPUを示す場合には、CPU間の直接接続を介して、メモリアクセスのリクエストを送信する。

20

【0032】

また、CPU21は、自身と接続されたメモリに対するリクエストを他のノードから受信した場合には、リクエストの対象となるデータを自身と接続されたメモリ22から読み出し、リクエスト元へ送信する。

【0033】

なお、各CPU21~21cは、実行中のアプリケーションが共有メモリの割り当てを要求した場合には、相互に通信を行い、アプリケーションが使用する共有メモリの割り当てを行う機能を有する。また、各CPU21~21cは、TLBを用いたアドレス変換を行うとともに、TLBミスが発生した際に、トラップ処理を実行するなど、従来のCPUと同様の処理を実行する機能を有するものとする。

30

【0034】

メモリ22~22cは、情報処理システム1が有する全てのCPUが共用するメモリである。また、情報処理システム1においては、全てのビルディングブロック10~10eが有するメモリに対して、各ビルディングブロック10~10eのサービスプロセッサが、同一の物理アドレス空間にマッピングされる物理アドレスを振分ける。すなわち、情報処理システム1が有する全てのメモリには、重複しない値の物理アドレスが割り当てられている。

【0035】

また、メモリ22~22cは、記憶領域の一部を、情報処理システム1が有する全てのCPUが共用する共有領域とし、他の部分を、自身にアクセスするCPU21~21cがカーネルデータやユーザデータを格納するローカル領域とする。また、メモリ22~22cには、情報処理システム1が用いる物理アドレス空間のうち、ある位置のビットが同一の値となる範囲の物理アドレスが有領域に振り分けられる。また、メモリ22~22cには、ある位置のビットが共有領域に振り分けた物理アドレスとは異なる値となる範囲の物理アドレスがローカル領域に振り分けられる。

40

【0036】

例えば、メモリ22~22cには、46ビット目のビットが「0」となる物理アドレスがローカル領域に振り分けられ、46ビット目のビットが「1」となる物理アドレスが共有領域に振り分けられる。詳細な例を挙げると、メモリ22~22cのローカル領域に対しては、物理アドレス空間のうち、「0」~「0x63ffffffffff」に含

50



まれる物理アドレスが振り分けられる。また、メモリ 22 ~ 22 c の共有領域に対しては、物理アドレス空間のうち、「0x6400 000 0000」~「0x1 27ff ffff ffff」に含まれる物理アドレスが振り分けられる。

【0037】

なお、情報処理システム1においては、各ビルディングブロック10 ~ 10 g ごとに、異なる範囲に含まれる物理アドレスをメモリに振り分ける。以下、図面を用いて、情報処理システム1において、各ビルディングブロック10 ~ 10 e ごとに、メモリに振り分ける物理アドレスの範囲を説明する。

【0038】

図3は、実施例1に係るビルディングブロックのメモリに振り分けられる物理アドレスの範囲を説明するための図である。なお、図3に示す例では、各ビルディングブロックをBB (Building Block) と記載した。また、BB # 0 とは、ビルディングブロック10を示し、BB # 1 とは、ビルディングブロック10 aを示し、BB # 15 とは、ビルディングブロック10 eを示す。すなわち、図3に示す例では、情報処理システム1は、16個のビルディングブロックを有するものとする。

10

【0039】

また、図3に示す例では、各ビルディングブロックには、最大で4TB (Terabyte) のメモリを搭載可能であるものとする。また、以下の説明においては、メモリアドレスの表記を簡易化するため、例えば「2<sup>42</sup>」となるアドレス番地を「4TB」と記載する。

【0040】

20

図3に示す例では、ビルディングブロック10が有するメモリ22 ~ 22 c においては、物理アドレス空間のうち、「0」~「4TB - 1」までの範囲に含まれる物理アドレスがローカル領域に振り分けられる。また、ビルディングブロック10が有するメモリ22 ~ 22 c においては、物理アドレス空間のうち、「64TB」~「68TB - 1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。

【0041】

また、ビルディングブロック10 aが有するメモリにおいては、物理アドレス空間のうち、「4TB」~「8TB - 1」までの範囲に含まれる物理アドレスがローカル領域に振り分けられる。また、ビルディングブロック10 aが有するメモリにおいては、物理アドレス空間のうち、「68TB」~「72TB - 1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。

30

【0042】

また、ビルディングブロック10 eが有するメモリにおいては、物理アドレス空間のうち、「60TB」~「64TB - 1」までの範囲に含まれる物理アドレスがローカル領域に振り分けられる。また、ビルディングブロック10 aが有するメモリにおいては、物理アドレス空間のうち、「124TB」~「128TB - 1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。

【0043】

この結果、情報処理システム1は、図4に示すように、物理アドレス空間を各ビルディングブロック10 ~ 10 e が有する全てのメモリに対して振り分けることとなる。図4は、実施例1に係る情報処理システムが各メモリに振り分ける物理アドレスを説明するための図である。

40

【0044】

具体的には、図4に示す例では、情報処理システム1は、「0」~「256TB - 1」までの物理アドレスのうち、「0」~「64TB - 1」までの範囲をローカル領域に振り分ける物理アドレスとする。また、情報処理システム1は、「64TB」~「128TB - 1」までの範囲を共有領域に振り分ける物理アドレスとする。

【0045】

すなわち、情報処理システム1は、最下位のビットを0ビット目として46ビット目のビットが「0」の範囲をローカル領域に振り分け、「1」の範囲を共有領域に振り分ける

50

。なお、情報処理システム 1 は、「128TB」～「256TB-1」までの範囲を I/O 空間として用いる。

【0046】

なお、図 3、4 に示す例は、あくまで一例であり、情報処理システム 1 は、異なる振り分け方を採用してもよい。以下、情報処理システム 1 が、物理アドレスを振り分けるバリエーションの例について図を用いて説明する。

【0047】

図 5 は、物理アドレスの振り分けのバリエーションを説明するための第 1 の図である。図 5 に示す例では、各ビルディングブロック 10 ~ 10 e が有するメモリにおいては、「0」～「4TB-1」までの範囲に含まれる物理アドレスがローカル領域に振り分けられる。また、図 5 に示す例では、ビルディングブロック 10 が有するメモリ 22 においては、「4TB」～「8TB-1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。

10

【0048】

また、図 5 に示す例では、ビルディングブロック 10 a が有するメモリにおいては、「8TB」～「12TB-1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。また、図 5 に示す例では、ビルディングブロック 10 e が有するメモリにおいては、「64TB」～「68TB-1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。

【0049】

この結果、図 5 に示す例では、情報処理システム 1 は、物理アドレス空間のうち、「0」～「4TB-1」までの範囲の物理アドレスをローカル領域に振り分け、「4TB」～「128TB-1」までの範囲の物理アドレスを共有領域に振り分ける。また、図 5 に示す例では、情報処理システム 1 は、「128TB」～「256TB-1」までの範囲を I/O 空間として用いる。すなわち、情報処理システム 1 は、最下位のビットを 0 ビット目として 42 ビット目のビットが「0」の範囲をローカル領域に振り分け、「1」の範囲を共有領域に振り分ける。

20

【0050】

また図 6 は、物理アドレスの振り分けのバリエーションを説明するための第 2 の図である。図 6 に示す例では、各ビルディングブロック 10 ~ 10 e が有するメモリにおいては、「0」～「4TB-1」までの範囲に含まれる物理アドレスが I/O 空間用に保存される。また、図 6 に示す例では、各ビルディングブロック 10 ~ 10 e が有するメモリにおいては、「4TB」～「8TB-1」までの範囲に含まれる物理アドレスがローカル領域に振り分けられる。

30

【0051】

また、図 6 に示す例では、ビルディングブロック 10 が有するメモリ 22 ~ 22 c においては、「8TB」～「12TB-1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。また、図 6 に示す例では、ビルディングブロック 10 a が有するメモリにおいては、「12TB」～「16TB-1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。また、図 6 に示す例では、ビルディングブロック 10 e が有するメモリにおいては、「68TB」～「72TB-1」までの範囲に含まれる物理アドレスが共有領域に振り分けられる。

40

【0052】

この結果、図 6 に示す例では、情報処理システム 1 は、物理アドレス空間のうち、「0」～「4TB-1」までの範囲を I/O 空間とし、「4TB」～「8TB-1」までの範囲の物理アドレスをローカル領域に振り分けることとなる。また、図 5 に示す例では、情報処理システム 1 は、「8TB」～「256TB-1」までの範囲の物理アドレスを共有領域に振り分けることとなる。すなわち、情報処理システム 1 は、最下位のビットを 0 ビット目として 43 ビット目のビットが「0」の範囲をローカル領域に振り分け、「1」の範囲を共有領域に振り分ける。

50

## 【 0 0 5 3 】

図 2 に戻って、制御部 2 5 は、ビルディングブロック 1 0 の制御を行う。例えば、制御部 2 5 は、ビルディングブロック 1 0 の電源管理や、ビルディングブロック 1 0 内の異常の監視や制御等を実行する。また、制御部 2 5 は、管理用ネットワークを介して、管理端末 3 や他のビルディングブロック 1 0 ~ 1 0 e が有するサービスプロセッサの制御部とも接続されており、管理端末 3 によって指示された制御や、各ビルディングブロック 1 0 ~ 1 0 e 間で連係した制御を実行できる。また、制御部 2 5 は、各 CPU 2 1 ~ 2 1 c が実行する OS と通信を行うことができる。

## 【 0 0 5 4 】

なお、実施例 1 では、各ビルディングブロック 1 0 ~ 1 0 e が有するサービスプロセッサは、管理用ネットワークを介して接続されているが、実施例はこれに限定されるものではない。たとえば、各ビルディングブロック 1 0 ~ 1 0 e を接続する X B を介して相互に通信しても良い。

10

## 【 0 0 5 5 】

また、制御部 2 5 は、通信部 2 6 と通信部 2 3 を介して、各 CPU 2 1 ~ 2 1 c にアクセスする。そして、制御部 2 5 は、後述するように、各ビルディングブロック 1 0 ~ 1 0 e が有するノードマップの更新や制御等を実行する。

## 【 0 0 5 6 】

なお、通信部 2 3 は、サービスプロセッサ 2 4 が有する通信部 2 6 を介して、制御部 2 5 による制御信号を各 CPU 2 1 ~ 2 1 c に伝達する。また、通信部 2 6 は、制御部 2 5 による制御信号をノード 2 0 が有する通信部 2 3 に伝達する。また、X B 接続部 2 7、2 7 a は、各 CPU 2 1 ~ 2 1 c を X B 2 と接続し、各ビルディングブロック 1 0 ~ 1 0 e が有する CPU 間の通信を中継する。また、P C I e 接続部 2 8 は、各 CPU 2 1 ~ 2 1 c による I / O ( Input Output ) 装置へのアクセスを中継する。

20

## 【 0 0 5 7 】

次に、図 7 を用いて、各 CPU 2 1 ~ 2 1 c が有する機能構成について説明する。図 7 は、実施例 1 に係る CPU の機能構成を説明するための図である。なお、CPU 2 1 a ~ 2 1 c は、CPU 2 1 と同様の機能を発揮するものとして、説明を省略する。また、図 7 に示す例では、サービスプロセッサ 2 4 と CPU 2 1 とを接続する接続部 2 3、2 6 については、記載を省略した。

30

## 【 0 0 5 8 】

図 7 に示す例では、CPU 2 1 は、演算処理部 3 0、ルータ 4 0、メモリアクセス部 4 1、P C I e 制御部 4 2 を有する。また、演算処理部 3 0 は、演算部 3 1、L 1 ( Level 1 ) キャッシュ 3 2、L 2 キャッシュ 3 3、ノードマップ 3 4、アドレス変換部 3 5、キャッシュディレクトリ管理部 3 6、パケット制御部 3 7 を有する。また、パケット制御部 3 7 は、リクエスト生成部 3 8、リクエスト受信部 3 9 を有する。また、P C I e 制御部 4 2 は、リクエスト生成部 4 3、P C I e バス制御部 4 4 を有する。

## 【 0 0 5 9 】

まず、演算処理部 3 0 が有するノードマップ 3 4 について説明する。ノードマップ 3 4 は、物理アドレスと、物理アドレスが示す記憶領域を有するメモリと接続された CPU の CPU I D とを対応付けて記憶する。以下、ノードマップ 3 4 が記憶する情報の例を図面を用いて説明する。

40

## 【 0 0 6 0 】

図 8 は、実施例 1 に係るノードマップが記憶する情報の一例を説明するための図である。図 8 に示す例では、ノードマップ 3 4 は、アドレス、バリッド、ノード I D、CPU I D とを対応付けたエントリを記憶する。ここで、各エントリのアドレスには、連続する複数の物理アドレスを含むアドレス域を示す情報が格納される。

## 【 0 0 6 1 】

例えば、情報処理システム 1 は、全てのメモリに対して振り分けた物理アドレス空間を均等な大きさのアドレス域に分割し、各アドレス域に # 0、# 1、# 2 等の識別子を付与

50

する。そして、情報処理システム 1 は、各アドレス域を示す識別子を、ノードマップ 3 4 が有する各エントリのアドレスに格納する。

【 0 0 6 2 】

また、各エントリのバリッドには、物理アドレスが示す記憶領域にアクセスすることができるか否かを示すバリッドビットが格納される。例えば、物理アドレスが示す記憶領域が、各 CPU で共有される共有領域である場合には、アクセスを行う事ができる旨のバリッドビット（例えば「 1 」）が格納される。

【 0 0 6 3 】

また、ノード ID とは、物理アドレスが振り分けられたメモリが存在するノードを示す識別子である。また、CPU ID とは、物理アドレスが振り分けられたメモリと接続された CPU を示す識別子である。すなわち、ノードマップ 3 4 は、アクセス対象となる物理アドレスが、どの CPU と接続されたメモリの物理アドレスであることを示す情報を記憶する。

10

【 0 0 6 4 】

例えば、図 8 に示す例では、ノードマップ 3 4 は、識別子が「 # 0 」のアドレス域が、ノード ID 「 0 」で示されるノードに存在し、CPU ID が「 0 」の CPU がアクセスを行う旨を示す。また、ノードマップ 3 4 は、識別子が「 # 1 」のアドレス域が、ノード ID 「 0 」で示されるノードに存在し、CPU ID が「 1 」の CPU がアクセスを行う旨を示す。また、ノードマップ 3 4 は、識別子が「 # 2 」のアドレス域が、CPU 2 1 がアクセスを行わない、又は、マッピングされていないアドレス域であるため、ノード ID と CPU ID とが設定されていない旨を示す。

20

【 0 0 6 5 】

なお、ノードマップ 3 4 は、アクセス対象となる物理アドレスがどの CPU と接続された物理アドレスであることを示すことができれば、本実施例以外の任意の形式で情報を登録することとしてよい。以下、ノードマップ 3 4 のバリエーションの例について、図 9 および図 1 0 を用いて説明する。

【 0 0 6 6 】

図 9 は、ノードマップが記憶する情報のバリエーションの一例を説明するための第 1 の図である。図 9 に示す例では、ノードマップ 3 4 は、バリッド、スタートアドレス、アドレスマスク、ノード ID、CPU ID を対応付けてエントリを記憶する。ここで、スタートアドレスとは、アドレス域に含まれる物理アドレスのうち、最若番の物理アドレスが格納される。

30

【 0 0 6 7 】

また、アドレスマスクには、CPU が管理する物理アドレスの範囲を示すアドレスマスクが格納される。例えば、あるエントリは、アドレスマスクが「 0 x f f f f f f f f 0 0 0 0 」である場合には、同一エントリのスタートアドレスと上位 4 8 ビットが一致するアドレス領域を、同一エントリの CPU ID が示す CPU が管理することを示す。

【 0 0 6 8 】

例えば、図 9 に示す例では、ノードマップ 3 4 は、最初のエントリとして、アドレス「 0 x 0 0 0 0 0 」からアドレスマスク「 0 x 3 f f f 」でマスクされる範囲、すなわち「 0 x 0 3 f f f 」までの範囲が 1 つのアドレス域である旨を示す。また、ノード 3 4 は、「 0 x 0 0 0 0 0 」から「 0 x 0 3 f f f 」のアドレス域が、ノード ID 「 0 」で示されるノードに存在し、CPU ID が「 0 」の CPU がアクセスするアドレス域である旨を示す。

40

【 0 0 6 9 】

同様に、ノードマップ 3 4 は、「 0 x 1 0 0 0 0 」から「 0 x 1 3 f f f 」のアドレス域が、ノード ID 「 1 」で示されるノードに存在し、CPU ID が「 4 」の CPU がアクセスするアドレス域である旨を示す。また、ノードマップ 3 4 は、「 0 x 1 4 0 0 0 」から「 0 x 1 7 f f f 」のアドレス域が、ノード ID 「 1 」で示されるノードに存在し、C

50

P U I D が「 5 」の C P U がアクセスするアドレス域である旨を示す。また、ノードマップ 3 4 は、「 0 x 2 0 0 0 0 」から「 0 x 2 1 f f f 」のアドレス域が、ノード I D 「 2 」で示されるノードに存在し、C P U I D が「 8 」の C P U がアクセスするアドレス域である旨を示す。

【 0 0 7 0 】

なお、ノードマップ 3 4 は、図 9 に示すように、アドレス域をスタートアドレスとアドレスマスクとで表現した場合には、物理アドレスが各アドレス域に含まれるか否かを論理和と論理積の組合せで実行することができるため、回路構成が容易となる。

【 0 0 7 1 】

また、図 1 0 は、ノードマップが記憶する情報のバリエーションの一例を説明するための第 2 の図である。図 1 0 に示す例では、ノードマップ 3 4 は、パリッド、スタートアドレス、レンジス、ノード I D、C P U I D を対応付けたエントリを記憶する。ここで、レンジスとは、アドレス域の大きさを設定する情報である。

【 0 0 7 2 】

例えば、スタートアドレスが「 0 x 1 2 0 0 0 0 」で、レンジスが「 0 x 1 f f f f 」とすると、同一エントリの C P U I D が示す C P U は、管理するメモリに対して、物理アドレス「 0 x 1 2 0 0 0 0 」から「 0 x 1 3 f f f f 」を割当てることとなる。

【 0 0 7 3 】

例えば、図 1 0 に示す例では、ノードマップ 3 4 は、最初のエントリとして、アドレス「 0 x 0 0 0 0 0 」からレンジスが「 0 x 3 f f f 」に含まれる範囲、すなわち「 0 x 0 3 f f f 」までの範囲が 1 つのアドレス域である旨を示す。また、ノード 3 4 は、「 0 x 0 0 0 0 0 」から「 0 x 0 3 f f f 」のアドレス域が、ノード I D 「 0 」で示されるノードに存在し、C P U I D が「 0 」の C P U がアクセスするアドレス域である旨を示す。

【 0 0 7 4 】

同様に、ノードマップ 3 4 は、「 0 x 1 0 0 0 0 」から「 0 x 1 3 f f f 」のアドレス域が、ノード I D 「 1 」で示されるノードに存在し、C P U I D が「 4 」の C P U がアクセスするアドレス域である旨を示す。また、ノードマップ 3 4 は、「 0 x 1 4 0 0 0 」から「 0 x 1 7 f f f 」のアドレス域が、ノード I D 「 1 」で示されるノードに存在し、C P U I D が「 5 」の C P U がアクセスするアドレス域である旨を示す。また、ノードマップ 3 4 は、「 0 x 2 0 0 0 0 」から「 0 x 2 0 2 e f 」のアドレス域が、ノード I D 「 2 」で示されるノードに存在し、C P U I D が「 8 」の C P U がアクセスするアドレス域である旨を示す。

【 0 0 7 5 】

なお、ノードマップ 3 4 は、図 1 0 に示すように、アドレス域をスタートアドレスとレンジスとで表現した場合には、各アドレス域の長さを柔軟に設定することができる。すなわち、ノードマップ 3 4 は、アドレス域をスタートアドレスとアドレスマスクとで表現した場合は、L S B (Least Significant Bit) から 1 が連続する範囲のアクセス域を指定することとなる。一方、各アドレス域をスタートアドレスとレンジスとで表現した場合には、各アドレス域の長さを任意の長さに設定することができる。

【 0 0 7 6 】

図 7 に戻って、演算部 3 1 は、演算処理を実行し、O S やアプリケーションを実行する演算装置のコアである。また、演算部 3 1 は、データの読み込みを行う場合には、読み込み対象となるデータが格納された記憶領域の論理アドレスをアドレス変換部 3 5 に出力する。

【 0 0 7 7 】

L 1 キャッシュ 3 2 は、データやディレクトリのうち頻繁に利用されるデータを一時的に記憶するキャッシュメモリである。L 2 キャッシュ 3 3 は、L 1 キャッシュ 3 2 と同様に、データやディレクトリのうち頻繁に利用されるデータを一時的に記憶するが、L 1 キャッシュ 3 2 よりも記憶容量が大きく、データを読み書きする速度が低速なキャッシュメモリである。なお、ディレクトリとは、メモリ 2 2 の各記憶領域に記憶されたデータをキ

10

20

30

40

50

キャッシュしたCPUや、キャッシュされたデータの更新状況を示す情報である。

【0078】

アドレス変換部35は、TLB(Translation Lookaside Buffer)を用いて、演算部31が出力した論理アドレスを物理アドレスに変換する。例えば、アドレス変換部35は、論理アドレスと物理アドレスとを対応付けたエントリを記憶するTLBを有し、演算部31から取得した論理アドレスと対応付けて記憶する物理アドレスをキャッシュディレクトリ管理部36に出力する。なお、アドレス変換部35は、TLBミスが発生した場合は、トラップ処理を実行し、TLBミスした物理アドレスと論理アドレスの組をTLBに登録する。

【0079】

また、アドレス変換部35は、CPU21が実行するアプリケーションから共有メモリへの割当てを要求された場合には、以下の処理を実行する。すなわち、アドレス変換部35は、各CPU21~21cが共用する共有領域にアクセスする際にアプリケーションが用いる論理アドレスと、共有領域に割当てられる範囲の物理アドレスとを対応付けたエントリをTLBに設定する。

【0080】

また、アドレス変換部35は、アプリケーションやOSからローカル領域の割当てを要求された場合は、以下の処理を実行する。すなわち、アドレス変換部35は、アプリケーションやOSがCPU21専用のローカル領域にアクセスするための論理アドレスと、ローカル領域に割当てられる範囲の物理アドレスとを対応付けたエントリをTLBに設定する。

【0081】

キャッシュディレクトリ管理部36は、キャッシュデータおよびディレクトリの管理を行う。具体的には、キャッシュディレクトリ管理部36は、アドレス変換部35から、演算部31が出力した論理アドレスを変換した物理アドレスを取得する。

【0082】

そして、キャッシュディレクトリ管理部36は、アドレス変換部35から物理アドレスを取得した場合には、ディレクトリをチェックし、物理アドレスが示すデータの状態が正常かチェックする。また、物理アドレスが示すデータをL1キャッシュ32またはL2キャッシュ33がキャッシュしている場合には、キャッシュしているデータを演算部31に出力する。

【0083】

一方、キャッシュディレクトリ管理部36は、物理アドレスが示すデータをL1キャッシュ32またはL2キャッシュ33がキャッシュしていない場合は、物理アドレスが示す記憶領域が、メモリ22に存在するか否かを判別する。そして、キャッシュディレクトリ管理部36は、物理アドレスが示す記憶領域がメモリ22に存在しない場合には、ノードマップ34を参照する。

【0084】

また、キャッシュディレクトリ管理部36は、ノードマップ34を参照し、取得した物理アドレスを含む範囲のエントリを識別する。そして、キャッシュディレクトリ管理部36は、識別したエントリのCPUIDがCPU21のCPUIDであるか否かを判別する。その後、キャッシュディレクトリ管理部36は、識別したエントリのCPUIDがCPU21のCPUIDである場合は、メモリアクセス部41に、物理アドレスを出力する。

【0085】

また、キャッシュディレクトリ管理部36は、識別したエントリのCPUIDがCPU21のCPUIDではない場合には、以下の処理を実行する。すなわち、キャッシュディレクトリ管理部36は、識別したエントリのCPUIDとノードIDとを取得する。そして、キャッシュディレクトリ管理部36は、パケット制御部37に対して、取得したCPUIDと物理アドレスとを出力する。

【0086】

10

20

30

40

50

なお、キャッシュディレクトリ管理部 3 6 は、出力した物理アドレスが示す記憶領域に格納されているデータをメモリアクセス部 4 1 やパケット制御部 3 7 から取得した場合には、取得したデータを L 1 キャッシュ 3 2 および L 2 キャッシュ 3 3 に格納する。そして、キャッシュディレクトリ管理部 3 6 は、L 1 キャッシュ 3 2 にキャッシュさせたデータを演算部 3 1 に出力する。

【 0 0 8 7 】

また、キャッシュディレクトリ管理部 3 6 は、パケット制御部 3 7 から物理アドレスを取得した場合、すなわち、他の CPU からのメモリアクセスのリクエストの対象となる物理アドレスを取得した場合には、以下の処理を実行する。すなわち、キャッシュディレクトリ管理部 3 6 は、取得した物理アドレスのうち、所定の位置のビットが「 0 」であるか  
10 「 1 」であるかに応じて、取得した物理アドレスがローカル領域に振り分けられた物理アドレスであるか否かを判別する。

【 0 0 8 8 】

例えば、キャッシュディレクトリ管理部 3 6 は、情報処理システム 1 の各メモリに対して、図 3、図 4 に例示した範囲の物理アドレスが振り分けられている場合には、最下位のビットを 0 ビット目として 4 6 ビット目が「 0 」であるか「 1 」であるかを判別する。そして、キャッシュディレクトリ管理部 3 6 は、4 6 ビット目が「 0 」である場合には、取得した物理アドレスがローカル領域に振り分けられた物理アドレスであると判別する。このような場合には、キャッシュディレクトリ管理部 3 6 は、パケット制御部 3 7 に対して、  
20 リクエスト元に否定応答（アクセスエラー）を送信するように指示する。

【 0 0 8 9 】

また、キャッシュディレクトリ管理部 3 6 は、4 6 ビット目が「 1 」である場合には、取得した物理アドレスが共有領域に振り分けられた物理アドレスであると判別する。このような場合には、キャッシュディレクトリ管理部 3 6 は、取得した物理アドレスが示す記憶領域に記憶されたデータを取得し、取得したデータをパケット制御部 3 7 に出力し、リクエスト元へ送信するよう指示する。

【 0 0 9 0 】

なお、キャッシュディレクトリ管理部 3 6 は、メモリ 2 2 に格納されたデータにアクセスする場合には、物理アドレスが示す記憶領域のデータと、キャッシュされたデータとのコピーレンスを保持する処理を行う。例えば、キャッシュディレクトリ管理部 3 6 は、  
30 キャッシュエントリごとにキャッシュデータの状態を示すキャッシュタグと、ディレクトリとを参照する。そして、キャッシュディレクトリ管理部 3 6 は、キャッシュタグとディレクトリとに基づいて、キャッシュコピーレンスを保持する処理、および、メモリアクセス処理を実行する。

【 0 0 9 1 】

ここで、図 1 1 A は、キャッシュタグの一例を説明するための図である。図 1 1 A に示す例では、キャッシュタグは、縮退フラグ、E C C (Error Check and Correct memory) チェックビット、I F (Instruction Fetch) / オプコード、L 1 キャッシュステート、L 2 キャッシュステート、A A とを有する。

【 0 0 9 2 】

ここで、縮退フラグとは、縮退するか否かを示すキャッシュライン縮退情報である。また、E C C チェックビットとは、冗長化のために付加されるチェックビットである。I F / オプコードとは、データがインストラクションであるかデータであることを示す情報である。  
40

【 0 0 9 3 】

また、A A とは、アドレス情報であり、詳細には、物理アドレスのフレームアドレスが格納される。また、L 1 キャッシュステート、および、L 2 キャッシュステートとは、L 1  
50 キャッシュ 3 2 および L 2 キャッシュ 3 3 に格納されたデータの状態を示す情報である。

【 0 0 9 4 】

10

20

30

40

50

例えば、L1 キャッシュステートやL2 キャッシュステートには、「M (Modified)」、「E (Exclusive)」、「S (Shared)」、「I (Invalid)」のいずれかを示すビットが格納される。ここで、Modifiedとは、いずれか1つのCPUがデータをキャッシュしており、かつ、キャッシュされたデータが更新されている状態を示す。なお、キャッシュされたデータの状態がModifiedである場合には、ライトバックを実行する必要がある。

【0095】

また、Exclusiveとは、いずれか1つのCPUがデータをキャッシュしており、かつ、キャッシュされたデータが更新されていない状態を示す。また、Sharedとは、複数のCPUがデータをキャッシュしており、かつ、キャッシュされたデータが更新されていないことを示す。なお、Invalidとは、キャッシュのステータスが登録されていないことを示す。

10

【0096】

一方、ディレクトリは、2ビットのCKビット、63ビットのPRC、4ビットのUEを管理する。ここで、CKビットとは、キャッシュされたデータの状態をコード化した情報である。また、PRCとは、当該キャッシュラインのデータをキャッシュしたCPUの位置をビットマップで示す情報である。また、UEとは、ディレクトリの異常と要因とを示す情報である。

【0097】

キャッシュディレクトリ管理部36は、取得した物理アドレスに格納されたデータをキャッシュするCPUや、キャッシュされたデータの状態等を識別する。そして、キャッシュディレクトリ管理部36は、キャッシュされたデータの状態に基づいて、スヌープを発行してメモリのデータを更新する等の処理を行い、キャッシュされたデータとメモリのデータとのコヒーレンスを保持する。その後、キャッシュディレクトリ管理部36は、データを要求元に出力する。

20

【0098】

ここで、キャッシュディレクトリ管理部36がキャッシュコヒーレンスを保持する処理の一例について説明する。例えば、キャッシュディレクトリ管理部36は、ステータスがM (Modified)であるデータをキャッシュしたCPUに対してライトバックを指示する命令を送信するようパケット生成部38に指示する。そして、キャッシュディレクトリ管理部36は、データのステータスを更新し、更新後のステータスに応じた処理を実行する。なお、キャッシュディレクトリ管理部36が送受信するリクエストや命令の種別については、後述する。

30

【0099】

リクエスト生成部38は、キャッシュディレクトリ管理部36から物理アドレスと、CPU IDとを取得した場合には、取得した物理アドレスと、CPU IDとを格納したパケット、すなわち、メモリアクセスのリクエストとなるパケットを生成する。そして、リクエスト生成部38は、生成したパケットをルータ40に送信する。

【0100】

ここで、図11Bは、実施例1に係るCPUが送信するパケットを説明するための図である。なお、図11Bに示す例では、物理アドレスをPA (Physical Address)と記載した。図11Bに示す例では、リクエスト生成部38は、CPU IDと物理アドレスと、リクエストの内容を示すデータとが格納されたリクエストを生成し、生成したリクエストをルータ40に出力する。このような場合には、ルータ40は、リクエスト生成部38が生成したリクエストをXB接続部27を解してXB2に出力する。すると、XB2は、リクエストに格納されたCPU IDが示すCPUへとリクエストを転送する。

40

【0101】

なお、リクエスト生成部38は、キャッシュディレクトリ管理部36からコヒーレンスを保持するためのリクエストや命令の発行の指示を受付けた場合には、指示されたリクエストや命令を生成する。そして、リクエスト生成部38は、生成したリクエストや命令を

50



ルータ40、XB接続部27、XB2を介して、指示されたCPUに送信する。なお、リクエスト生成部38は、I/O装置からデータを取得する場合は、I/Oに対するアクセス要求をルータ40に出力する。

【0102】

図7に戻って、リクエスト受信部39は、XB2、XB接続部27、ルータ40を介して、他のCPUが出力したパケットを受信すると、受信したパケットに含まれる物理アドレスを取得する。そして、リクエスト受信部39は、取得した物理アドレスをキャッシュディレクトリ管理部36に出力する。また、リクエスト受信部39は、他のCPUが送信したデータを受信した場合には、受信したデータをキャッシュディレクトリ管理部36に出力する。

10

【0103】

なお、リクエスト受信部39は、コヒーレンスを保持するためのリクエストや命令を受信した場合には、受信したリクエストや命令をキャッシュディレクトリ管理部36に出力する。また、リクエスト受信部39は、I/Oに対するアクセス要求の応答やデータをルータ40から受信した場合は、受信した応答やデータをキャッシュディレクトリ管理部36に出力する。このような場合には、キャッシュディレクトリ管理部36は、例えば、取得したデータをメモリアクセス部41に出力し、メモリ22に格納する処理を行う。

【0104】

ルータ40は、パケット制御部37が有するリクエスト生成部38が出力したパケットを受信した場合には、受信したリクエストをXB接続部27に出力する。また、ルータ40は、XB接続部27を介して、他のCPUが送信したパケットやデータをリクエスト受信部39に出力する。また、ルータ40は、パケット制御部37がI/O等に対して出力したパケットをPCIe制御部42に出力する。また、ルータ40は、I/Oからの応答等をPCIe制御部42から受信した場合には、受信した応答等をパケット制御部37に出力する。

20

【0105】

メモリアクセス部41は、いわゆるMAC (Memory Access Controller) であり、メモリ22に対するアクセスの制御を行う。例えば、メモリアクセス部41は、キャッシュディレクトリ管理部36から物理アドレスを受信した場合には、受信した物理アドレスに格納されたデータをメモリ22から取得し、取得したデータをキャッシュディレクトリ管理部36に出力する。なお、メモリアクセス部41は、メモリーミラー機能を用いて、共有領域を冗長化してもよい。

30

【0106】

PCIe制御部42が有するリクエスト生成部43は、ルータ40を介してI/Oに対するアクセス要求を取得した場合には、アクセス要求の対象となるI/O装置に送信するリクエストを生成し、生成したリクエストをPCIeバス制御部44に出力する。PCIeバス制御部44は、リクエスト生成部43が生成したリクエストを取得した場合には、PCIe接続部28を介して、I/O装置にリクエストを送信する。

【0107】

次に、図12を用いて、CPU21が他のCPUに対してリクエストを送信する処理の一例について説明する。図12は、実施例1に係るCPUがリクエストを送信する処理の一例を説明するための図である。例えば、図12中(A)に示すように、サービスプロセッサ24からノードマップ34に対して、物理アドレスが振り分けられるメモリにアクセスするCPUのCPUIDと物理アドレスとを対応付けたエントリの設定が行われる。

40

【0108】

また、演算部31は、演算処理を実行し、図12中(B)に示すように、アクセス対象となる論理アドレスをアドレス変換部35に出力する。すると、アドレス変換部35は、論理アドレスを物理アドレスに変換し、変換した物理アドレスを図12中(C)に示すように、キャッシュディレクトリ管理部36に出力する。

【0109】

50

ここで、キャッシュディレクトリ管理部 36 は、アドレス変換部 35 から物理アドレスを取得すると、図 12 中 (D) に示すように、ノードマップ 34 を参照し、取得した物理アドレスと対応付けられた CPU ID を取得する。そして、キャッシュディレクトリ管理部 36 は、取得した CPU ID が CPU 21 の CPU ID ではない場合には、図 12 中 (E) に示すように、取得した CPU ID と物理アドレスとをパケット制御部 37 に出力する。

【0110】

このような場合には、リクエスト生成部 38 は、キャッシュディレクトリ管理部 36 から取得した物理アドレスと CPU ID とを格納したパケットを生成し、図 12 中 (F) に示すように、生成したパケットをルータ 40 に出力する。すると、図 12 中 (G) に示すように、ルータ 40 は、リクエスト生成部 38 から取得したパケットを XB 接続部 27 に出力する。その後、図 12 中 (H) に示すように、XB 接続部 27 は、取得したパケットを XB 2 に出力する。すると、XB 2 は、パケットに格納された CPU ID が示す CPU へパケットを伝達することとなる。

10

【0111】

次に、図 13 を用いて、CPU 21 が他の CPU からパケットを受信した際に実行する処理の一例について説明する。図 13 は、実施例 1 に係る CPU がパケットを受信した際に実行する処理の一例を説明するための図である。例えば、図 13 中 (I) に示すようにリクエスト受信部 39 は、他の CPU から CPU 21 の CPU ID とメモリ 22 に振り分けられた物理アドレスとが格納されたパケットを受信する。

20

【0112】

このような場合には、リクエスト受信部 39 は、受信したパケットから物理アドレスを取得し、図 13 中 (J) に示すように、取得した物理アドレスをキャッシュディレクトリ管理部 36 に出力する。すると、キャッシュディレクトリ管理部 36 は、取得した物理アドレスの 46 ビット目が「0」であるか「1」であるかを判別する。

【0113】

すなわち、キャッシュディレクトリ管理部 36 は、情報処理システム 1 が図 3、図 4 に示すように、共有領域とローカル領域に振り分ける物理アドレスを設定している場合には、物理アドレスの全ビットを識別せずともよい。すなわち、キャッシュディレクトリ管理部 36 は、46 ビット目が「0」であるか「1」であるかを判別するだけで、物理アドレスが示す記憶領域が、共有領域であるかローカル領域であるかを正確に判別することができる。

30

【0114】

そして、キャッシュディレクトリ管理部 36 は、受信した物理アドレスの 46 ビット目が「1」である場合には、共有領域に対するアクセスであると判別する。このような場合には、キャッシュディレクトリ管理部 36 は、図 13 中 (K) に示すように、物理アドレスが示す記憶領域のデータが L1 キャッシュ 32 および L2 キャッシュ 33 にキャッシュされているか判別する。

【0115】

また、キャッシュディレクトリ管理部 36 は、データがキャッシュされていないと判別した場合には、図 13 中 (L) に示すように、物理アドレスをメモリアクセス部 41 に出力する。すると、図 13 中 (M) に示すように、メモリアクセス部 41 は、メモリ 22 から物理アドレスが示す記憶領域のデータを取得し、キャッシュディレクトリ管理部 36 に出力する。

40

【0116】

そして、キャッシュディレクトリ管理部 36 は、L1 キャッシュ 32、L2 キャッシュ 33、またはメモリアクセス部 41 からデータを取得した場合には、取得したデータをパケット制御部 37 に出力し、リクエスト元の CPU に送信するよう指示する。

【0117】

例えば、CPU 21 ~ 21c、通信部 23、サービスプロセッサ 24、制御部 25、通

50

信部 26、XB 接続部 27、PCIe 接続部 28 とは、電子回路である。また、演算部 31、アドレス変換部 35、キャッシュディレクトリ管理部 36、パケット制御部 37、リクエスト生成部 38、リクエスト受信部 39 とは、電子回路である。

【0118】

また、ルータ 40、メモリアクセス部 41、PCIe 制御部 42、リクエスト生成部 43、PCIe バス制御部 44 とは、電子回路である。ここで、電子回路の例として、ASIC (Application Specific Integrated Circuit) や FPGA (Field Programmable Gate Array) などの集積回路、または CPU (Central Processing Unit) や MPU (Micro Processing Unit) などを適用する。

【0119】

また、メモリ 22 ~ 22a とは、RAM (Random Access Memory)、ROM (Read Only Memory)、フラッシュメモリ (flash memory) などの半導体メモリ素子である。また、L1 キャッシュ 32、L2 キャッシュ 33 は、SRAM (Static Random Access Memory) 等の高速な半導体メモリ素子である。

【0120】

次に、各 CPU 21 ~ 21c がキャッシュコヒーレンスを保持する処理について簡単に説明する。なお、以下の説明においては、情報処理システム 1 の各 CPU はイリノイプロトコルを用いて、キャッシュコヒーレンスを保持するものとする。

【0121】

なお、以下の説明においては、情報処理システム 1 が有する各メモリは、全ての CPU からキャッシュ可能な空間を有するメモリとして識別されるものとする。また、以下の説明においては、キャッシュ対象となるデータを記憶するメモリに、その CPU 内の MAC を介して物理的に直接接続されている CPU をホーム CPU とし、キャッシュを要求した CPU をローカル CPU と記載する。

【0122】

また、ホーム CPU に対して既にリクエストを送信し、データをキャッシュ済みである CPU をリモート CPU と記載する。なお、ローカル CPU とホーム CPU が同一の CPU となる場合や、ローカル CPU とリモート CPU とは同一の CPU となる場合も存在する。

【0123】

例えば、ローカル CPU は、自身のノードマップを参照し、アクセス対象となる物理アドレスがホーム CPU がアクセスするメモリに振り分けられていると判別する。そして、ローカル CPU は、物理アドレスを格納したリクエストをホーム CPU に対して発行する。なお、ローカル CPU が発行するリクエストには、複数の種別のリクエストが存在する。このため、ホーム CPU が有するキャッシュディレクトリ管理部は、取得したリクエストの種別に応じたキャッシュコヒーレンス制御を実行することとなる。

【0124】

例えば、ローカル CPU が発行するリクエストの種別としては、共有型フェッチアクセス、排他型フェッチアクセス、キャッシュ無効化要求、キャッシュリプレース要求等が存在する。共有型フェッチアクセスとは、Move Into Share の実行要求であり、ホーム CPU がアクセスするメモリからデータの読出しを行う際に発行されるリクエストである。

【0125】

また、排他型フェッチアクセスとは、例えば Move In Exclusively の実行要求であり、ホーム CPU がアクセスするメモリへデータストアを行う際の、キャッシュへのデータロードを行う際に発行される。また、キャッシュ無効化要求とは、例えば Move Out の実行要求であり、キャッシュラインの無効化をホーム CPU に対して要求する際に発行される。なお、ホーム CPU は、キャッシュ無効化要求を受信すると、リモート CPU に対してキャッシュ無効化要求を発行する場合や、キャッシュを Invalidation とさせる命令を発行する場合がある。

10

20

30

40

50

## 【0126】

キャッシュリプレース要求とは、例えばWrite Backの実行要求であり、更新されたキャッシュデータ、すなわちModified状態のキャッシュデータをホームCPUがアクセスするメモリに書き戻す際に発行される。なお、キャッシュリプレース要求には、例えば、Flush Backの実行要求であり、更新されていないキャッシュデータ、すなわち、Shared又はExclusive状態のキャッシュの破棄を行う際に発行される。

## 【0127】

ホームCPUは、上述したリクエストをローカルCPUから受信した場合には、リクエストを処理するために、ローカルCPUやリモートCPUに対して、命令を発行する。ここで、ホームCPUは、取得したリクエストの種別に応じてキャッシュコヒーレンス制御を実行するため、複数の種別の命令を発行することとなる。例えば、ホームCPUは、リモートCPUがキャッシュしているデータをローカルCPUにロードさせるMove Out and Bypass to Shareを発行する。

10

## 【0128】

また、例えば、ホームCPUは、ローカルCPU以外のすべてのリモートCPUのキャッシュを無効化し、その後、ホームCPUがローカルCPUにデータを送信するためのMove Out and Bypass Exclusivelyを発行する。また、ホームCPUは、リモートCPUにキャッシュの無効化を要求するMove Out WITH Invalidationを発行する。なお、ホームCPUがMove Out WITH Invalidationを発行した場合には、全てのCPUのキャッシュが、対象となるアドレスについてInvalidate状態となる。

20

## 【0129】

また、ホームCPUは、リモートCPUにキャッシュラインの無効化を要求するMove Out for Flushを発行する。なお、ホームCPUがMove Out for Flushを発行した場合には、対象となるデータを、ホームCPUのみがキャッシュした状態となる。また、ホームCPUは、対象となるデータの状態がSharedであるときに、リモートCPUにキャッシュの破棄を要求するBuffer Invalidationを発行する。

30

## 【0130】

ホームCPUは、リクエストの種別に応じて、上述した命令を発行し、各CPUがキャッシュしたデータのステートを遷移させる。また、ローカルCPUやリモートCPUは、命令を受信した場合には、命令が示す処理を実行し、自身がキャッシュしたデータのステートを遷移させる。

## 【0131】

その後、ローカルCPUやリモートCPUは、命令に対する完了応答やデータ付の完了応答をホームCPUに送信する。また、ホームCPUやリモートCPUは、命令処理を実行した後に、ローカルCPUに対して、データ付のリクエスト応答を送信することとなる。

## 【0132】

## [CPUの処理の流れ]

次に、図14を用いて、情報処理システム1において各CPUが有するノードマップ34を設定する処理の流れについて説明する、図14は、ノードマップを設定する処理の流れを説明するためのフローチャートである。なお、以下の説明においては、1つのCPUとCPUがアクセスするメモリとの組をノードとして記載する。また、以下の説明においては、新たなノードを情報処理システム1に追加する例について説明する。

40

## 【0133】

まず、情報処理システム1のオペレータは、ノードの新規増設を行う(ステップS101)。次に、各ビルディングブロック10~10eのサービスプロセッサが追加されたノードのハードウェアの構成の読み取りを行う(ステップS102)。次に、情報処理シス

50

テム 1 のオペレータは、新たなノードが有するメモリの共有領域の割り当てをサービスプロセッサに指示する (ステップ S 1 0 3 )。

【 0 1 3 4 】

次に、情報処理システム 1 のオペレータは、新たなノードのサービスプロセッサに電源投入を指示する (ステップ S 1 0 4 )。すると、各ビルディングブロック 1 0 ~ 1 0 e のサービスプロセッサは、読み取った構成の情報を元に、各ビルディングブロック 1 0 ~ 1 0 e が有する CPU のノードマップ 3 4 を I 2 C を用いて設定する (ステップ S 1 0 5 )。その後、情報処理システム 1 は、各ビルディングブロック 1 0 ~ 1 0 e の電源投入を行い (ステップ S 1 0 6 )、処理を終了する。

【 0 1 3 5 】

次に、図 1 5 を用いて、情報処理システム 1 が共有領域を制御する処理の流れについて説明する。図 1 5 は、共有領域を制御する処理の流れを説明するためのフローチャートである。まず、情報処理システム 1 は、アプリケーションの要求に応じて、ノード間の共有メモリの割当処理を実行する (ステップ S 2 0 1 )。次に、情報処理システム 1 は、ノード間で共有する共有メモリのアタッチ処理を実行する (ステップ S 2 0 2 )。

【 0 1 3 6 】

その後、情報処理システム 1 が有する各 CPU によって実行されるアプリケーションが各メモリを使用する (ステップ S 2 0 3 )。次に、情報処理システム 1 は、共有メモリのデタッチ処理を実行する (ステップ S 2 0 4 )。その後、情報処理システム 1 は、共有メモリの解放処理を実行し (ステップ S 2 0 5 )、処理を終了する。なお、ステップ S 2 0 1、および、ステップ S 2 0 5 は、その共有メモリのホームノード上のアプリケーションのみが実施してもよいし、実際の処理は、nop となるものの、その共有メモリのホームノード以外のノード上のアプリケーションも実施するものとしても良い。

【 0 1 3 7 】

次に、図 1 6 を用いて、図 1 5 中ステップ S 2 0 1 で示した共有メモリの割当処理を実行する処理の流れについて説明する。図 1 6 は、共有メモリの割当処理を説明するためのフローチャートである。図 1 6 に示す例では、例えば、CPU 2 1 が実行するアプリケーションが OS に対して、ノード間の共有メモリ割当処理の実行を要求する (ステップ S 3 0 1 )。

【 0 1 3 8 】

すると、CPU 2 1 が実行する OS が共有領域用の物理アドレスの領域から要求されたサイズのメモリ割当を行う (ステップ S 3 0 2 )。次に、OS が割り当てた共有メモリの管理用 ID をアプリケーションに引渡し (ステップ S 3 0 3 )、共有メモリの割当処理を終了する。

【 0 1 3 9 】

次に、図 1 7 を用いて、図 1 5 中ステップ S 2 0 2 で示したノード間の共有メモリアタッチ処理の流れについて説明する。図 1 7 は、共有メモリアタッチ処理を説明するためのフローチャートである。まず、アプリケーションは、OS に対して管理用 ID を引渡し、ノード間の共有メモリアタッチ処理を要求する (ステップ S 4 0 1 )。このような場合には、OS は、他のノードで実行されている OS と通信を行い、管理用 ID に対応する物理アドレスを獲得する (ステップ S 4 0 2 )。

【 0 1 4 0 】

ここで、OS が他のノードで実行されている OS と通信を行う場合には、LAN (Local Area Network) などによる通信、サービスプロセッサ 2 4 を介した各ノード間の通信等を用いる。また、例えば、各ノードで実行される OS は、特定の共有領域を、ノード間通信に用いる領域として設定し、設定した領域に対する情報の格納や読み取りを行う事で、通信を行うこととしても良い。

【 0 1 4 1 】

次に、OS は、物理アドレスに対応する論理アドレス (Virtual Address) を決定し、割当を行う (ステップ S 4 0 3 )。例えば、CPU 2 1 で実行される OS は、物理アドレ

10

20

30

40

50

スと論理アドレスとのTLBをアドレス変換部35に設定する。

【0142】

なお、各CPU21~21cが用いる論理アドレスは、重複する範囲であっても良く、また、CPUごとに異なる範囲でもよい。また、各CPU21~21cが用いる論理アドレスは、アプリケーションがOSに指定できるようにしてもよい。その後、OSは、論理アドレスの値をアプリケーションに引渡し(ステップS404)、処理を終了する。

【0143】

次に、図18を用いて、図15中ステップS203で示したアプリケーションがノード間の共有メモリを使用する処理の流れについて説明する。図18は、アプリケーションが共有メモリを使用する処理を説明するためのフローチャートである。例えば、CPU21が実行するアプリケーションは、論理アドレスを発行し、論理アドレスが示す記憶領域へのアクセスを行う(ステップS501)。

10

【0144】

すると、CPU21は、TLBミスが発生したか否かを判別する(ステップS502)。そして、CPU21は、TLBミスが発生した場合は(ステップS502肯定)、トラップ処理を実行し、TLBに論理アドレスと物理アドレスとの組のエントリを設定する(ステップS503)。

【0145】

次に、アプリケーションは、再度論理アドレスを発行し、TLBによる物理アドレスへの変換を経て、正常に共有メモリに対するアクセスを実行する(ステップS504)。一方、TLBミスが発生しなかった場合は(ステップS502否定)、正常に共有メモリに対するアクセスが実行され(ステップS505)、処理が終了する。

20

【0146】

次に、図19を用いて、図15中ステップS204で示したノード間の共有メモリデータタッチ処理の流れについて説明する。図19は、ノード間の共有メモリデータタッチ処理を説明するためのフローチャートである。例えば、CPU21が実行するアプリケーションは、OSに対して、ノード間共有メモリの論理アドレス、または管理用IDを指定して、データタッチ処理を要求する(ステップS601)。

【0147】

すると、CPU21が実行するOSは、キャッシュのフラッシュを行う(ステップS602)。すなわち、OSは、共有メモリの割り当て解除後、再度共有メモリとして割り当てを行った場合に、共有メモリとして割り当てが行われていない際に共有メモリの実メモリにアクセスするCPUがリブートすると、キャッシュと実メモリの状態が食い違う恐れがある。このため、OSは、キャッシュのフラッシュを行い、キャッシュと実メモリの状態とが食い違う状態を防止する。

30

【0148】

そして、OSは、ノード間共有メモリ、すなわち、アプリケーションが利用していた範囲の論理アドレスの割当を解除し、解除した論理アドレスに関連するTLBのエントリを削除する(ステップS603)。また、OSは、ノード間で通信を行い、本アプリケーションが対象PAの使用を完了したことを通知する(ステップS604)。そして、OSは、ノード間通信により、解放済みの共有メモリについて、最後の利用者がデータタッチを行ったことをホームノードが認識した場合、指定された共有メモリ用のメモリ割当て解除を行う(ステップS605)。なお、ステップS605の処理は、図20に示すステップS702の処理と関連する。

40

【0149】

なお、ステップS603以降は、OSは、本ノード上で、データタッチが完了しているメモリアドレスについてTLBミス(ステップS502肯定)が発生しても、データタッチが完了している論理アドレスに対応する物理アドレスをTLBに設定しない。このような場合には、ステップS504の処理は、正常に終了せず、アクセスエラーとなる。また、データタッチ完了後、ステップS402と逆に、OSがノード間で通信し、本アプリケーションがこ

50

の共有メモリのPAに対してアクセスを完了したことを通知する。もし、この共有メモリがホームノード上で解放済みで、かつ、このアプリケーションがこの共有メモリの最後の利用者であった場合は、ホームノードに解放処理を依頼する。

**【0150】**

次に、図20を用いて、図15中ステップS205で示したノード間共有メモリの解放処理の流れについて説明する。図20は、ノード間共有メモリの解放処理を説明するためのフローチャートである。例えば、CPU21が実行するアプリケーションは、OSに対してノード間共有メモリの解放処理を要求する(ステップS701)。すると、OSは、指定された共有領域の利用者が全てでタッチしていた場合は、割当てを解放し(ステップS702)、処理を終了する。もし、デタッチが完了していなければ、割当ての解放処理は行わず、処理を終了する。なお、実際の割当ての完了処理は、ステップS605で行われる。

10

**【0151】**

次に、図21を用いて、CPU21が他のCPUに対して、メモリアクセスのリクエストを送信する処理の流れについて説明する。図21は、リクエストを発行する処理の流れを説明するためのフローチャートである。例えば、CPU21の演算部は、論理アドレスを発行する(ステップS801)。

**【0152】**

すると、アドレス変換部35において、論理アドレスから物理アドレスへの変換が行われる(ステップS802)。次に、キャッシュディレクトリ管理部36が、物理アドレスを取得し、キャッシュディレクトリ管理を実行する(ステップS803)。すなわち、キャッシュディレクトリ管理部36は、取得した物理アドレスが示す記憶領域についてのキャッシュステータスを遷移させる。

20

**【0153】**

次に、キャッシュディレクトリ管理部36は、ノードマップ34を参照し、取得した物理アドレスが他ノードのメモリに振り分けられた物理アドレスであるか否かを判別する(ステップS804)。そして、キャッシュディレクトリ管理部36は、取得した物理アドレスが他ノードのメモリに振り分けられた物理アドレスではないと判別した場合には(ステップS804否定)、取得した物理アドレスを用いてメモリアクセスを実行する(ステップS805)。

30

**【0154】**

一方、キャッシュディレクトリ管理部36は、取得した物理アドレスが他ノードのメモリに振り分けられた物理アドレスである場合には(ステップS804肯定)、ノードマップ34から物理アドレスと対応付けられたCPUIDを取得する(ステップS806)。そして、パケット送信部が、CPUIDと物理アドレスとを格納したパケット、すなわち、メモリアクセスのリクエストを生成し、XB2に送出し(ステップS807)、処理を終了する。

**【0155】**

次に、図22を用いて、CPU21が他のCPUからメモリアクセスのリクエストを受信した際に行う処理の流れについて説明する。図22は、リクエストを受信した際に行う処理の流れを説明するためのフローチャートである。なお、図22に示す例では、CPU21が、他のCPUからMove Into ShareやMove Into Exclusiveを受信した際に行う処理の流れについて説明する。例えば、CPU21は、他のCPUからXB2を介してリクエストを受信する(ステップS901)。

40

**【0156】**

このような場合には、CPU21は、リクエストの対象となる物理アドレスの所定のビットが「1」であるか否かを判別することで、リクエストの対象となる物理アドレスがローカル領域であるか否かを判別する(ステップS902)。そして、CPU21は、リクエストの対象となる物理アドレスがローカル領域であると判別した場合には(ステップS902肯定)、リクエスト元のCPUに否定応答を返信し(ステップS903)、処理を

50

終了する。

【0157】

また、CPU 21は、リクエストの対象となる物理アドレスがローカル領域でない場合には（ステップS902否定）、キャッシュディレクトリ管理を実行する（ステップS904）。また、CPU 21は、物理アドレスが示す記憶領域のステータスを判定する（ステップS905）。

【0158】

そして、CPU 21は、判定したステータスに応じた命令を他のCPUに対して発行し（ステップS906）、ステータスを遷移させる（ステップS907）。その後、CPU 21は、物理アドレスが示す記憶領域のデータをリクエスト元のCPUに送信する応答を行い（ステップS908）、処理を終了する。

10

【0159】

次に、図23を用いて、CPU 21が応答を受信した際に実行する処理の流れについて説明する。図23は、CPUが応答を受信した際に実行する処理の流れを説明するためのフローチャートである。例えば、CPU 21は、応答を受信する（ステップS1001）。このような場合には、CPU 21は、応答の内容が正常な応答であるか否かを判別する（ステップS1002）。

【0160】

そして、CPU 21は、応答の内容が正常である場合、すなわち、リクエスト対象となるデータを受信した場合には（ステップS1002肯定）、データを用いた正常な処理を実行し（ステップS1003）、処理を終了する。一方、CPU 21は、否定応答を受信した場合は（ステップS1002否定）、否定応答の理由がアクセスエラーであるか否かを判別する（ステップS1004）。

20

【0161】

そして、CPU 21は、否定応答の理由がアクセスエラーではない場合には（ステップS1004否定）、通常のエラー処理を実行し（ステップS1005）、処理を終了する。一方、CPU 21は、否定応答の理由がアクセスエラーではない場合には（ステップS1004肯定）、エラーが発生した物理アドレスをエラーレジスタに設定して、トラップ処理を実行し（ステップS1006）、処理を終了する。

【0162】

[実施例1の効果]

上述したように、情報処理システム1は、CPU 21～21cとメモリ22～22cと、各CPU 21～21cを接続するXB2とを有する。また、CPU 21は、論理アドレスと物理アドレスとの変換を行うアドレス変換部と、物理アドレスとCPU IDとを変換するノードマップ34を有する。

30

【0163】

そして、CPU 21は、物理アドレスとCPU IDとを有するリクエストの packets を送信する。また、CPU 21は、他のCPUからリクエストの packets を受信した場合には、受信した packets に格納された物理アドレスに基づいて、アクセス対象となる記憶領域が、共有領域であるかローカル領域であるかを判別する。

40

【0164】

このようにすることで、情報処理システム1は、効率的、かつ小さなハードウェア物量でノード間共有メモリに対するメモリアクセスを行うことができる。すなわち、情報処理システム1は、CPU 21が物理アドレスとCPU IDとを変換するノードマップ34を用いてアドレス変換を行うので、効率的なメモリアクセスを行うことができる。

【0165】

また、CPU 21は、他のCPUがアクセスするメモリの共有領域にアクセスする場合には、物理アドレスとCPU IDを格納した packets をXB2に送出するのみでよい。このため、情報処理システム1は、効率的なメモリアクセスを行うことができる。

【0166】

50



また、情報処理システム1は、CPU21が他のCPUからリクエストの packets を受信した場合には、受信した packets に格納された物理アドレスに基づいて、アクセス対象となる記憶領域が、共有領域であるかローカル領域であるかを判別する。このため、情報処理システム1は、ローカル領域に格納するカーネルデータやユーザデータのセキュリティレベルを高く保つことができる。また、情報処理システム1は、全てのメモリをキャッシュ可能とするので、メモリアクセスにおけるレイテンシを容易に隠蔽することができる。

【0167】

また、CPU21は、他のCPUがアクセスするメモリの共有領域に対して、メモリ22にアクセスする場合と同様の方法でアクセスする。すなわち、CPU21が有する演算部31は、アクセス対象となる記憶領域がメモリ22上に存在する場合にも、他のメモリ上に存在する場合にも、論理アドレスを出力するだけでよい。

10

【0168】

このため、情報処理システム1は、I/Oの排他制御等の処理やプログラミング等を実行せずとも、容易に共有領域にアクセスできるため、メモリアクセス性能を向上させることができる。また、CPU21は、実行するプログラムやOSに改変を行わずとも、共有メモリを適切に利用することができる結果、プリフェッチ処理の実行を従来と同様に行う事ができるため、メモリアクセスの性能を向上させることができる。

【0169】

また、情報処理システム1は、所定のビットが「1」となる物理アドレスを共有領域に割り当て、所定のビットが「0」となる物理アドレスをローカル領域に割り当てる。このため、CPU21は、物理アドレスのうち、所定の1ビットが「1」であるか否かを判別するだけで、アクセス対象の物理アドレスが共有領域の物理アドレスであるか否かを容易に判別することができる。この結果、情報処理システム1は、効率的なメモリアクセスを行うことができる。

20

【0170】

また、CPU21は、他のCPUからのメモリアクセスの対象がローカル領域へのアクセスであると判定した場合には、否定応答を返信する。このため、情報処理システム1は、共有領域以外へのアクセスを防止する結果、エラーを防ぐことができる。

【0171】

また、キャッシュディレクトリ管理部36は、ノードマップ34を用いて、物理アドレスをノードマップ34に対応付けて記憶されたCPUIDに変換する。このため、CPU21は、アクセス対象となる物理アドレスが振り分けられたメモリにアクセスするCPUを識別することができる。

30

【0172】

また、各ビルディングブロック10~10eは、ノードマップ34の書き換えを行うサービスプロセッサを有する。このため、情報処理システム1は、メモリ22~22cごとに、ローカル領域と共有領域とを自由に割り当てることができる。例えば、情報処理システム1は、メモリ22が4TBの容量を有する場合に、ローカル領域に1TBを割り当て、共有領域に3TBを割り当てるというように、任意の容量の記憶領域をノード間で共有することができる。

40

【0173】

また、情報処理システム1は、新たなCPUとメモリを追加した場合やCPUやメモリの削除を行った場合にも、サービスプロセッサを介して容易にローカル領域と共有領域との割り当てを行うことができる。

【0174】

また、CPU21は、メモリ22に記憶されたデータをキャッシュしたCPUを管理するディレクトリを用いて、キャッシュコヒーレンスの制御を行う。このため、情報処理システム1は、情報処理システム1が有するCPUの数が増加した場合にも、XB2のトラフィックを増加させることなく、効率的にキャッシュコヒーレンスを保持することができ

50

る。

【0175】

具体的には、情報処理システム1においては、各CPU間の通信が、リモートCPUとホームCPU間、または、リモートCPUとホームCPUと更新したデータをキャッシュするローカルCPU間に限定される。このため、情報処理システム1は、効率的にキャッシュコヒーレンスを保持することができる。

【0176】

また、CPU21は、キャッシュミスが発生した場合に、キャッシュミスした物理アドレスが他のCPUがアクセスするメモリに割り当てられた物理アドレスであるか否かを判別する。そして、CPU21は、キャッシュミスした物理アドレスが他のCPUがアクセスするメモリに割り当てられた物理アドレスであると判別した場合には、物理アドレスをCPUIDに変換し、物理アドレスとCPUIDとを格納したパケットの生成および送出行う。このため、CPU21は、無駄なアドレス変換処理を行うことなくメモリアクセスを行うことができる。

10

【0177】

また、CPU21は、実行するアプリケーションが共有領域の獲得を要求した場合には、アプリケーションが利用する論理アドレスと、共有領域に割り当てられる物理アドレスとを変換するTLBを設定する。このため、CPU21は、実行するアプリケーションやOSに共有領域やローカル領域へのアクセスを意識した改変を加えずとも、メモリアクセスを行うことができる。

20

【実施例2】

【0178】

これまで本発明の実施例について説明したが実施例は、上述した実施例以外にも様々な異なる形態にて実施されてよいものである。そこで、以下では実施例2として本発明に含まれる他の実施例を説明する。

【0179】

(1)ビルディングブロックについて

上述した情報処理システム1は、4つのCPUを有するビルディングブロック10~10eを有していた。しかし、実施例はこれに限定されるものではなく、ビルディングブロック10~10eは、任意の数のCPU及び各CPUがアクセスするメモリを有することができる。また、CPUとメモリは、1対1で対応している必要はなく、メモリに直接アクセスするCPUは全体の一部であってもよい。

30

【0180】

(2)共有領域とローカル領域の割り当てについて

上述した共有領域とローカル領域に対する物理アドレスの割り当ては、あくまで一例であり、情報処理システム1は、任意の物理アドレスを各領域に割当てることができる。

【0181】

例えば、情報処理システム1は、物理アドレスの最下位1ビットが「0」となる物理アドレスを共有領域に割り当て、物理アドレスの最下位1ビットが「1」となる物理アドレスを共有領域に割当てることとしてもよい。このような場合には、各CPUは、物理アドレスの最下位1ビットが「0」であるか「1」であるかを判別することで、アクセス対象が共有領域であるかを容易に判別できる。

40

【0182】

また、情報処理システム1は、物理アドレス空間の前半に含まれる任意の物理アドレスを共有領域に割り当て、物理アドレス空間の後半に含まれる任意の物理アドレスをローカル領域に割当ててもよい。このような場合には、各CPUは、物理アドレスの最上位1ビットが「0」であるか「1」であるかを判別することで、アクセス対象が共有領域であるかを容易に判別できる。なお、情報処理システム1は、物理アドレス空間の前半に含まれる任意の物理アドレスをローカル領域に割当て、後半に含まれる任意の物理アドレスを共有領域に割当ててもよい。

50

## 【0183】

すなわち、情報処理システム1は、任意の物理アドレスを共有領域とローカル領域に割り当てることができるが、所定のビットが同一の値となる物理アドレスを共有領域に割り当て、所定のビットが共有領域とは異なる値の物理アドレスをローカル領域に割り当てることで、アクセス対象が共有領域であるかローカル領域であるかを容易に判別できる。

## 【0184】

(3) CPUが送信するパケットについて

上述したCPU21は、CPUIDとPAとを有するパケットをメモリアクセスのリクエストとして送信した。しかし、実施例は、これに限定されるものではない。すなわち、CPU21は、アクセス対象となるメモリにアクセスするCPUを一意に識別できるのであれば、任意の情報を格納したパケットを出力してよい。

10

## 【0185】

また例えば、CPU21は、CPUIDからVC(Virtual Connection)IDに変換し、VCIDを格納することとしてもよい。また、CPU21は、パケットに、データ長を示すレングス等の情報を格納することとしてもよい。

## 【0186】

(4) CPUが発行する命令について

上述したように、各CPU21~21cは、リクエストや命令を発行して、キャッシュのコヒーレンスを保持した。しかし、上述したリクエストや命令は、あくまで一例であり、例えばCPU21~21cは、CAS(Compare And Swap)命令を発行してもよい。

20

## 【0187】

このように、CPU21~21cがCAS命令を発行した場合には、排他制御のコンテンツンションが複数のCPU間で頻発しても、各CPUのキャッシュ上で処理が行われる。この結果、CPU21~21cは、メモリアクセスの発生による遅延を防止するとともに、各CPU間のドランザクションが混雑するのを防ぐことができる。

## 【0188】

(5) ハイパーバイザを経由した制御について

上述した情報処理システム1では、OSによってハードウェアであるアドレス変換部35にアクセスを行う例について説明した。しかし、実施例はこれに限定されるものではなく、たとえば、仮想マシンを動作させるハイパーバイザ(HPV: Hypervisor)がアドレス変換部35にアクセスを行っても良い。

30

## 【0189】

すなわち、ハイパーバイザが動作するノードにおいては、OSは、キャッシュやMMUなどのCPU21~21cのハードウェア資源に対して直接の操作を行わず、操作をハイパーバイザに依頼することとなる。このように、各CPU21~21cは、ハイパーバイザを介した制御を受付ける場合には、仮想アドレスを実アドレス(RA: Real Address)に変換し、その後、実アドレスを物理アドレスに変換することとなる。

## 【0190】

また、ハイパーバイザが動作するノードにおいては、割り込み処理は、OSには直接割り込まず、HPVに対して割り込みを行う。このような場合には、ハイパーバイザが、OSの割り込み処理ハンドラを読み出すことで割り込みを行う。なお、上述したハイパーバイザが実行する処理は、仮想マシンを動作させるために実行される公知な処理である。

40

## 【0191】

(6) パーティションを用いた処理について

上述した情報処理システム1では、各CPU21~21cは、1つのノードマップを用いてメモリアクセスを送信していた。しかし、実施例はこれに限定されるものではない。例えば、各ビルディングブロック10~10eは、複数のノード群として動作し、各ノード群毎に、同一のファームウェア(ハイパーバイザ)を動作させる1つの論理パーティションを構成しても良い。

50

## 【 0 1 9 2 】

このような場合には、各 CPU 2 1 ~ 2 1 c は、アクセス先の CPU を示すノードマップと、同一論理パーティション内の CPU を示すノードマップとを有する。このように、各 CPU 2 1 ~ 2 1 c は、同一論理パーティション内に含まれる CPU を示すノードマップを有することで、エラー発生通知、ダウン要求、リセット要求パケット等の、論理パーティションを超えて転送すべきではない特殊パケットの転送範囲を識別することができる。

## 【 0 1 9 3 】

以下、同一論理パーティション内に含まれる CPU を示すノードマップを有する CPU について説明する。図 2 4 は、実施例 2 に係る情報処理システムを説明するための図である。図 2 4 に示すように、ビルディングブロック 1 0、1 0 a は、論理パーティション # A を動作させ、ビルディングブロック 1 0 b ~ 1 0 d は、論理パーティション # B を動作させる。

10

## 【 0 1 9 4 】

ここで、論理パーティション # A では、複数のドメイン # A ~ # C と、ファームウェア # A が動作する。また、論理パーティション # B では、複数のドメイン # D ~ # G とファームウェア # B が動作する。なお、ファームウェア # A およびファームウェア # B とは、例えばハイパーバイザである。また、ドメイン # A では、アプリケーションと OS とが動作しており、他のドメイン # B ~ # G もドメイン # A と同様に、アプリケーションと OS とが動作する。

20

## 【 0 1 9 5 】

つまり、各ドメイン # A ~ # G は、それぞれ独立してアプリケーションと OS が動作する仮想マシンである。ここで、ビルディングブロック 1 0 が有する各 CPU 2 1 ~ 2 1 c は、パーティション # A に含まれる各 CPU に対して上述した特殊パケットを送信してもよいが、パーティション # B に含まれる各 CPU に対しては特殊パケットを送信すべきではない。

## 【 0 1 9 6 】

このため、各ビルディングブロック 1 0 ~ 1 0 d の CPU は、同一の論理パーティションに含まれる CPU の CPU ID を示すノードマップを有する。例えば、CPU 2 1 は、物理アドレスと、物理アドレスが示す記憶領域を有するメモリと接続された CPU の CPU ID とを対応付けて記憶するノードマップ 3 4 を有する。また、CPU 2 1 は、CPU 2 1 と同一のパーティション、すなわち、パーティション # A に含まれる CPU の CPU ID を記憶するノードマップ 3 4 a を有する。なお、ノードマップ 3 4 a は、ノードマップ 3 4 と同様に、サービスプロセッサ 2 4 によって設定されるものとする。

30

## 【 0 1 9 7 】

以下、図面を用いて、同一の論理パーティションに含まれる CPU の CPU ID を示すノードマップの一例について説明する。図 2 5 は、パーティションの一例を説明するための図である。例えば、図 2 5 に示す例では、パーティション # A は、ビルディングブロック # 0 を有する。また、ビルディングブロック # 0 は、CPU # 0 とアドレス域「 # 0 」が割当てられたメモリとを有する。

40

## 【 0 1 9 8 】

また、パーティション # B は、ビルディングブロック # 1 とビルディングブロック # 2 とを有する。また、ビルディングブロック # 1 は、CPU # 4、CPU # 5、アドレス域「 # 1 」が割当てられたメモリ、アドレス域「 # 2 」が割当てられたメモリを有する。なお、アドレス域「 # 1 」が割当てられたメモリには、CPU # 4 がアクセスし、アドレス域「 # 2 」が割当てられたメモリには、CPU # 5 がアクセスする。また、ビルディングブロック # 2 は、CPU # 8 とアドレス域「 # 3 」が割当てられたメモリを有する。

## 【 0 1 9 9 】

次に、図 2 6 A ~ 2 6 C を用いて、図 2 5 に示す CPU # 0 が有するノードマップと、CPU # 4 が有するノードマップとについて説明する。まず、図 2 6 A および図 2 6 B を

50

用いて、パーティション# AのCPUが記憶するノードマップについて説明する。なお、図26Aは、パーティション# AのCPUが記憶するノードマップの一例を説明するための図である。また、図26Bは、パーティション# Aを示すノードマップの一例を説明するための図である。

【0200】

なお、以下の説明では、ノードID「0」は、ビルディングブロック# 0を示し、ノードID「1」は、ビルディングブロック# 1を示す、ノードID「2」は、ビルディングブロック# 2を示す。また、CPUID「0」は、CPU# 0のCPUIDであり、CPUID「4」は、CPU# 4のCPUIDであり、CPUID「5」は、CPU# 5のCPUIDであり、CPUID「8」は、CPU# 8のCPUIDであるものとする。

10

【0201】

例えば、図26Aに示す例では、ノードマップ34は、アドレス域「# 0」が、ビルディングブロック# 0に存在し、CPU# 0がアクセスを行う旨を示す。また、ノードマップ34は、アドレス域「# 1」が、ビルディングブロック# 1に存在し、CPU# 4がアクセスを行う旨を示す。また、ノードマップ34は、アドレス域「# 2」が、ビルディングブロック# 1に存在し、CPU# 5がアクセスを行う旨を示す。また、ノードマップ34は、アドレス域「# 3」がビルディングブロック# 2に存在し、CPU# 8がアクセスを行う旨を示す。

【0202】

また、図26Bには、パーティション# Aを示すノードマップを示した。図26Bに示すように、パーティション# Aを示すノードマップは、各エントリに、バリッドと、ノードIDとCPUIDとを有する。例えば、図26Bに示す例では、ノードマップは、パーティション# Aにビルディングブロック# 0のCPU# 0が含まれる旨を示す。

20

【0203】

例えば、図25に示す例では、CPU# 0は、図26Aおよび図26Bに示すノードマップを有する。そして、CPU# 0は、メモリアクセスを行う場合には、図26Aに示すノードマップを用いて、アクセス先のCPUを識別する。一方、CPU# 0は、同一パーティション内のCPUのみに特殊パケットを送信する場合には、図26Bに示すノードマップを用いて、送信先のCPUを識別する。すなわち、CPU# 0は、図26Bに例示するノードマップが示すパーティション# A内のCPUに対して、特殊パケットを送信する。

30

【0204】

一方、CPU# 4は、メモリアクセスを行うために、図26Aに示すノードマップと、図26Cに示すノードマップとを有する。ここで、図26Cは、パーティション# Bを示すノードマップの一例を説明するための図である。図26Cに示す例では、パーティション# Bを示すノードマップは、パーティション# Bに、ビルディングブロック# 1のCPU# 4およびCPU# 5、ビルディングブロック# 2のCPU# 38が存在することを示す。CPU# 4は、図26Cに例示するノードマップが示すパーティション# B内のCPUに対して、特殊パケットを送信する。

【0205】

40

このように、CPU# 1およびCPU# 4は、アドレス域とCPUIDとを対応付けたノードマップと、パーティションを示すノードマップとを記憶する。そして、CPU# 1およびCPU# 4は、アドレス域とCPUIDとを対応付けたノードマップを用いて、他のノードが有するメモリに対して直接メモリアクセスを行う。また、CPU# 1は、パーティション# Aを示すノードマップを用いて、特殊パケットの送信を行う。また、CPU# 4は、パーティション# Bを示すノードマップを用いて、特殊パケットの送信を行う。

【0206】

このように、各CPUは、自身を含むパーティションごとに、異なる値を有するノードマップを有してもよい。また、各CPUは、自身を含むパーティションごとに異なる値を有するノードマップを有する場合は、特殊パケットをパーティションを超えた送信を行う

50

ことを防ぐことができる。

【0207】

なお、各CPUは、実施例1と同様、スタートアドレスとアドレスマスク、又は、スタートアドレスとレンジでアクセス対象となるアドレス域を示しても良い。すなわち、CPU#1とCPU#4とは、スタートアドレスとアドレスマスク、又は、スタートアドレスとレンジとを用いて、アクセス対象となるアドレス域を示すノードマップを用いて、アクセス対象となるノードを識別する。また、CPU#1とCPU#4とは、それぞれ異なるパーティションを示すノードマップを用いて、特殊パケットの送信を行う。

【符号の説明】

【0208】

- 1 情報処理システム
- 2 X B
- 10 ~ 10 e ビルディングブロック
- 20 ノード
- 21 ~ 21 c CPU
- 22 ~ 22 c メモリ
- 23、26 通信部
- 24 サービスプロセッサ
- 25 制御部
- 27、27 a X B 接続部
- 28 P C I e 接続部
- 30 演算処理部
- 31 演算部
- 32 L 1 キャッシュ
- 33 L 2 キャッシュ
- 34 ノードマップ
- 35 アドレス変換部
- 36 キャッシュディレクトリ管理部
- 37 パケット制御部
- 38 リクエスト生成部
- 39 リクエスト受信部
- 40 ルータ
- 41 メモリアクセス部
- 42 P C I e 制御部
- 43 リクエスト生成部
- 44 P C I e バス制御部

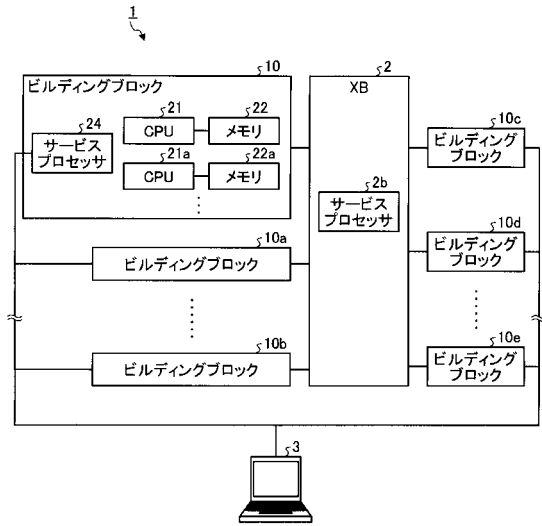
10

20

30

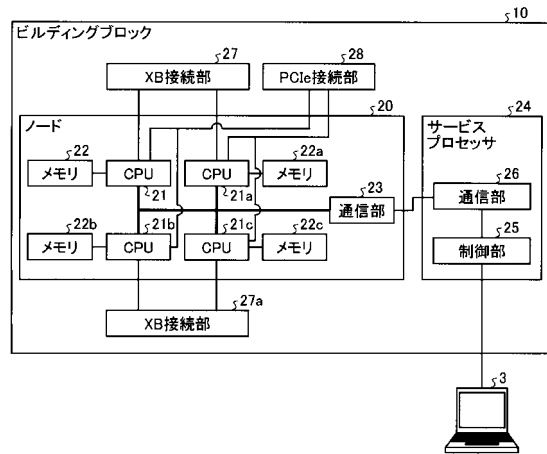
【 図 1 】

実施例1に係る情報処理システムの一例を説明するための図



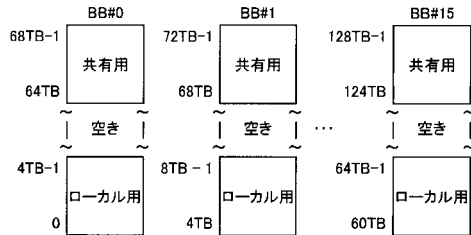
【 図 2 】

実施例1に係るビルディングブロックの機能構成を説明するための図



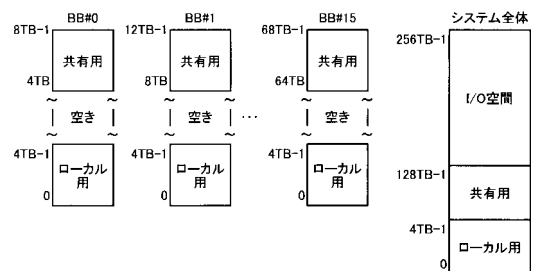
【 図 3 】

実施例1に係るビルディングブロックのメモリに振り分けられる物理アドレスの範囲を説明するための図



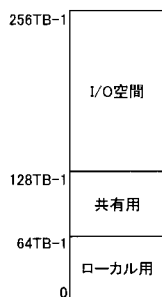
【 図 5 】

物理アドレスの振り分けのバリエーションを説明するための第1の図



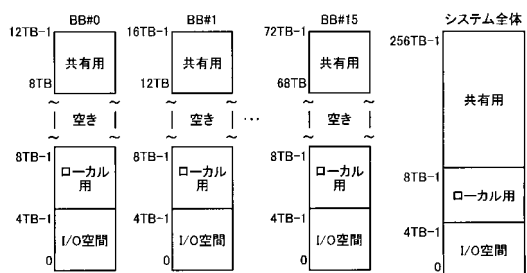
【 図 4 】

実施例1に係る情報処理システムが各メモリに振り分ける物理アドレスを説明するための図

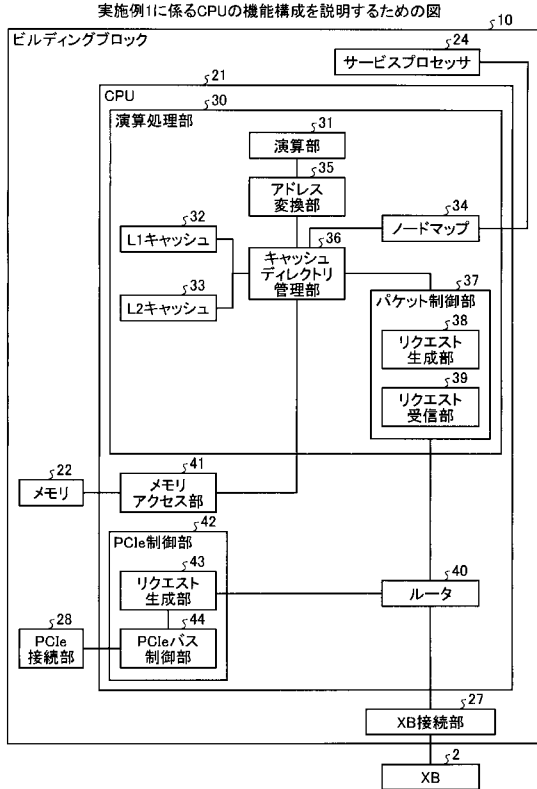


【 図 6 】

物理アドレスの振り分けのバリエーションを説明するための第2の図



【図7】



【図8】

実施例1に係るノードマップが記憶する情報の一例を説明するための図

アドレス	バリッド	ノードID	CPUID
#0	1	0	0
#1	1	0	1
#2	0		

【図9】

ノードマップが記憶する情報のバリエーションの一例を説明するための第1の図

バリッド	スタートアドレス	アドレスマスク	ノードID	CPUID
1	0x00000	0x3fff	0	0
1	0x10000	0x3fff	1	4
1	0x14000	0x3fff	1	5
1	0x20000	0x1fff	2	8

【図10】

ノードマップが記憶する情報のバリエーションの一例を説明するための第2の図

バリッド	スタートアドレス	レンジ	ノードID	CPUID
1	0x00000	0x3fff	0	0
1	0x10000	0x3fff	1	4
1	0x14000	0x3fff	1	5
1	0x20000	0x2ef	2	8

【図11A】

キャッシュタグの一例を説明するための図

縮退フラグ	ECCチェックビット	IF/オブコート	L1 キャッシュステート	L2 キャッシュステート	AA
-------	------------	----------	--------------	--------------	----

【図11B】

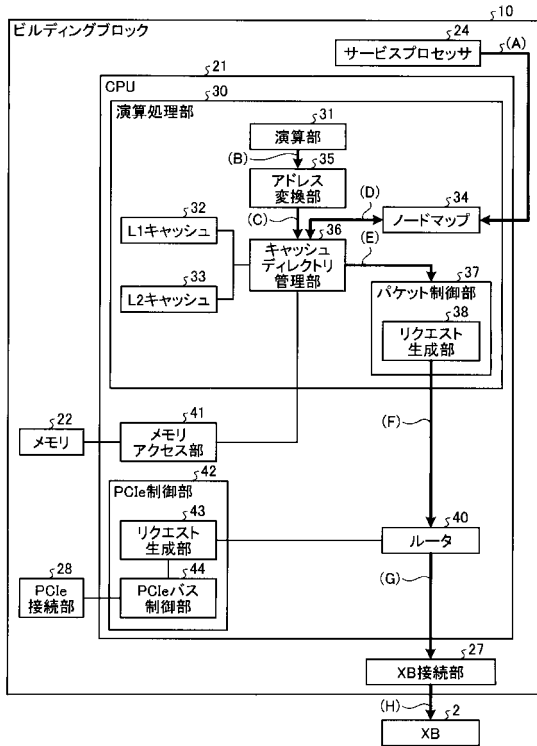
実施例1に係るCPUが送信するパケットを説明するための図

CPUID	PA	データ
-------	----	-----



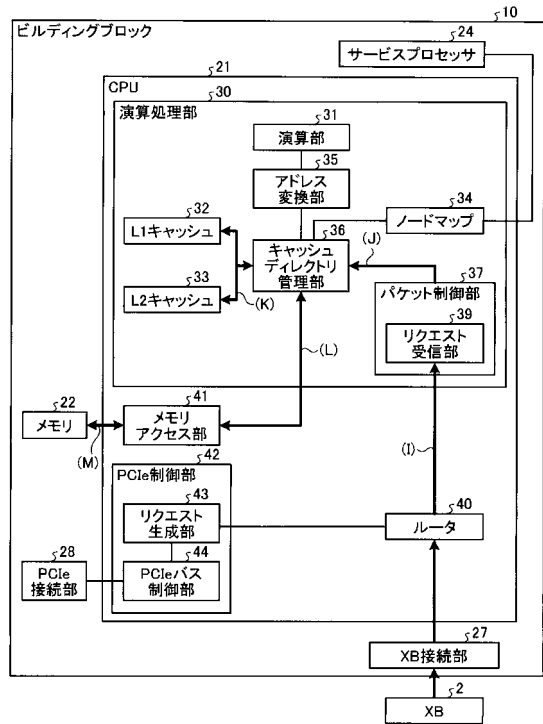
【図12】

実施例1に係るCPUがリクエストを送信する処理の一例を説明するための図



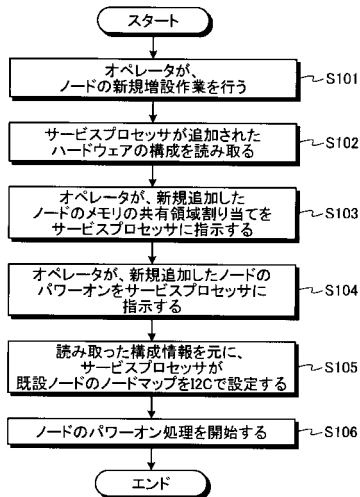
【図13】

実施例1に係るCPUがパケットを受信した際に行う処理の一例を説明するための図



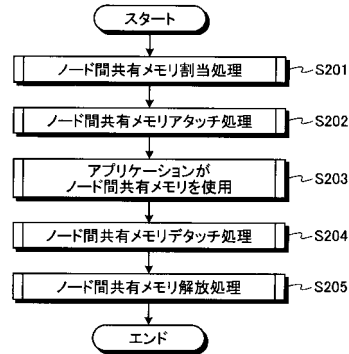
【図14】

ノードマップを設定する処理の流れを説明するためのフローチャート



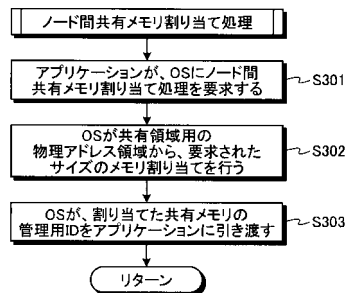
【図15】

共有領域を制御する処理の流れを説明するためのフローチャート



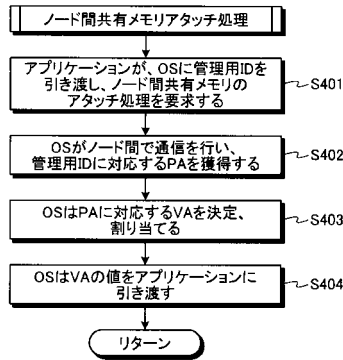
【図16】

共有メモリの割当処理を説明するためのフローチャート



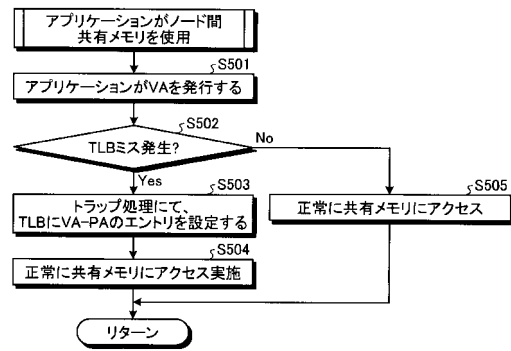
【図17】

共有メモリアタッチ処理を説明するためのフローチャート



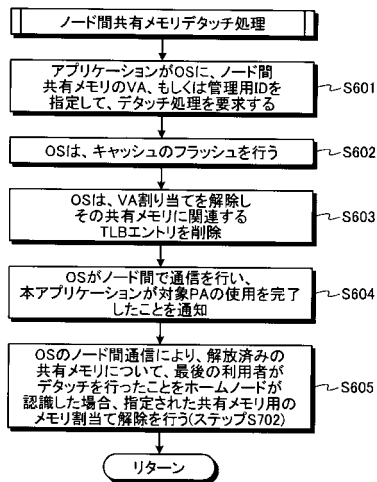
【図18】

アプリケーションが共有メモリを使用する処理を説明するためのフローチャート



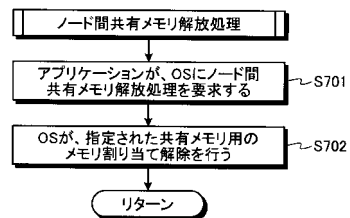
【図19】

ノード間の共有メモリデタッチ処理を説明するためのフローチャート



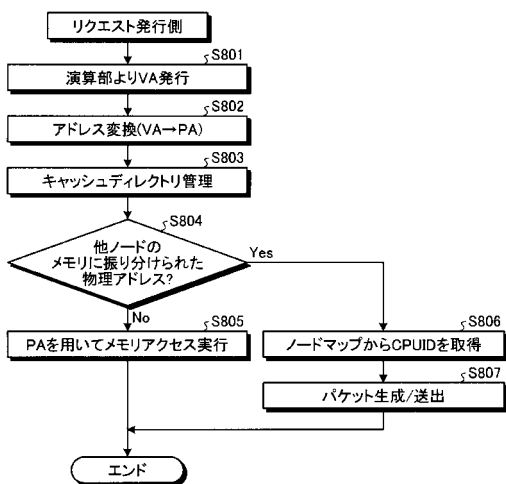
【図20】

ノード間共有メモリの解放処理を説明するためのフローチャート



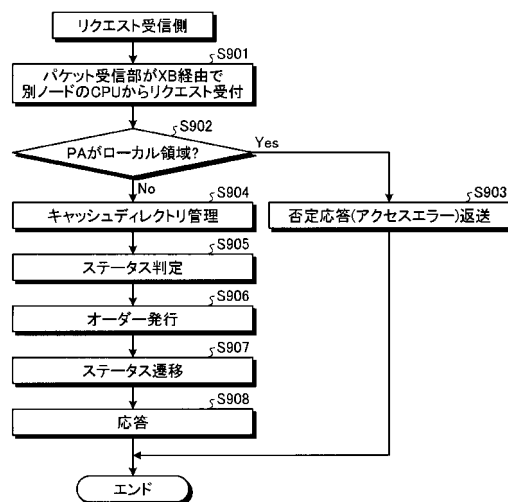
【図 2 1】

リクエストを発行する処理の流れを説明するためのフローチャート



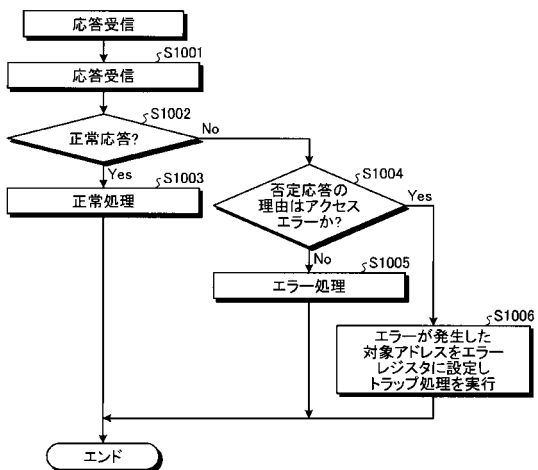
【図 2 2】

リクエストを受信した際に実行する処理の流れを説明するためのフローチャート



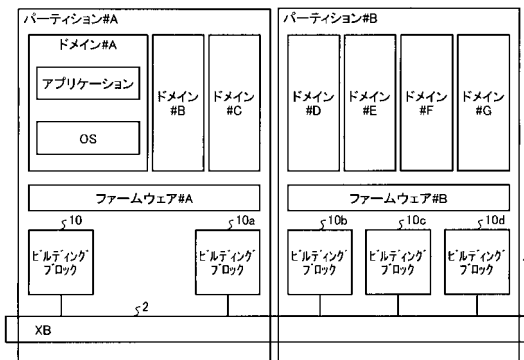
【図 2 3】

CPUが応答を受信した際に実行する処理の流れを説明するためのフローチャート



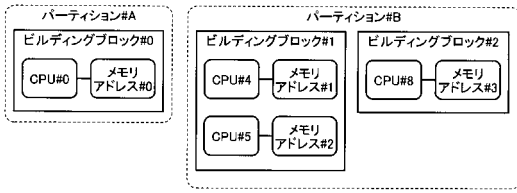
【図 2 4】

実施例2に係る情報処理システムを説明するための図



【図 25】

パーティションの一例を説明するための図



【図 26 A】

パーティション#AのCPUが記憶するノードマップの一例を説明するための図

アドレス	バリッド	ノードID	CPUID
#0	1	0	0
#1	1	1	4
#2	1	1	5
#3	1	2	8

⋮

【図 26 B】

パーティション#Aを示すノードマップの一例を説明するための図

エントリー	バリッド	ノードID	CPUID
#0	1	0	0
#1	0		
#2	0		
#3	0		

⋮

【図 26 C】

パーティション#Bを示すノードマップの一例を説明するための図

エントリー	バリッド	ノードID	CPUID
#0	1	1	4
#1	1	1	5
#2	1	2	8
#3	0		

⋮

---

フロントページの続き

審査官 野田 佳邦

- (56)参考文献 特開平10-240707(JP,A)  
特開2000-067009(JP,A)  
特開平02-244253(JP,A)  
特開平08-030568(JP,A)  
特開2007-199999(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F 12/00 - 12/06  
G06F 12/08 - 12/12