

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2010-165022

(P2010-165022A)

(43) 公開日 平成22年7月29日 (2010.7.29)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 15/167 (2006.01)	G06F 15/167 610A	5B045
G06F 9/54 (2006.01)	G06F 9/46 480A	
G06F 15/17 (2006.01)	G06F 15/17	

審査請求 未請求 請求項の数 10 O L (全 24 頁)

(21) 出願番号 特願2009-4678 (P2009-4678)
 (22) 出願日 平成21年1月13日 (2009.1.13)

(71) 出願人 000006747
 株式会社リコー
 東京都大田区中馬込 1 丁目 3 番 6 号
 (74) 代理人 100110607
 弁理士 間山 進也
 (72) 発明者 本橋 弘臣
 東京都大田区中馬込 1 丁目 3 番 6 号 株式
 会社リコー内
 Fターム(参考) 5B045 BB02 BB12 BB28 BB29 BB32
 DD02

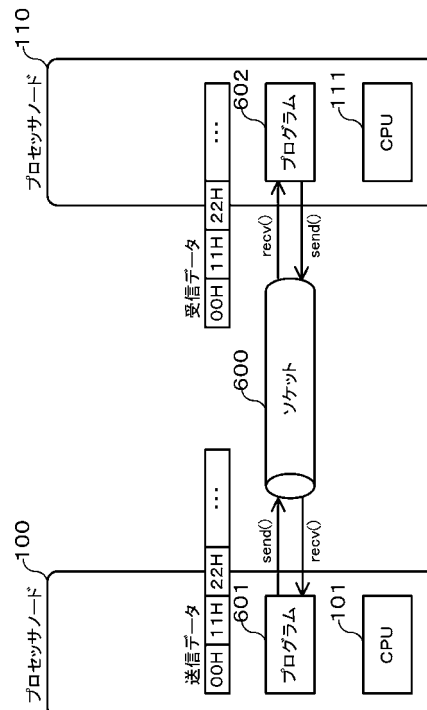
(54) 【発明の名称】 プロセッサ間通信装置、プロセッサ間通信方法、プログラムおよび記録媒体

(57) 【要約】

【課題】 各プロセッサノード間での共有メモリを用いた通信処理で発生する通信オーバーヘッドを低減し、高速通信が可能なプロセッサ間通信装置、プロセッサ間通信方法、プログラムおよび記録媒体を提供すること。

【解決手段】 本発明に係るプロセッサ間通信装置は、第1および第2のプロセッサノードからアクセス可能に共有した共有メモリを備え、第1のプロセッサノードでのプログラムに従って、通信データを共有メモリに書き込む書込手段と、第1および第2のプロセッサノードを接続するソケットと、書込通知情報をソケットを介して第2のプロセッサノードに送信する第1の通知手段と、書込通知情報に応じて共有メモリから通信データを読み出す読出手段と、読出通知情報をソケットを介して第1のプロセッサノードに送信する第2の通知手段とを備えている。

【選択図】 図6



【特許請求の範囲】**【請求項 1】**

第 1 および第 2 のプロセッサノードから構成され、前記第 1 および第 2 のプロセッサノードのそれぞれからアクセスすることが可能に共有して設けられた共有メモリを備えたプロセッサ間通信装置において、

前記第 1 のプロセッサノードで動作しているプログラムに従って、通信データを前記共有メモリに書き込む書込手段と、

前記書込手段により前記通信データが前記共有メモリに書き込まれたことを通知する書込通知情報を、ソケット通信により前記第 2 のプロセッサノードに送信する第 1 の通知手段と、

前記第 1 の通知手段により送信された前記書込通知情報に応じて、前記共有メモリから通信データを読み出す読出手段と、

前記読出手段により前記通信データが読み出されたことを通知する読出通知情報を、前記ソケット通信により前記第 1 のプロセッサノードに送信する第 2 の通知手段とを備える、プロセッサ間通信装置。

【請求項 2】

前記共有メモリは、

前記通信データに対して、循環的に書き込み処理を行うことが可能なリングバッファと、

前記リングバッファ内での通信データを書き込む際の開始位置を示すライトポイントと

前記リングバッファ内での通信データを読み出す際の開始位置を示すリードポイントを含む、請求項 1 に記載のプロセッサ間通信装置。

【請求項 3】

前記共有メモリは、前記第 1 および第 2 のプロセッサノード内にそれぞれ設けられたメインメモリの記憶領域の一部を用いて成る、請求項 1 または 2 に記載のプロセッサ間通信装置。

【請求項 4】

前記書込手段により前記共有メモリに書き込まれた前記通信データのデータ量が、所定量に達した場合に、前記通信データが書き込まれたことを通知するデータ量通知情報を、前記ソケット通信を介して前記第 2 のプロセッサノードに送信するデータ量通知手段を備える、請求項 1 ~ 3 に記載のプロセッサ間通信装置。

【請求項 5】

前記第 1 の通知手段が、前記ソケット通信により前記第 2 のプロセッサノードに前記書込通知情報を送信するための条件を設定する条件設定手段を備える、請求項 1 ~ 4 に記載のプロセッサ間通信装置。

【請求項 6】

前記共有メモリ内のうち空き領域を検出して、これらの検出した各領域から、前記書込手段が通信データを書き込むための領域を確保して割り当てる割り当て手段を備える、請求項 1 ~ 5 に記載のプロセッサ間通信装置。

【請求項 7】

前記割り当て手段により割り当てられた、前記書込手段が通信データを書き込むための領域のアドレスを通知するための割り当て通知情報を、前記ソケット通信を介して前記第 2 のプロセッサノードに送信する送信手段を備える、請求項 1 ~ 6 に記載のプロセッサ間通信装置。

【請求項 8】

プロセッサ間通信装置が実行するプロセッサ間通信方法であって、前記プロセッサ間通信方法は、プロセッサ間通信装置が、

前記プロセッサ間通信装置に構成されている第 1 および第 2 のプロセッサノードのうち、前記第 1 のプロセッサノードで動作しているプログラムに従って、前記第 1 および第 2

10

20

30

40

50

のプロセッサノードのそれぞれからアクセスすることが可能に共有して設けられた共有メモリに通信データを書き込むステップと、

前記通信データが前記共有メモリに書き込まれたことを通知する書込通知情報を、前記第1および第2のプロセッサノードを、前記プログラム間で通信データの送受信が可能に接続するソケット通信を介して、前記第2のプロセッサノードに送信するステップと、

前記送信された前記書込通知情報に応じて、前記共有メモリから通信データを読み出すステップと、

前記通信データが読み出されたことを通知する読出通知情報を、前記ソケット通信を介して前記第1のプロセッサノードに送信するステップと

を備える、プロセッサ間通信方法。

10

【請求項9】

プロセッサ間通信装置がプロセッサ間通信を行うための装置実行可能なプログラムであって、前記プログラムは、プロセッサ間通信装置を、

前記プロセッサ間通信装置に構成されている第1および第2のプロセッサノードのうち、前記第1のプロセッサノードで動作している動作プログラムに従って、前記第1および第2のプロセッサノードのそれぞれからアクセスすることが可能に共有して設けられた共有メモリに通信データを書き込む書込手段と、

前記書込手段により前記通信データが前記共有メモリに書き込まれたことを通知する書込通知情報を、前記第1および第2のプロセッサノードを前記動作プログラム間で通信データの送受信が可能に接続するソケット通信により、前記第2のプロセッサノードに送信する第1の通知手段と、

20

前記第1の通知手段により送信された前記書込通知情報に応じて、前記共有メモリから通信データを読み出す読出手段と、

前記読出手段により前記通信データが読み出されたことを通知する読出通知情報を、前記ソケット通信を介して前記第1のプロセッサノードに送信する第2の通知手段として機能させる、装置実行可能なプログラム。

【請求項10】

請求項9に記載のプログラムを格納した情報処理装置可読な記録媒体。

【発明の詳細な説明】

【技術分野】

30

【0001】

本発明は、複数のプロセッサノードから構成され、各プロセッサノードのそれぞれからアクセスすることが可能に共有して設けられた共有メモリを備えたプロセッサ間通信装置、プロセッサ間通信方法、プログラムおよび記録媒体に関する。

【背景技術】

【0002】

従来より、複数のプロセッサノードを備えたプロセッサ間通信装置の1例である疎結合型のマルチプロセッサシステムの動作効率は、各プロセッサノード間の通信処理で発生する通信オーバーヘッドからのシステムパフォーマンスへの影響をどのようにして低減させるかに左右される。それは、マルチプロセッサシステムが実行する各処理のジョブのうち、プロセッサノードの負荷が大きいジョブが存在すると、このジョブを複数の各処理毎に細分化して各プロセッサノードにそれぞれ割り当てたととしても、大抵はそれぞれの細分化されたタスクなどの各処理の間で何らかの関連性や依存性があるために、各段階の処理を経ていくに従って各プロセッサノード間で演算結果や計算結果などの情報を相互に受け渡す通信処理を行うからである。

40

【0003】

このような各プロセッサノード間での通信処理では、従来ではEthernet（登録商標）やInfiBand、Myrinetなどのインタフェースが利用されてきた。しかし、Ethernet（登録商標）には、送信元のプロセッサノードが動作している処理でデータを送信してから、受信側のプロセッサノードが動作している処理でデータを

50

受け取るまでの時間、即ち、通信時のレイテンシが長い、また、TCP/IPなどのプロトコル処理が重いという欠点があった。

【0004】

通信時のレイテンシが長いと、プロセッサノード間で頻繁にデータを相互に受け渡す場合には、通信オーバーヘッドが増大してシステム全体のパフォーマンス、動作効率が低下する。また、プロトコル処理が重いと、貴重なCPU性能が本来の演算処理や計算処理などの本来の目的以外の処理で浪費されてしまう。

【0005】

一方、InfiniBandやMyrinetは、レイテンシが短く、プロトコル処理がハードウェア化されているためにCPUの負担が軽いという利点があるが、これらのインタフェースカードは、Ethernet（登録商標）と比べると高機能・高性能であるがゆえに非常に高価であり、組み込み機器ではローコストの実現が困難であった。

10

【0006】

また、マルチプロセッサシステムで用いられている組み込み機器では、ローコストが要求されていることもありCPUの個数が多くても高々2～4個程度に留まっていることから、各プロセッサノード間での通信処理は共有メモリを利用して行われる場合が多かった。このような場合の共有メモリには、従来より、デュアルポートメモリとして参照されるメモリ素子が使用されていたが、このデュアルポートメモリの機能では、1個に対して2個のCPUまでに限って接続可能であるという制限があった。

【0007】

最近ではデュアルポートメモリ以外にも、2個のプロセッサノードのPCIバス同士を接続するPCI-PCIブリッジや、2個以上の複数のプロセッサノードのPCI-Expressバスをクロスバーで接続するPCI-Expressスイッチを用いたLSIが使用されている。

20

【0008】

これらのPCI-PCIブリッジやPCI-Expressスイッチでは、通常1つのプロセッサノードが占有利用するメインメモリに対して、他のプロセッサノードがダイレクトにアクセスすることができるため、ローコストで容易に共有メモリの機能を有することができる。また、これらのPCI-PCIブリッジやPCI-Expressスイッチは、2個以上のCPUを備えたマルチプロセッサシステムで利用されている。

30

【0009】

特許第3743381号公報（特許文献1）に記載されたコンピュータシステムでは、複数のホストと、これらのホストのいずれからもアクセス可能な共有メモリと、ホスト間をつなぐ通信経路とを備えるコンピュータシステムにおいて、自ホスト内の、ホスト間で連携して動作するプログラムから発行されるホスト間処理同期要求を受理し、その旨を示す情報を共有メモリに記録し、前記旨を前記通信経路を介して他ホストに通知し、自ホストのホスト間処理同期要求の受理や他ホストから通知される、ホスト間処理同期要求の受理通知を契機とし、共有メモリを参照してホスト間処理同期成立の判断・認識を行う各ホスト内のホスト間処理同期手段と、ホスト間の通信が失われた場合にも対処できるように、他ホストからのホスト間処理同期要求の受理通知がなくても、一定時間毎に共有メモリの参照に基づくホスト間処理同期成立判断を前記ホスト間処理同期手段に行わせるべく処理同期確認要求を発行する各ホスト内のタイマ起動手段とを有している。

40

【発明の概要】

【発明が解決しようとする課題】

【0010】

しかし、共有メモリは、他の一般的な通信処理と、ソフトウェア上での利便性が全く異なっており、共有メモリを利用するためには各プロセッサノードのプログラム間での同期や排他などの仕組みをプログラマが何らかの方法で実現しなければならないといった利便性での問題があった。さらには、密結合型のマルチプロセッサシステム向けに開発されソケット通信を利用しているソフトウェアを、共有メモリを備えた疎結合型のマルチプロセ

50

ッサシステムに移行しようとする場合には、プログラム内のソースコード中でソケット通信を利用している箇所の記述を、共有メモリを利用する処理に書き換えなければならないため、プログラマにとってソフトウェアの開発負担が非常に増大するという問題があった。

【0011】

本発明は、このような課題に鑑みてなされたものであり、各プロセッサノード間での共有メモリを用いた通信処理で発生する通信オーバーヘッドを低減し、高速通信が可能なプロセッサ間通信装置、プロセッサ間通信方法、プログラムおよび記録媒体を提供することを目的とする。

【課題を解決するための手段】

【0012】

以上の課題を解決するために本発明では、第1および第2のプロセッサノードから構成され、第1および第2のプロセッサノードのそれぞれからアクセスすることが可能に共有して設けられた共有メモリを備えたプロセッサ間通信装置において、第1のプロセッサノードで動作しているプログラムに従って、書込手段が通信データを前記共有メモリに書き込むと、第1および第2のプロセッサノードをプログラム間で通信データの送受信が可能に接続するソケット通信を介して、第1の通知手段が書込通知情報を第2のプロセッサノードに送信する。そして、読出手段が、第1の通知手段により送信された書込通知情報に応じて共有メモリから通信データを読み出し、第2の通知手段が、読出通知情報をソケット通信を介して第1のプロセッサノードに送信する。

【0013】

このため、第1および第2の各プロセッサノード間で、同一の共有メモリを共有利用することによって、この共有メモリが分散共有メモリとして機能する。また、ソケットを介して書込通知情報や読出通知情報を送受信することによって、この共有メモリが複数のプロセッサノード間で利用できるというソケット通信の特徴と、広帯域であるという共有メモリの特徴とをそれぞれ有し、通信オーバーヘッドを低減し、高速通信を可能とすることができる。

【発明の効果】

【0014】

本発明によれば、各プロセッサノード間での共有メモリを用いた通信処理で発生する通信オーバーヘッドを低減し、高速通信が可能なプロセッサ間通信装置、プロセッサ間通信方法、プログラムおよび記録媒体を提供することができる。

【図面の簡単な説明】

【0015】

【図1】第1の実施形態におけるマルチプロセッサシステムの全体構成を示す説明図である。

【図2】第2の実施形態におけるマルチプロセッサシステムの全体構成を示す説明図である。

【図3】第2の実施形態におけるマルチプロセッサシステムのメインメモリの記憶領域の使用状態を示す説明図である。

【図4】第3の実施形態におけるマルチプロセッサシステムの全体構成を示す説明図である。

【図5】第1の実施形態におけるマルチプロセッサシステムのプロセッサノードから他のプロセッサノード内の共有メモリに対してCPUが書き込みを行った処理を示す説明図である。

【図6】第1の実施形態におけるマルチプロセッサシステムの複数のプロセッサノード間がソケットで接続されこのソケット600を利用して通信を行っている処理を示す説明図である。

【図7】第1の実施形態におけるマルチプロセッサシステムの複数のプロセッサノードのうち、2つのプロセッサノードが通信を行っている処理を示す説明図である。

10

20

30

40

50

【図 8】第 1 の実施形態におけるマルチプロセッサシステムの複数のプロセッサノードのうち、2 つのプロセッサノードが通信を行っている処理を示す説明図である。

【図 9】第 1 の実施形態におけるマルチプロセッサシステムの複数のプロセッサノード間がソケットで接続されている場合の通信データの送受信処理を示すシーケンス図である。

【図 10】第 1 の実施形態におけるマルチプロセッサシステムの通信データの送受信処理のコネクション確立後の具体的な例を示すシーケンス図である。

【図 11】第 1 の実施形態におけるマルチプロセッサシステムの通信データの送受信処理のコネクション確立後の他の具体的な例を示すシーケンス図である。

【図 12】第 1 の実施形態におけるマルチプロセッサシステムのクライアントで実行されるコネクション確立および切断の処理を示すシーケンス図である。

【図 13】第 1 の実施形態におけるマルチプロセッサシステムのソケットの動作を指示するためのソケット指示情報を記憶したテーブルを示す説明図である。

【図 14】第 1 の実施形態におけるマルチプロセッサシステムのソケット指示情報に従って共有メモリに対して行った処理の例を示す説明図である。

【図 15】第 1 の実施形態におけるマルチプロセッサシステムのソケット指示情報に従って共有メモリに対して行った処理の他の例を示す説明図である。

【図 16】第 1 の実施形態におけるマルチプロセッサシステムの共有メモリ内の記憶領域で割り当てられたメモリウィンドウのデータ構成を示す説明図である。

【図 17】第 1 の実施形態におけるマルチプロセッサシステムの共有メモリ内のリングバッファに通信データを用いて行った処理を示す説明図である。

【図 18】第 1 の実施形態におけるマルチプロセッサシステムの複数のプロセッサノードのうち 2 つのプロセッサノードが共有メモリに通信データを用いて行った処理を示す説明図である。

【図 19】第 1 の実施形態におけるマルチプロセッサシステムの複数のプロセッサノードのうち 2 つのプロセッサノードが共有メモリに通信データを用いて行った他の処理を示す説明図である。

【図 20】第 1 の実施形態におけるマルチプロセッサシステムのプロセッサノードのハードウェア構成を示す説明図である。

【発明を実施するための形態】

【0016】

〔第 1 の実施形態〕

以下、本発明に係る第 1 の実施形態をもって説明するが、本発明は、実施形態に限定されるものではない。図 1 は、本発明によるプロセッサ間通信装置の 1 つの例であるマルチプロセッサシステム 10 の第 1 の実施形態での全体構成を示す説明図である。このマルチプロセッサシステム 10 は、例えば共有メモリを備えた疎結合型のマルチプロセッサシステムであり、複数のプロセッサノードから構成され、共有メモリが各プロセッサノードのそれぞれからアクセス可能に共有して設けられている。

【0017】

マルチプロセッサシステム 10 は、図 1 に示すように、複数のプロセッサノード 100、110、120 と、これらの各プロセッサノード 100、110、120 間でデータなどの送受信を行うためのネットワークの 1 例である Ethernet（登録商標）130 と、各プロセッサノード 100、110、120 のそれぞれからアクセス可能に共有して設けられた共有メモリ 140 とから構成されている。

【0018】

プロセッサノード 100 は、アプリケーションなどの各種の処理を行うためのプログラムを実行する CPU 101 と、CPU 101 の処理に従って後述するようにメインメモリ 103 にアクセスしデータなどの読み出しを行うメモリコントローラ 102 と、CPU 101 の処理によって得られた各種のデータを記憶するメインメモリ 103 とを備えている。

【0019】

10

20

30

40

50

CPU101は、予め用意されたプログラムを実行することにより、例えばアプリケーション上での各種の処理を行い、メモリコントローラ102や後述するホストPCIブリッジ104、Ethernet（登録商標）カード105を介して外部から入力された通信データを用いて演算処理や判定処理などを行う。

【0020】

メモリコントローラ102は、メインメモリ103にアクセスして、CPU101での処理によって得られた演算結果や判定結果などのデータをメインメモリ103内に形成された共有メモリに書き込んだり、また、メモリコントローラ102は、この共有メモリから通信データを始めとする各データを読み出し、ホストPCIブリッジ104、Ethernet（登録商標）カード105、Ethernet（登録商標）130を介して他のプロセッサノード110、120に送信するための処理を行う。

10

【0021】

また、プロセッサノード100は、プロセッサノード100および外部との間のデータの入出力を行うホストPCIブリッジ104と、Ethernet（登録商標）130に接続されEthernet（登録商標）130を介して外部との間でデータの送受信を行うEthernet（登録商標）カード105とを備えている。

【0022】

ホストPCIブリッジ104は、プロセッサノード100の内部側がメモリコントローラ102に接続され、外部側がEthernet（登録商標）カード105、共有メモリ140に接続されている。ホストPCIブリッジ104は、外部側の接続対象をEthernet（登録商標）カード105または共有メモリ140に切り替えて、プロセッサノード100およびEthernet（登録商標）130間と、プロセッサノード100および他のプロセッサノード110、120間でのデータの入出力を行う。

20

【0023】

Ethernet（登録商標）カード105は、例えばプロセッサノード100のPCIカードスロットに装着されてEthernet（登録商標）130と直接接続され、Ethernet（登録商標）130を介して外部との間でデータの送受信を行う。

【0024】

Ethernet（登録商標）カード105は、他のプロセッサノード110、120のうちのいずれかから要求された通信データを要求するための要求情報を受信する。また、Ethernet（登録商標）カード105は、メモリコントローラ102により通信データが共有メモリ140に書き込まれたことを通知する書込通知情報を、他のプロセッサノード110、120にそれぞれ送信したり、メモリコントローラ102により通信データが読み出されたことを通知する読出通知情報を、他のプロセッサノード110、120にそれぞれ送信したり、また、他のプロセッサノード110、120のEthernet（登録商標）カード105により送信された読出通知情報に従って、メモリコントローラ102により読み出された通信データを、通信データを要求したプロセッサノードに送信する処理を行う。

30

【0025】

プロセッサノード110、120は、プロセッサノード100と同様の構成となっており、CPU101、メモリコントローラ102、メインメモリ103、ホストPCIブリッジ104、Ethernet（登録商標）カード105と同様の機能を有するCPU111、121、メモリコントローラ112、122、メインメモリ113、123、ホストPCIブリッジ114、124、Ethernet（登録商標）カード115、125を備えている。

40

【0026】

図20は、プロセッサノード100のハードウェア構成を示す説明図である。プロセッサノード100は、図20に示すように、CPU101にメモリコントローラ102が接続され、このメモリコントローラ102に対してメインメモリ103がデータの書き込みおよび読み出しが可能に、また、ホストPCIブリッジ104が通信データなどの各デ

50

ータの送受信が可能にそれぞれ接続されている。

【0027】

また、ホスト PCIブリッジ104には、複数種類のPCIデバイス機器が接続可能なPCIバス401が接続され、このPCIバス401にPCIバスを介して共有メモリ140に接続されたPCI PCIブリッジ402と、複数種類の各PCIデバイス403と、Ethernet（登録商標）130に接続されたEthernet（登録商標）カード105とがそれぞれ接続されている。

【0028】

メモリコントローラ102は、CPU101のホストバスに接続され、CPU101からメインメモリ103やEthernet（登録商標）カード105、PCI PCIブリッジ402、各種PCIデバイス403に対するリードまたはライトの要求を受け取ると、対象のデバイスに対してリクエストを振り分ける役目がある。マルチプロセッサシステム100の具体的な例としてパソコンを用いた場合、一般的には例えばIntel製のx86 CPUを搭載しているが、その際にはマザーボードというメイン基板の上にノースブリッジやMCH（Memory Controller Hub）として参照されるチップセットが搭載されており、図20のメモリコントローラ102はこれらのチップセットに相当している。

10

【0029】

一方、ホスト PCIブリッジ104は、サウスブリッジやICH（I/O Controller Hub）として参照されるチップセットに相当し、CPU101から各種Ethernet（登録商標）カード105、PCI PCIブリッジ402、各種PCIデバイス403に対するアクセス要求を発行したり、あるいはPCIデバイス403からメインメモリ103に対するDMAリクエストを受け取ってメモリコントローラ102に送信する機能を有する。

20

【0030】

図5は、例えばプロセッサノード110から他のプロセッサノード100内の共有メモリ140に対してCPU101が書き込みを行った処理を示す説明図である。図5に示すように、プロセッサノード100のメインメモリ103の一部の記憶領域がプロセッサノード110のメインメモリ113内のメモリウィンドウにマッピングされている。この時に、プロセッサノード110で動作しているプロセスがメモリウィンドウの記憶領域に対して書き込みを行うと、この書き込まれたデータと記憶領域を示すアドレスがプロセッサノード100に送信され、リモートメモリとして利用されているプロセッサノード100内のメインメモリ103に対して書き込み処理が行われる。そのメインメモリ103に対してプロセッサノード100で動作しているプロセスが読み出しを行うことにより、プロセッサノード110からプロセッサノード100に対して通信データを伝達することが可能となる。

30

【0031】

図5は、別々のプロセッサノードで動作しているプロセスの間で同一の共有メモリ140を共有利用することでデータの伝達を行うことが可能となっているが、このようなメモリを分散共有メモリとして参照する。分散共有メモリは、複数のプロセッサノード間で利用できるというソケット通信の特徴と、広帯域であるという共有メモリの特徴とをそれぞれ有している。

40

【0032】

マルチプロセッサシステム10が備えているOSのうち、例えばUNIX（登録商標）のような各プロセッサノード100、110、120でのプログラム内のプロセス間でのメモリ保護機能を持っているOSではプロセス間でグローバル変数等を用いてデータの伝達を行うことができない。このようなOSとしては、プロセス間でデータの送受信などのコミュニケーションを行うための機能として、ソケット通信や共有メモリといったIPC（Inter-Process Communication）を実現する手段を備えている。ソケット通信は、ネットワーク透過な機能なので、各プロセッサノード100、1

50

10、120内のみならず、各プロセッサノード100、110、120間でデータをやり取りする場合にも使用することができるという長所があるが、その反面、通信オーバーヘッドが多いために大量のデータ送受信の際には処理効率が低下するという欠点もある。共有メモリ140は、通信オーバーヘッドが低く大量のデータ送受信でも効率良く実行できるが（広帯域）、通常はそれぞれのプロセッサノード100、110、120内でしか利用できないため、複数のプロセッサノード100、110、120によって構成されるマルチプロセッサシステムではあまり利用されることがなかった。

【0033】

図6は、複数のプロセッサノード100、110、120間がソケット600で接続され、例えばプロセッサノード100で動作しているプログラムがこのソケット600を利用してソケット通信を行っている処理を示す説明図である。ここで、ソケット600は、複数のプロセッサノード100、110、120のそれぞれの間でIPアドレス、ポート番号が設定されて、このソケット600によって接続された各プロセッサノードが相互にソケット通信を行うことが可能となっている状態を示す。

10

【0034】

各プロセッサノード100、110、120のうち、プロセッサノード100、110間は、図6に示すように、例えばEthernet（登録商標）130などの各種のネットワークを透過するソケット600で接続された二つのプログラム601、602の間で双方向に通信データを送受信することができる。

【0035】

このソケット600で接続されることによって示されるソケット通信では、各プロセッサノード100、110、120のうち、任意の2つのプロセッサノードで動作しているそれぞれのプログラム間で通信データを送受信することができる。また、ソケット600を通して送受信可能なデータは、バイトストリームであり、ここで、バイトストリームとは、図6に示すように、送受信のそれぞれのデータでバイトデータが列に並んだデータであり、送信元が送出した各バイトデータがその送出された送出順序のまま受信側のプロセッサノードに届くようになっている。

20

【0036】

図7、図8は、複数のプロセッサノード100、110、120のうち、例えばプロセッサノード100、110で動作しているプログラムが通信を行っている処理を示す説明図である。プロセッサノード100、110間では、図7に示すように、プロセッサノード100、110でそれぞれ内部のメインメモリ103、113からプログラムを読み出して動作し、各プログラムの処理に従って相互に通信データの書き込みを要求する。

30

【0037】

プロセッサノード100からプロセッサノード110への書き込みの要求では、図8に示すように、プロセッサノード100のメインメモリ103の一部の記憶領域がプロセッサノード110のメインメモリ113内のメモリウィンドウにマッピングされている。この時に、プロセッサノード110のプログラムで動作しているプロセスBがメモリウィンドウの記憶領域に対して書き込みを行うと、この書き込まれたデータと記憶領域を示すアドレスがプロセッサノード100に送信され、リモートメモリとして利用されているプロセッサノード100内のメインメモリ103に対して書き込み処理が行われる。そのメインメモリ103に対してプロセッサノード100で動作しているプロセスが読み出しを行うことにより、プロセッサノード110からプロセッサノード100に対して通信データを伝達することが可能となる。

40

【0038】

一方、プロセッサノード110からプロセッサノード100への書き込みの要求では、プロセッサノード110のメインメモリ113の一部の記憶領域がプロセッサノード100のメインメモリ103内のメモリウィンドウにマッピングされている。そして、同様にプロセッサノード100からプロセッサノード110に対して通信データが伝達される。

50

【0039】

図9は、図6に示すように、複数のプロセッサノード100、110、120間がソケット600で接続されている場合の通信データの送受信処理を示すシーケンス図である。この場合の送受信処理では、図9に示すように、一方のサーバ90内のプロセッサノード100と、もう一方のクライアント91内のプロセッサノード110とが、それぞれソケット600を介して接続されている。

【0040】

サーバ90は、ステップS901でクライアント91から送信されたソケット600使用要求(socket)のデータを受け取ると、これに応じて同様のデータと、(bind)、(listen)、(accept)の各データを生成し、サーバ90およびクライアント91間の接続を確立するための処理を行う。

10

【0041】

サーバ90、クライアント91は、ステップS902でソケット600を初期化する処理を行い、各種のデータを送受信可能な状態に設定する。

【0042】

クライアント91は、ステップS903でソケット600を介して送信された信号(p_cond_signal)をCPU111で動作するプログラムに一度戻して、この信号に従って通信データを共有メモリ140に送信し、通信データを共有メモリ140に書き込む処理を行う。

20

【0043】

クライアント91は、ステップS904でステップS903において共有メモリ140に通信データを書き込んだことを通知するための書込通知情報を、ソケット600を介してサーバ90に送信する。

【0044】

一方、サーバ90は、ステップS903においてクライアント91と同様に、ソケット600を介して送信された信号(p_cond_signal)をCPU101で動作するプログラムに一度戻す。そして、サーバ90は、ステップS904においてクライアント91から送信された書込通知情報を受信し、ステップS905でこの書込通知情報に従って、通信データを共有メモリ140から読み出す処理を行う。

30

【0045】

サーバ90は、ステップS906でステップS905において共有メモリ140から通信データを読み出したことを通知するための読出通知情報を、ソケット600を介してクライアント91に送信する。

【0046】

クライアント91は、ステップS907でサーバ90およびクライアント91間の接続を切断することを通知するための切断通知信号を、ソケット600を介してサーバ90に送信する。一方、サーバ90は、ステップS908で切断通知信号を受信すると、これに応じて接続を切断することに応答することを通知するための切断通知信号を、ソケット600を介してクライアント91に送信する。

40

【0047】

サーバ90、クライアント91は、ステップS909でソケット600の動作を終了する処理を行い、各種のデータを送受信不可能な状態に設定する。また、サーバ90は、サーバ90およびクライアント91間の接続を切断するための処理を行う。

【0048】

図10は、図9に示すように、通信データの送受信処理の接続確立後の具体的な例を示すシーケンス図である。この例の送受信処理では、図10に示すように、クライアント91がステップS1003で送信された信号(p_cond_signal)をCPU101で動作するプログラムに一度戻して、この信号に従って256KBの2つの通信データを共有メモリ140に送信する。クライアント91は、2つのそれぞれの通信データを共有メモリ140の記憶領域上にコピーして書き込む処理を2回行う。

50

【 0 0 4 9 】

クライアント 9 1 は、ステップ S 1 0 0 4 でステップ S 1 0 0 3 において共有メモリ 1 4 0 に通信データを書き込んだことを通知するための書込通知情報を、サーバ 9 0 に送信する。

【 0 0 5 0 】

一方、サーバ 9 0 は、ステップ S 1 0 0 3 においてクライアント 9 1 と同様に、送信された信号 (p _ c o n d _ s i g n a l) を CPU 1 0 1 で動作するプログラムに一度戻し、ステップ S 1 0 0 5 で通信データの待ち状態となる。そして、サーバ 9 0 は、ステップ S 1 0 0 4 においてクライアント 9 1 から送信された書込通知情報を受信し、ステップ S 1 0 0 6 でこの書込通知情報に従って、5 1 2 K B の通信データを共有メモリ 1 4 0 の記憶領域上からコピーして読み出す処理を行う。

10

【 0 0 5 1 】

サーバ 9 0 は、ステップ S 1 0 0 7 でステップ S 1 0 0 6 において共有メモリ 1 4 0 から通信データを読み出したことを通知するための読出通知情報を、クライアント 9 1 に送信する。

【 0 0 5 2 】

図 1 1 は、図 9 に示すように、通信データの送受信処理の接続確立後の他の具体的な例を示すシーケンス図である。この例の送受信処理では、図 1 1 に示すように、クライアント 9 1 がステップ S 1 1 0 3 で送信された信号 (p _ c o n d _ s i g n a l) を CPU 1 1 1 で動作するプログラムに一度戻して、この信号に従って 1 . 5 M B の通信データを共有メモリ 1 4 0 に送信する。クライアント 9 1 は、この通信データを 5 1 2 K B のそれぞれ 3 つの同一容量の通信データに分割し、これらのうち 2 つの通信データを共有メモリ 1 4 0 の記憶領域上にコピーして書き込む処理を 2 回行う。

20

【 0 0 5 3 】

クライアント 9 1 は、ステップ S 1 1 0 5 でステップ S 1 1 0 3 において共有メモリ 1 4 0 に通信データを書き込んだことを通知するための書込通知情報を、それぞれ 2 回の書き込み処理毎にサーバ 9 0 に送信する。ここで、例えば共有メモリ 1 4 0 の容量が 1 M B の場合には、共有メモリ 1 4 0 の記憶領域が満たされ書き込む処理が不可能となるため、ステップ S 1 1 0 4 で通信データの待ち状態となる。

【 0 0 5 4 】

一方、サーバ 9 0 は、上述のステップ S 1 1 0 3 でクライアント 9 1 と同様に、送信された信号 (p _ c o n d _ s i g n a l) を CPU 1 0 1 で動作するプログラムに一度戻す。そして、サーバ 9 0 は、ステップ S 1 1 0 5 においてクライアント 9 1 から送信された書込通知情報をそれぞれ受信し、ステップ S 1 1 0 6 でこの書込通知情報に従って、1 つ目の通信データを共有メモリ 1 4 0 の記憶領域上からコピーして読み出す処理を行う。

30

【 0 0 5 5 】

サーバ 9 0 は、ステップ S 1 1 0 7 でステップ S 1 1 0 6 において共有メモリ 1 4 0 から通信データを読み出したことを通知するための読出通知情報を、クライアント 9 1 に送信する。

【 0 0 5 6 】

クライアント 9 1 は、ステップ S 1 1 0 7 においてサーバ 9 0 から送信された書込通知情報を受信し、これに応じて通信データの待ち状態を解除する。そして、クライアント 9 1 は、ステップ S 1 1 0 8 で共有メモリ 1 4 0 の記憶領域に書き込まれた通信データを消去した上で、残りの 3 つ目の通信データを共有メモリ 1 4 0 の記憶領域上にコピーして書き込む処理を行う。

40

【 0 0 5 7 】

クライアント 9 1 は、ステップ S 1 1 0 9 でステップ S 1 1 0 8 において共有メモリ 1 4 0 に通信データを書き込んだことを通知するための書込通知情報をサーバ 9 0 に送信する。

【 0 0 5 8 】

50

図12は、図11に示すように、クライアント91で実行されるコネクション確立および切断の処理を示すシーケンス図である。これらの処理では、図12に示すように、マルチプロセッサシステム10の起動時にまず、クライアント91が、ステップS1201でメモリコントローラ102、ホスト-PCIBリッジ104を介して、共有メモリ140内で通信データを書き込むための記憶領域を確保する。そして、クライアント91は、ステップS1202、S1203でソケット600使用要求(socket)のデータをCPU111からメモリコントローラ102、ホスト-PCIBリッジ104に送信して初期化する。また、コネクションを確立し各種のデータが送受信可能な状態に設定する。

【0059】

クライアント91は、ステップS1204で共有メモリ140内で通信データを書き込むための記憶領域を割り当てるための割り当て要求をメモリコントローラ102、ホスト-PCIBリッジ104に送信する。また、クライアント91は、ステップS1205で共有メモリ140を書き込みや読み出しなどのアクセス待ち状態に設定する。

10

【0060】

クライアント91は、ステップS1206でメモリコントローラ102、ホスト-PCIBリッジ104により共有メモリ140にアクセスし、通信データを書き込むための記憶領域を割り当てるための処理を行う。例えば、クライアント91は、共有メモリ140内の記憶領域のうち、データが書き込まれていない空き領域、またはデータが消去済みの領域を検出して、これらの検出した各領域から通信データを書き込むための領域を確保し、この確保した領域内の各アドレスなどを割り当てる処理を行う。クライアント91は、ステップS1206において記憶領域を割り当てた割り当て結果を、ステップS1207でメモリコントローラ102、ホスト-PCIBリッジ104によりCPU111に送信する。

20

【0061】

クライアント91は、ステップS1208でステップS1207において送信された割り当て結果や割り当てた各アドレスを通知するための割当通知情報を、サーバ90に送信する。そして、クライアント91は、ステップS1209で、サーバ90から送信された、サーバ90で同様に共有メモリ140内の記憶領域で割り当てられた割り当て結果や割り当てた各アドレスを通知するための割当通知情報を、受信する。

30

【0062】

クライアント91は、ステップS1210、S1211でソケット600を介して送信された信号(p_cond_signal)をCPU111で動作するプログラムに一度戻して、この信号に従って通信データを共有メモリ140に送信し、通信データを共有メモリ140に書き込む処理を行う。ここで、通信データは、共有メモリ140内の記憶領域のうち、ステップS1206において割り当てた領域に書き込まれる。

【0063】

クライアント91は、ステップS1212でステップS1211において共有メモリ140に通信データを書き込んだことを通知するための書込通知情報を、ソケット600を介してサーバ90に送信する。

40

【0064】

一方、サーバ90では、ステップS1210においてクライアント91と同様に、ソケット600を介して送信された信号(p_cond_signal)をCPU101で動作するプログラムに一度戻す。そして、サーバ90では、ステップS1211においてクライアント91から送信された書込通知情報を受信し、この書込通知情報に従って、通信データを共有メモリ140から読み出す処理が行われる。

【0065】

サーバ90は、ステップS1213でステップS1211において共有メモリ140から通信データを読み出したことを通知するための読出通知情報を、ソケット600を介してクライアント91に送信する。

【0066】

50

クライアント 9 1 は、ステップ S 1 2 1 4 でサーバ 9 0 およびクライアント 9 1 間のコネクションを切断することを通知するための切断通知信号を、ソケット 6 0 0 を介してサーバ 9 0 に送信する。一方、サーバ 9 0 では、ステップ S 1 2 1 5 で切断通知信号を受信すると、これに応じてコネクションを切断することに応答することを通知するための応答通知信号を、ソケット 6 0 0 を介してクライアント 9 1 に送信する。

【 0 0 6 7 】

クライアント 9 1 は、ステップ S 1 2 1 6 で共有メモリ 1 4 0 内の記憶領域を解放するための解放要求をメモリコントローラ 1 0 2、ホスト - P C I ブリッジ 1 0 4 に送信する。また、クライアント 9 1 は、ステップ S 1 2 1 7 で共有メモリ 1 4 0 を書き込みや読み出しなどのアクセス待ち状態に設定する。

10

【 0 0 6 8 】

クライアント 9 1 は、ステップ S 1 2 1 8 でメモリコントローラ 1 0 2、ホスト - P C I ブリッジ 1 0 4 により共有メモリ 1 4 0 にアクセスし、ステップ S 1 2 0 6 において割り当てた記憶領域を解除し、解放するための処理を行う。クライアント 9 1 は、ステップ S 1 2 1 8 において記憶領域を解放した解放結果を、ステップ S 1 2 1 9 でメモリコントローラ 1 0 2、ホスト - P C I ブリッジ 1 0 4 により C P U 1 1 1 に送信する。

【 0 0 6 9 】

クライアント 9 1 は、ステップ S 1 2 2 0 でソケット 6 0 0 の動作を終了する処理を行い、各種のデータを送受信不可能な状態に設定する。そして、処理を終了するための信号 (p _ c o n d _ s i g n a l) を C P U 1 0 1 で動作するプログラムに戻して、この信号に従ってプログラムの実行を終了する。

20

【 0 0 7 0 】

図 1 3 は、ソケット 6 0 0 の動作を指示するためのソケット指示情報を記憶したテーブル 1 3 0 0 を示す説明図である。このテーブル 1 3 0 0 では、図 1 3 に示すように、複数種類の各ソケット指示情報と、各ソケット指示情報のそれぞれで指示されるソケット 6 0 0 の動作の内容の情報とが関連付けられて記憶されている。例えば、テーブル 1 3 0 0 では、図 1 3 に示すように、ソケット指示情報である「 s o c k f d 」と、ソケット 6 0 0 の動作の内容を示す「ソケットディスクリプタ」とが関連付けられて記憶されている。また、他のソケット指示情報である「 s t a t e 」と、ソケット 6 0 0 の動作の内容を示す「ソケットステート」とが関連付けられて記憶されている。

30

【 0 0 7 1 】

図 1 4 は、ソケット指示情報に従って共有メモリ 1 4 0 に対して行った処理の例を示す説明図である。サーバ 9 0 は、図 1 4 に示すように、共有メモリ 1 4 0 にアクセスし、共有メモリ 1 4 0 内の記憶領域から通信データを書き込むための領域を確保し、この確保した領域内の各アドレスなどを割り当てるための処理を行う。そして、サーバ 9 0 は、記憶領域を割り当てたことを示すソケット指示情報「 l o c a l . b u f f _ a d r s 」をソケット 6 0 0 を介してホスト - P C I ブリッジ 1 0 4 により C P U 1 0 1 に送信する。

【 0 0 7 2 】

また、サーバ 9 0 は、図 1 4 に示すように、共有メモリ 1 4 0 にアクセスし、共有メモリ 1 4 0 内の記憶領域から、クライアント 9 1 により書き込まれた通信データを読み出すための領域を確保し、この確保した領域内の各アドレスなどを割り当てるための処理を行う。そして、サーバ 9 0 は、記憶領域を割り当てたことを示すソケット指示情報「 r e m o t e . b u f f _ a d r s 」をソケット 6 0 0 を介してホスト - P C I ブリッジ 1 0 4 により C P U 1 0 1 に送信する。

40

【 0 0 7 3 】

一方、クライアント 9 1 は、図 1 4 に示すように、共有メモリ 1 4 0 にアクセスし、共有メモリ 1 4 0 内の記憶領域から通信データを書き込むための領域を確保し、この確保した領域内の各アドレスなどを割り当てるための処理を行う。そして、クライアント 9 1 は、記憶領域を割り当てたことを示すソケット指示情報「 l o c a l . b u f f _ a d r s 」をソケット 6 0 0 を介してホスト - P C I ブリッジ 1 1 4 により C P U 1 1 1 に送信す

50

る。

【0074】

また、クライアント91は、図14に示すように、共有メモリ140にアクセスし、共有メモリ140内の記憶領域から、サーバ90により書き込まれた通信データを読み出すための領域を確保し、この確保した領域内の各アドレスなどを割り当てるための処理を行う。そして、クライアント91は、記憶領域を割り当てたことを示すソケット指示情報「remote_buf_adrs」をソケット600を介してホスト-PCIブリッジ114によりCPU111に送信する。

【0075】

図15は、ソケット指示情報に従って共有メモリ140に対して行った処理の他の例を示す説明図である。図15に示すように、例えばサーバ90やクライアント91などが共有メモリ140にアクセスし、共有メモリ140内の記憶領域から通信データを書き込むための複数の領域を確保し、これらの確保したそれぞれの領域内の各アドレスなどを割り当てるための処理を行う。

10

【0076】

そして、図15に示すように、1つ目(メモリウィンドウ1)の記憶領域を割り当てたことを示すソケット指示情報をソケット600を介してCPU101に送信する。ソケット指示情報には、割り当てた領域の容量(256MB)を示す「local_buf_size」と、割り当てた領域のアドレスを示す「local_buf_adrs」とが含まれている。

20

【0077】

また、2つ目(メモリウィンドウ2)の記憶領域を割り当てたことを示すソケット指示情報をソケット600を介してCPU101に送信する。ソケット指示情報には、割り当てた領域の容量(512MB)を示す「local_buf_size」と、割り当てた領域のアドレスを示す「local_buf_adrs」とが含まれている。

【0078】

また、3つ目(メモリウィンドウ3)の記憶領域を割り当てたことを示すソケット指示情報をソケット600を介してCPU101に送信する。ソケット指示情報には、割り当てた領域の容量(128MB)を示す「local_buf_size」と、割り当てた領域のアドレスを示す「local_buf_adrs」とが含まれている。

30

【0079】

図16は、共有メモリ140内の記憶領域で割り当てられたメモリウィンドウのデータ構成を示す説明図である。このメモリウィンドウには、図16(a)に示すように、マジックナンバーと、バッファサイズと、ライトポイントと、リードポイントと、リングバッファの各データが含まれている。マジックナンバーは、メモリウィンドウの先頭位置を識別するための識別情報である。

【0080】

バッファサイズは、メモリウィンドウの共有メモリ140内での全体の容量を示す例えばバイト数などのデータである。ライトポイントは、リングバッファ内での通信データを書き込む際の開始位置を示すポイントであり、リングバッファ内でのオフセットアドレスが含まれている。リードポイントは、リングバッファ内での通信データを読み出す際の開始位置を示すポイントであり、リングバッファ内でのオフセットアドレスが含まれている。

40

【0081】

メモリウィンドウに含まれるリングバッファは、記憶領域内の開始位置が空きとなっている場合には、通信データの内容をリングバッファの領域内に全て書き込むと、これに応じて、リングバッファ内の先頭位置に戻って書き込み処理を続けて行うことが可能であり、循環的に書き込み処理が可能なバッファである。リングバッファには、図16(b)に示すように、ソケット600を介して送受信する際の複数にそれぞれ分けられたデータ群であるパケットが含まれている。

50

【 0 0 8 2 】

それぞれの各パケットには、図 1 6 (c) に示すように、パケットデータサイズと、パケットデータが含まれている。パケットデータサイズは、パケットデータ全体の容量を示す例えばバイト数などのデータであり、4 b y t e から成る。パケットデータには、通信データの内容が複数の各パケット毎に分けて含まれている。

【 0 0 8 3 】

図 1 7 は、共有メモリ 1 4 0 内のリングバッファに通信データを用いて行った処理を示す説明図である。図 1 7 に示すように、まず、(1) で例えばサーバ 9 0 やクライアント 9 1 などが共有メモリ 1 4 0 にアクセスし、共有メモリ 1 4 0 のリングバッファを初期状態に設定する。即ち、ライトポインタおよびリードポインタの位置をリングバッファ内で先頭位置に設定する処理を行う。

10

【 0 0 8 4 】

また、(2) で、通信データの内容を含む 1 つ目のパケットをリングバッファに書き込み、その後、これに伴ってライトポインタの位置をこのパケットの終端位置に移動させる。ここで、通信データの内容を書き込む処理は、4 b y t e 単位毎に行う。5 b y t e などのデータを書き込むときには 2 つの 4 b y t e 単位毎に 8 b y t e 分まで書き込む処理を行う。(3) で、続けて通信データの内容を含む 2 つ目のパケットを書き込み、その後、これに伴ってライトポインタの位置をこのパケットの終端位置に移動させる。

【 0 0 8 5 】

次に、(4) で、他のサーバ 9 0 やクライアント 9 1 などが共有メモリ 1 4 0 にアクセスし、共有メモリ 1 4 0 のリングバッファから、上述の (1) で書き込まれた 1 つ目のパケットを読み出す処理を行う。このとき、この読み出したパケットを消去し、その後、これに伴ってリードポインタの位置を 2 つ目のパケットの開始位置に移動させる。

20

【 0 0 8 6 】

また、(5) で、サーバ 9 0 やクライアント 9 1 などが共有メモリ 1 4 0 にアクセスし、通信データの内容を含む 3 つ目のパケットをリングバッファに書き込む処理を行う。このとき、通信データの内容を書き込んでいき、リングバッファの領域内に全て書き込んだ場合には、リングバッファ内の先頭位置に戻って上述の (4) で消去した領域に続けて書き込む処理を行う。その後、これに伴ってライトポインタの位置も同様に、リングバッファ内の先頭位置に戻ってこのパケットの終端位置に移動させる。

30

【 0 0 8 7 】

(6) で、通信データの内容を含む 4 つ目のパケットを書き込み、その後、これに伴ってライトポインタの位置をこのパケットの終端位置に移動させる。ここで、通信データの内容を書き込んでいき、リングバッファの領域内に空きが残り 4 b y t e となった場合には、この時点でリングバッファ内がフル状態として、書き込み不可に設定する。

【 0 0 8 8 】

図 1 8 は、複数のプロセッサノード 1 0 0、1 1 0、1 2 0 のうち、例えばプロセッサノード 1 0 0、1 1 0 が共有メモリ 1 4 0 に通信データを用いて行った処理を示す説明図である。図 1 8 に示すように、まず、プロセッサノード 1 0 0 の C P U 1 0 1 は、ステップ S 1 8 0 1 で共有メモリ 1 4 0 にアクセスして、リングバッファ内で先頭位置 (0 x E 8 0 0 1 0 0 8) を指定する。そして、C P U 1 0 1 は、1 つ目の通信データをリングバッファに書き込み、その後、これに伴ってライトポインタの位置をこの通信データの終端に含まれるパケットの終端位置 (0 x 0 0 0 0 F F F 0) に移動させる。

40

【 0 0 8 9 】

一方、プロセッサノード 1 1 0 の C P U 1 1 1 は、ステップ S 1 8 0 2 で共有メモリ 1 4 0 にアクセスし、共有メモリ 1 4 0 のリングバッファから、ステップ S 1 8 0 1 において書き込まれた通信データをパケット毎に順次読み出す処理を行う。このとき、この読み出したパケットを消去し、その後、これに伴ってリードポインタの位置を、この読み出したパケットの次のパケットの開始位置に移動させる。

【 0 0 9 0 】

50

CPU111は、ステップS1803、S1804で共有メモリ140に続けてアクセスし、共有メモリ140のリングバッファから、ステップS1801において書き込まれた通信データをパケット毎に順次読み出す処理を繰り返し行っていく。このとき、これらの読み出した各パケットをそれぞれ消去し、その後、これに伴ってリードポインタの位置を、これらの読み出したパケットの次のパケットの開始位置に移動させる。

【0091】

次に、CPU101は、ステップS1805で共有メモリ140にアクセスして、リングバッファ内で先頭位置(0xE800 1008)を指定する。そして、CPU101は、2つ目の通信データをリングバッファに書き込み、その後、これに伴ってライトポインタの位置をこの通信データの終端に含まれるパケットの終端位置(0x0001 0010)に移動させる。

10

【0092】

CPU111は、ステップS1806、S1807で共有メモリ140に続けてアクセスし、共有メモリ140のリングバッファから、ステップS1806、S1807において書き込まれた通信データをパケット毎に順次読み出す処理を行う。このとき、これらの読み出した各パケットをそれぞれ消去し、その後、これに伴ってリードポインタの位置を、これらの読み出したパケットの次のパケットの開始位置に移動させる。

【0093】

図19は、複数のプロセッサノード100、110、120のうち、例えばプロセッサノード100、110が共有メモリ140に通信データを用いて行った他の処理を示す説明図である。図19に示すように、まず、プロセッサノード100のCPU101は、ステップS1901、S1902で共有メモリ140にアクセスして、共有メモリ140内の2つのリングバッファ内でそれぞれ先頭位置(0xE800 1008)、(0xE800 100A)を指定する。そして、CPU101は、2つの通信データをこれらの各リングバッファにそれぞれ書き込み、その後、これに伴って各ライトポインタの位置をこれらの各通信データの終端に含まれるパケットの終端位置(0x0000 FFF0)、(0x0000 0000)に移動させる。

20

【0094】

一方、プロセッサノード110のCPU111は、ステップS1903で共有メモリ140にアクセスし、共有メモリ140のリングバッファから、ステップS1901、S1902において書き込まれた通信データをパケット毎に順次読み出す処理を行う。このとき、これらの読み出した各パケットをそれぞれ消去し、その後、これに伴ってリードポインタの位置を、これらの読み出したパケットの次のパケットの開始位置に移動させる。

30

【0095】

次に、CPU101は、ステップS1904、S1906で共有メモリ140にアクセスして、共有メモリ140内の2つのリングバッファ内でそれぞれ先頭位置(0xE800 1008)、(0xE800 100A)を指定する。そして、CPU101は、2つの通信データをこれらの各リングバッファにそれぞれ書き込み、その後、これに伴って各ライトポインタの位置をこれらの各通信データの終端に含まれるパケットの終端位置(0x0000 0010)、(0x0000 0001)に移動させる。

40

【0096】

ここで、上述のステップS1904およびS1906での処理の間、ステップS1905でCPU111は、共有メモリ140にアクセスし、共有メモリ140のリングバッファから、ステップS1904、S1906において書き込まれた通信データをパケット毎に順次読み出す処理を行う。このとき、これらの読み出した各パケットをそれぞれ消去し、その後、これに伴ってリードポインタの位置を、これらの読み出したパケットの次のパケットの開始位置に移動させる。

【0097】

また、CPU111は、ステップS1907で共有メモリ140に続けてアクセスし、共有メモリ140のリングバッファから、ステップS1904、S1906において書き

50

込まれた通信データをパケット毎に順次読み出す処理を繰り返し行っていく。このとき、これらの読み出した各パケットをそれぞれ消去し、その後、これに伴ってリードポインタの位置を、これらの読み出したパケットの次のパケットの開始位置に移動させる。

【0098】

以上のように、第1の実施形態におけるマルチプロセッサシステム10では、各プロセッサノード100、110、120がネットワークを透過するソケット600で互いに接続されて、双方向に通信データが送受信可能となっている。共有メモリ140に対して通信データを用いて書き込みまたは読み出しの各処理を行う際には、クライアント91のCPU111が信号に従って通信データを共有メモリ140に送信し、通信データを共有メモリ140に書き込む処理を行う。そして、CPU111は、ソケット600を介してサーバ90に書込通知情報を送信し、共有メモリ140に通信データを書き込んだことを通知する。

10

【0099】

一方、サーバ90は、この書込通知情報に従って通信データを共有メモリ140から読み出す処理を行う。また、サーバ90は、ソケット600を介してクライアント91に読出通知情報を送信し、共有メモリ140から通信データを読み出したことを通知する。

【0100】

このため、各プロセッサノード100、110、120の間をソケット600で接続して、各プロセッサノード100、110、120で動作しているプログラム上のプロセスの間で、同一の共有メモリ140を共有利用することによって、この共有メモリ140が分散共有メモリとして機能する。この共有メモリ140は、複数のプロセッサノード間で利用できるというソケットの特徴と、広帯域であるという共有メモリ140の特徴とをそれぞれ有し、通信処理で発生する通信オーバーヘッドを低減し、高速通信を可能とすることができる。

20

【0101】

〔第2の実施形態〕

以下、本発明に係る第2の実施形態をもって説明する。図2は、本発明におけるプロセッサ間通信装置の1つの例であるマルチプロセッサシステム10の第2の実施形態での全体構成を示す説明図である。このマルチプロセッサシステム10は、複数のプロセッサノード200、210、220と、これらの各プロセッサノード200、210、220間でデータなどの送受信を行うためのEthernet（登録商標）230と、各プロセッサノード200、210、220のそれぞれをデータなどの送受信が可能に接続するPCIバス240とから構成されている。

30

【0102】

プロセッサノード200は、CPU201と、メモリコントローラ202と、メインメモリ203と、ホストPCIブリッジ204と、Ethernet（登録商標）カード205とを備えている。これらの各構成要素は、第1の実施形態におけるプロセッサノード100が備えているCPU101、メモリコントローラ102、メインメモリ103、ホストPCIブリッジ104、Ethernet（登録商標）カード105と同様の機能を有しており説明を省略する。

40

【0103】

また、プロセッサノード200は、上述の各構成要素に加えて、PCIバス240に接続され、このPCIバス240を介して外部との間でデータの送受信を行うPCI-PCIブリッジ206を備えている。PCI-PCIブリッジ206は、例えばプロセッサノード200内で、ホストPCIブリッジ204およびPCIバス240間に直接接続され、CPU201に対してPCIバス240を介して外部との間でデータの送受信を行う。

【0104】

プロセッサノード210、220は、プロセッサノード200と同様の構成となっており、CPU201、メモリコントローラ202、メインメモリ203、ホストPCIブ

50

リッジ 204、Ethernet（登録商標）カード 205 と同様の機能を有する CPU 211、221、メモリコントローラ 212、222、メインメモリ 213、223、ホスト PCIブリッジ 214、224、Ethernet（登録商標）カード 215、225 を備えている。

【0105】

そして、図示していないが、第 2 の実施形態においてもマルチプロセッサシステム 10 では、各プロセッサノード 200、210、220 がネットワークを透過するソケット 600 で互いに接続されて、双方向に通信データが送受信可能となっている。

【0106】

図 3 は、メインメモリ 203 の記憶領域の使用状態を示す説明図である。この第 2 の実施形態では、図 3 に示すように、メインメモリ 203、213、223 内の記憶領域の一部を用いて共有メモリとして利用する。このような構成とすることによって、上述の第 1 の実施形態において単一の共有メモリ 140 を各プロセッサノード 100、110、120 と別個に設ける構成と比較して、コストの低減を図ることが可能となる。

【0107】

〔第 3 の実施形態〕

以下、本発明に係る第 3 の実施形態をもって説明する。図 4 は、本発明におけるプロセッサ間通信装置の 1 つの例であるマルチプロセッサシステム 10 の第 3 の実施形態での全体構成を示す説明図である。このマルチプロセッサシステム 10 は、複数のプロセッサノード 300、310、320 と、これらの各プロセッサノード 300、310、320 間でデータなどの送受信を行うための Ethernet（登録商標）330 と、各プロセッサノード 300、310、320 のそれぞれをデータなどの送受信が可能に接続する PCI-Express スイッチ 340 とから構成されている。

【0108】

プロセッサノード 300 は、CPU 301 と、メモリコントローラ 302 と、メインメモリ 303 と、ホスト PCIブリッジ 304 と、Ethernet（登録商標）カード 305 とを備えている。これらの各構成要素は、第 1 の実施形態におけるプロセッサノード 100 が備えている CPU 101、メモリコントローラ 102、メインメモリ 103、ホスト PCIブリッジ 104、Ethernet（登録商標）カード 105 と同様の機能を有しており説明を省略する。

【0109】

プロセッサノード 310、320 は、プロセッサノード 300 と同様の構成となっており、CPU 301、メモリコントローラ 302、メインメモリ 303、ホスト PCIブリッジ 304、Ethernet（登録商標）カード 305 と同様の機能を有する CPU 311、321、メモリコントローラ 312、322、メインメモリ 313、323、ホスト PCIブリッジ 314、324、Ethernet（登録商標）カード 315、325 を備えている。

【0110】

そして、図示していないが、第 3 の実施形態においてもマルチプロセッサシステム 10 では、各プロセッサノード 300、310、320 がネットワークを透過するソケット 600 で互いに接続されて、双方向に通信データが送受信可能となっている。

【0111】

この第 3 の実施形態では、メインメモリ 303、313、323 内の記憶領域の一部を用いて共有メモリとして利用する。このような構成とすることによって、上述の第 1 の実施形態において単一の共有メモリ 140 を各プロセッサノード 100、110、120 と別個に設ける構成と比較して、コストの低減を図ることが可能となる。

【0112】

〔他の実施形態〕

上述の実施形態において、ソケット 600 を介して書込通知情報を送信することによって、共有メモリ 140 に通信データを書き込んだことを通知しているが、共有メモリ 14

10

20

30

40

50

0に書き込まれた通信データのデータ量が所定量に達した場合に、通知するようにしても良い。

【0113】

この場合には、プロセッサノード100、110、120を始めとする各プロセッサノードが予め共有メモリ140に書き込まれたデータ量を検出する検出部を備え、共有メモリ140への通信データの書き込み処理が開始されると、これに応じて書き込まれたデータ量を検出する。そして、このデータ量が予め設定された所定量に達すると、通信データが書き込まれたこと、および書き込まれたデータ量を通知するデータ量通知情報を生成し、ソケット600を介して送信する。

【0114】

また、上述の実施形態において、プロセッサノード100、110、120を始めとする各プロセッサノードが、共有メモリ140に書き込まれた通信データのデータ量や、メモリウィンドウに含まれるリングバッファの空きの領域が例えば4byteの所定量まで減少した場合など、通知情報で通知する条件をユーザの要望に応じて、任意に設定することが可能であっても良い。

【0115】

また、本発明の上記機能は、C、C++、Java（登録商標）、Java（登録商標）Applet、Java（登録商標）Script、Perl、Rubyなどのレガシープログラミング言語、オブジェクト指向プログラミング言語などで記述された装置実行可能なプログラムにより実現でき、装置可読な記録媒体に格納して頒布することができる。

【0116】

これまで本発明を図1～図20に示した第1～第3の実施形態をもって説明してきたが、本発明はこれに限定されるものではない。他の実施の形態、追加、変更、削除など、本発明の要旨を変更しない範囲内で変更することができ、いずれの態様においても本発明の作用・効果を奏する限り、本発明の範囲に含まれるものである。

【符号の説明】

【0117】

10...マルチプロセッサシステム、100、110、120、200、210、220、300、310、320...プロセッサノード、130、230、330...Ethernet（登録商標）、140...共有メモリ、240...PCIバス、340...PCI-Expressスイッチ

【先行技術文献】

【特許文献】

【0118】

【特許文献1】特許第3743381号公報

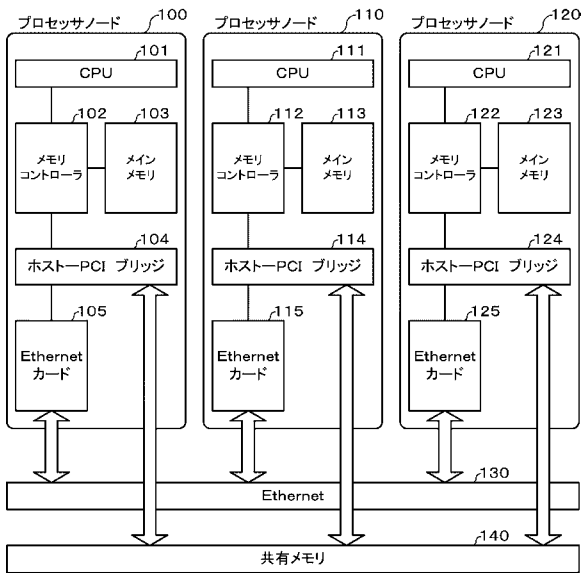
10

20

30

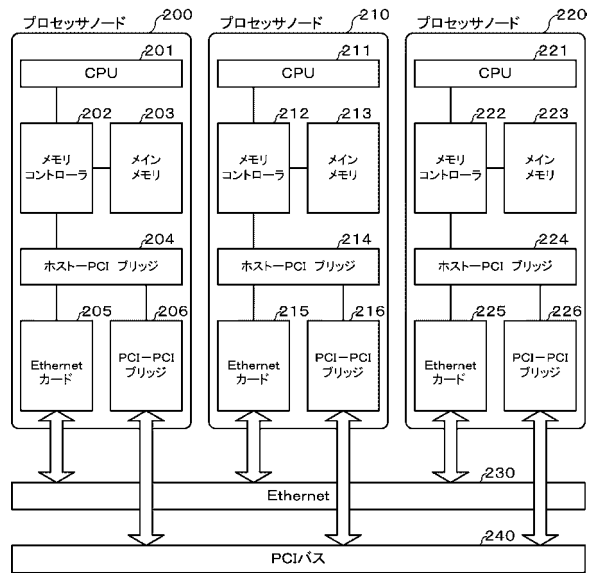
【 図 1 】

10

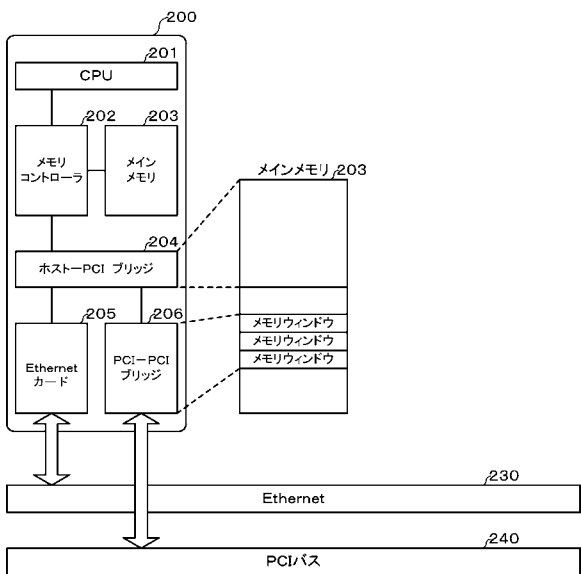


【 図 2 】

10

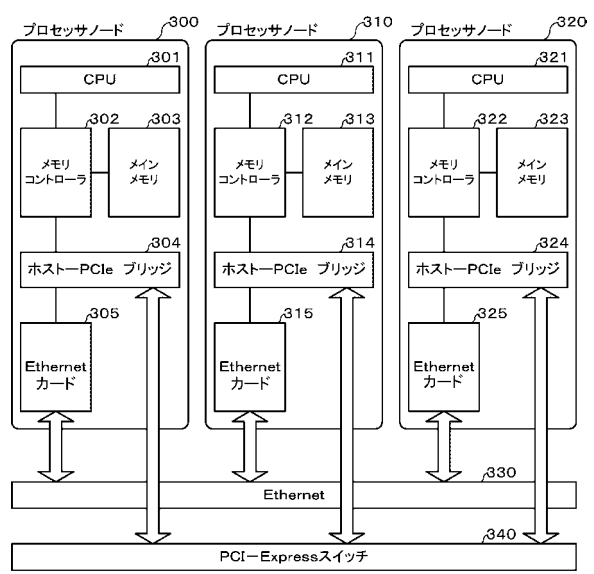


【 図 3 】

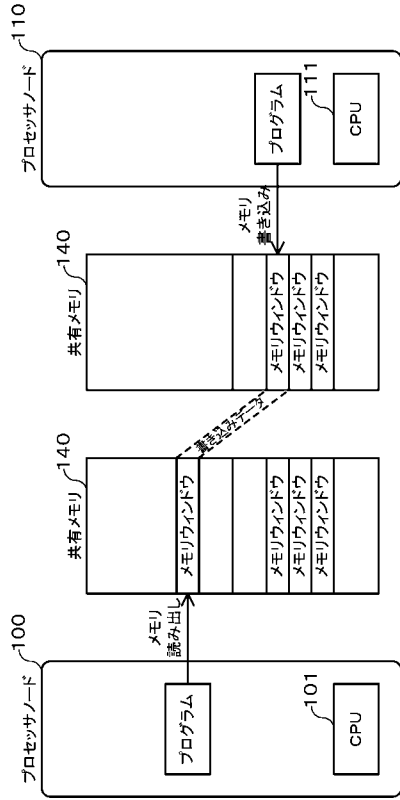


【 図 4 】

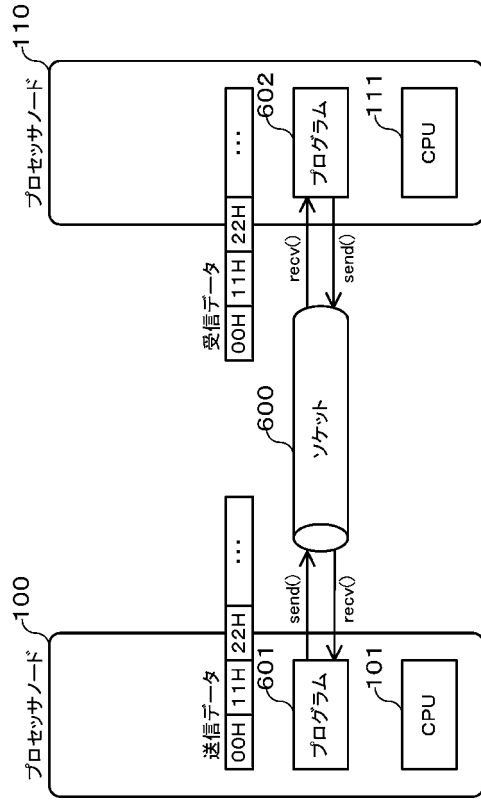
10



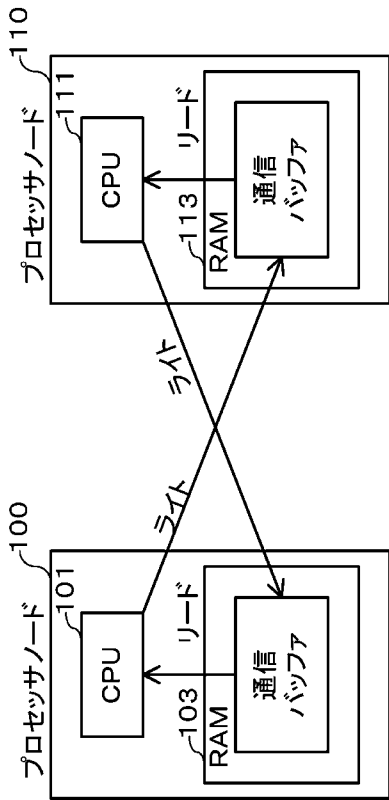
【 図 5 】



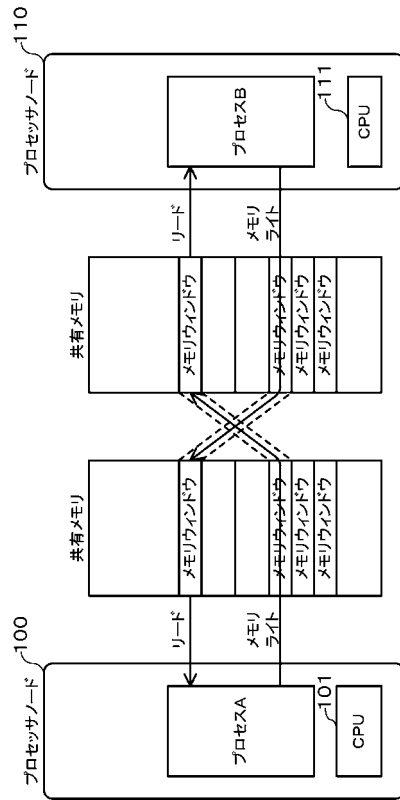
【 図 6 】



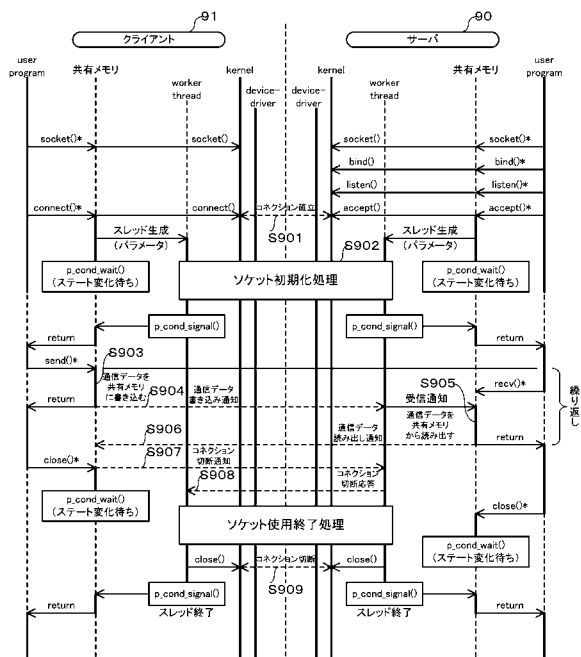
【 図 7 】



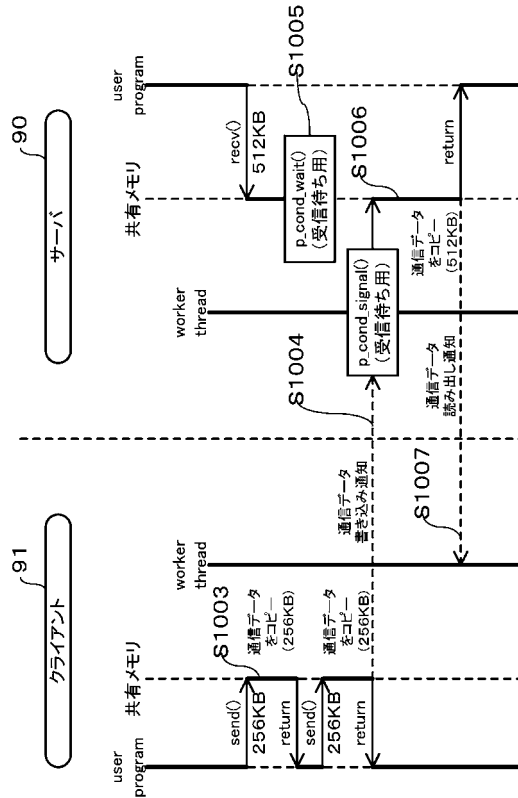
【 図 8 】



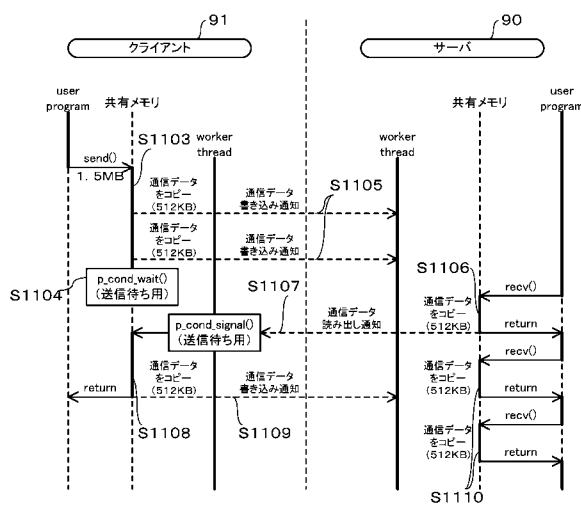
【図 9】



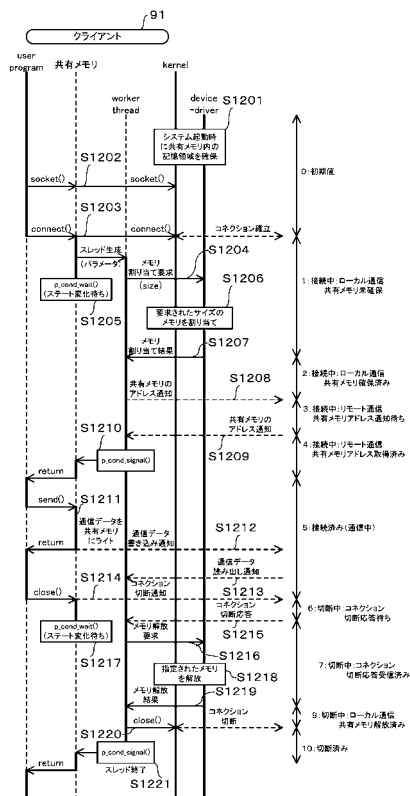
【図 10】



【図 11】



【図 12】

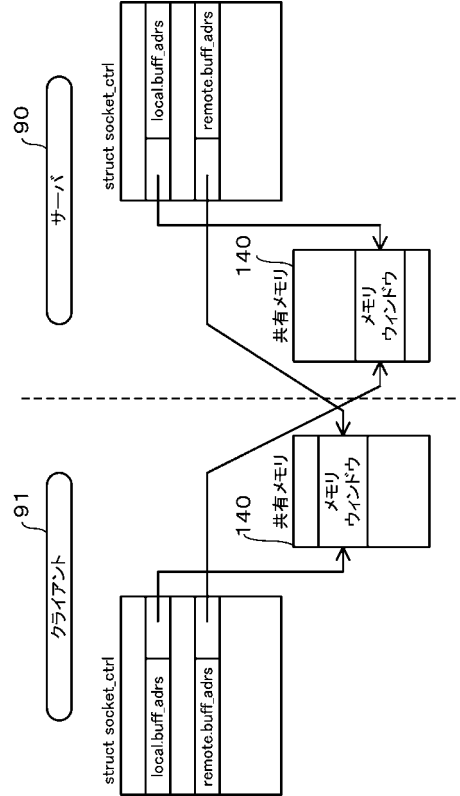


【 図 1 3 】

変数名	内容	値
sockfd	ソケットディスクリプタ	
state	ソケットステータス	
tx_write_notify_mode	送信データ書き込み通知モード	
rx_packet_count	通信パケットカウンタ	
tx_write_byte	送信データの書き込み済みバイト数	
rx_read_byte	送信データの読み込み済みバイト数	
local_buff_size	ローカル通信メモリ:メモリサイズ	
local_buff_adrs	ローカル通信メモリ:メモリアドレス	
local_ring_buff_size	ローカル通信バッファ:リングバッファサイズ	
local_ring_buff_adrs	ローカル通信バッファ:リングバッファアドレス	
remote_buff_size	リモート通信メモリ:メモリサイズ	
remote_buff_adrs	リモート通信メモリ:メモリアドレス	
remote_ring_buff_size	リモート通信バッファ:リングバッファサイズ	
remote_ring_buff_adrs	リモート通信バッファ:リングバッファアドレス	

1300

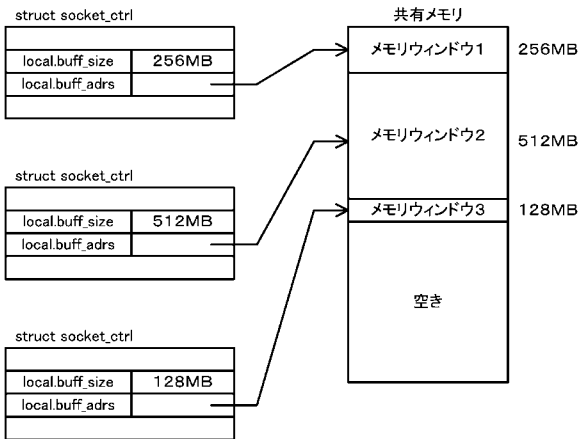
【 図 1 4 】



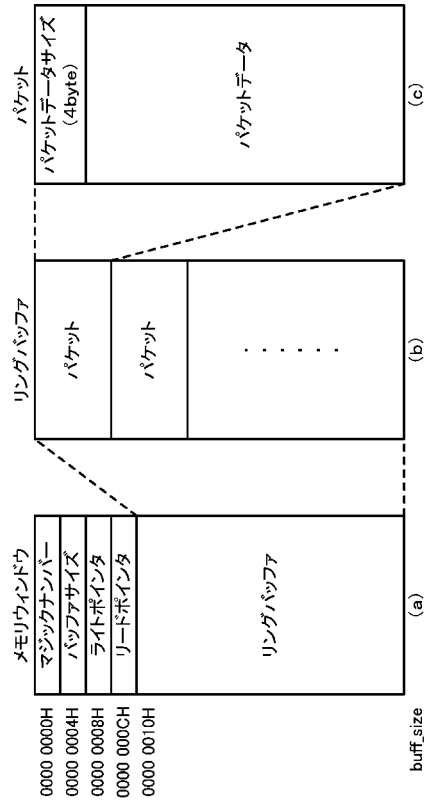
サーバ

クライアント

【 図 1 5 】



【 図 1 6 】



メモリウィンドウ

リングバッファ

パケット

パケット

パケットデータ

パケットデータサイズ (4byte)

リングバッファ

メモリウィンドウ

マジックナンバー

バッファサイズ

ライトポインタ

リードポインタ

0000 0000H

0000 0004H

0000 0008H

0000 000CH

0000 0010H

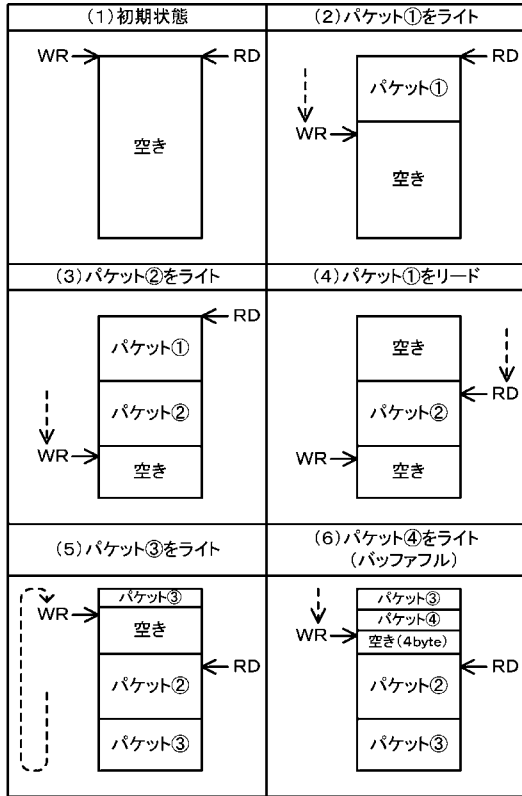
buf_size

(a)

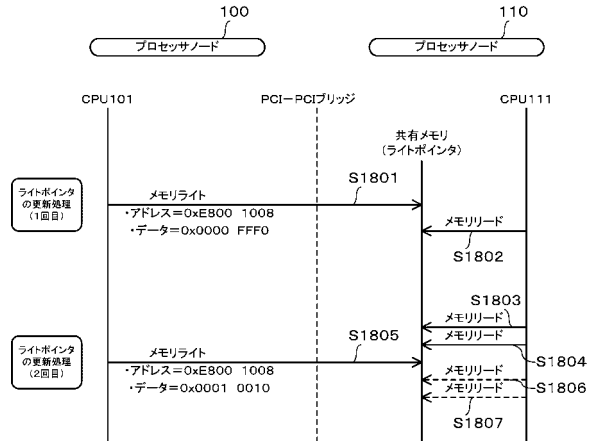
(b)

(c)

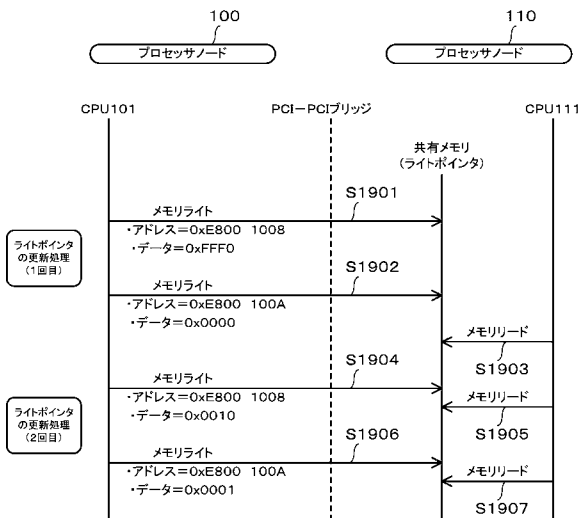
【 図 1 7 】



【 図 1 8 】



【 図 1 9 】



【 図 2 0 】

