

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第3601955号  
(P3601955)

(45) 発行日 平成16年12月15日(2004.12.15)

(24) 登録日 平成16年10月1日(2004.10.1)

(51) Int. Cl.<sup>7</sup>

F I

G06F 15/17

G06F 15/17

G06F 15/173

G06F 15/173

請求項の数 9 (全 25 頁)

(21) 出願番号	特願平9-290597	(73) 特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22) 出願日	平成9年10月23日(1997.10.23)	(74) 代理人	100068504 弁理士 小川 勝男
(65) 公開番号	特開平11-126196	(74) 代理人	100061893 弁理士 高橋 明夫
(43) 公開日	平成11年5月11日(1999.5.11)	(74) 代理人	100086656 弁理士 田中 恭助
審査請求日	平成16年3月30日(2004.3.30)	(72) 発明者	保田 淑子 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内
		(72) 発明者	藤井 啓明 東京都国分寺市東恋ヶ窪一丁目280番地 株式会社日立製作所中央研究所内 最終頁に続く

(54) 【発明の名称】 データ転送方法およびそれに適した計算機システム

(57) 【特許請求の範囲】

【請求項1】

相互結合ネットワークで接続された複数の要素プロセッサを有し、各要素プロセッサ内には、ユーザプロセスと交信するメッセージパッシングライブラリと、そのメッセージパッシングライブラリと交信するメモリ間直接転送ライブラリとが組み込まれ、各要素プロセッサは、プロセッサと、メモリと、上記相互結合ネットワークとの間でメッセージを交換するためのネットワークインタフェース回路とを有する計算機システムにおいて、送信側の要素プロセッサで走行中のユーザプロセスにより、その要素プロセッサ内のメッセージパッシングライブラリに対して発行されたデータ送信要求が要求する送信データと、そのデータ送信要求が要求する、メッセージパッシングライブラリにより定められた、受信側の要素プロセッサに送信されるべき付加情報との送信を、そのメッセージパッシングライブラリからその要素プロセッサ内のメモリ間直接転送ライブラリに対して要求し、そのメモリ間直接転送ライブラリにより、上記要求された送信データおよび付加情報と、それらの受信の完了を示す制御情報を書き込むべき領域を指定するための、受信側の要素プロセッサがあらかじめ決定した受信側アドレス情報とを含むメッセージの送信を上記ネットワークインタフェース回路に対して要求し、上記ネットワークインタフェース回路により、上記メッセージを組立て、受信側の要素プロセッサに宛てて上記相互結合ネットワークに送信し、上記メッセージの送信後に、そのメモリ間直接転送ライブラリが決定したアドレスを有する、上記メモリ内の記憶位置に送信完了を示す制御情報を上記ネットワークインタフェー

10

20

ス回路により書き込むステップからなるデータ転送方法。

【請求項 2】

受信側の要素プロセッサ内のネットワークインタフェース回路により、上記メッセージ内の送信データと付加情報を、受信側の要素プロセッサのメモリ内の、上記受信側アドレス情報により定まる領域に書き込み、

上記書き込みの終了後に、受信側の要素プロセッサ内のネットワークインタフェース回路により、受信完了を示す制御情報を、上記メモリ内の、上記受信側アドレス情報により定まる領域に書き込むステップをさらに有する請求項 1 記載のデータ転送方法。

【請求項 3】

上記データ送信要求は、上記送信データに関する第 1 のアドレス情報とデータ長情報と上記付加情報を指定し、

上記メッセージの送信を要求するステップは、

送信側の要素プロセッサ内の上記メモリ間直接転送ライブラリにより、上記要求された付加情報を、そのメモリ間直接転送ライブラリが付加情報と送信完了を示す制御情報を書き込むための領域として決定した、上記メモリ内の領域に書き込み、

そのメモリ間直接転送ライブラリにより、上記第 1 のアドレス情報、上記送信データ長情報、上記付加情報が書き込まれた上記領域を指定する第 2 のアドレス情報および上記受信側アドレス情報とを指定するメッセージ送信要求を送信側の要素プロセッサ内の上記ネットワークインタフェース回路に対して発行するステップからなり、

上記メッセージを送信するステップは、

上記ネットワークインタフェース回路により、そのメッセージ送信要求で指定された上記第 1 のアドレス情報と上記データ長情報との組および上記第 2 のアドレス情報にそれぞれ基づいて、上記送信データおよび上記付加情報を送信側の要素プロセッサのメモリより読み出し、

上記ネットワークインタフェース回路により、そのメッセージ送信要求が指定した上記受信側アドレス情報と上記データ長情報、上記読み出された送信データおよび付加情報を含むメッセージを生成し、上記ネットワークを介して受信側の要素プロセッサに送信するステップを有し、

上記送信完了を示す制御情報を書き込むステップは、上記第 2 のアドレス情報に基づいて、上記付加情報が記憶された領域と異なる記憶位置に送信完了を示す制御情報を書き込むステップを有する請求項 1 記載のデータ送信方法。

【請求項 4】

受信側の要素プロセッサ内のネットワークインタフェース回路により、上記メッセージ内の送信データと付加情報を、受信側の要素プロセッサのメモリ内の、上記受信側アドレス情報により定まる領域に書き込み、

上記書き込みの終了後に、受信側の要素プロセッサ内のネットワークインタフェース回路により、受信完了を示す制御情報を、上記メモリ内の、上記受信側アドレス情報により定まる領域に書き込むステップをさらに有し、

上記受信側のアドレス情報は、上記メッセージ内の送信データを書き込むべき第 3 のアドレス情報と受信完了を示す制御情報を書き込む第 4 のアドレス情報からなり、

上記メッセージ内の送信データと付加情報を書き込むステップは、上記第 3、第 4 のアドレス情報に基づいて、上記送信データおよび上記付加情報を、受信側の要素プロセッサのメモリに書き込むステップからなり、

上記受信完了を示す制御情報を書き込むステップは、上記第 4 のアドレス情報に基づいて、受信側の要素プロセッサのメモリ内の、上記付加情報が書き込まれた領域と異なる領域に書き込むステップからなる請求項 3 記載のデータ転送方法。

【請求項 5】

送信側の要素プロセッサで走行中の他のユーザプロセスによりその要素プロセッサ内の上記メモリ間直接転送ライブラリに対して発行された他のデータ送信要求が要求する送信データと、受信側の要素プロセッサがあらかじめ決定した、受信データと受信完了を示す制

10

20

30

40

50

御情報を書き込むべき領域を指定するための受信側アドレス情報とを含むメッセージの送信を、そのメモリ間直接転送ライブラリにより、送信側の要素プロセッサの上記ネットワークインタフェース回路に対して要求し、

上記メッセージを上記ネットワークインタフェース回路により組立て、受信側の要素プロセッサに宛てて上記相互結合ネットワークに送信し、

上記メッセージの送信後に、そのメモリ間直接転送ライブラリが決定したアドレスを有する、上記メモリ内の記憶位置に送信完了を示す制御情報を上記ネットワークインタフェース回路により書き込み、

受信側の要素プロセッサ内のネットワークインタフェース回路により、上記メッセージ内の送信データを、受信側の要素プロセッサのメモリ内の、上記受信側アドレス情報により 10  
定まる領域に書き込み、

上記書き込みの終了後に、受信側の要素プロセッサ内のネットワークインタフェース回路により、受信完了を示す制御情報を、上記メモリ内の、上記受信側アドレス情報により定まる領域に書き込むステップをさらに有する請求項 2 記載のデータ転送方法。

【請求項 6】

相互結合ネットワークで接続された複数の要素プロセッサを有し、

各要素プロセッサは、プロセッサと、メモリと、上記相互結合ネットワークとの間でメッセージを交換するためのネットワークインタフェース回路とを有し、

上記ネットワークインタフェース回路は、

転送すべきメッセージに関する情報として上記プロセッサにより供給される、第 1 のアドレス情報とデータ長情報との組および第 2 のアドレス情報とに基づいて、それぞれ送信されるべき第 1 のデータおよびそのデータとともに送信されるべき第 2 のデータを上記メモリから読み出すメモリアクセス回路と、 20

読みだされた第 1、第 2 のデータと、転送すべきメッセージに関する他の情報として上記プロセッサにより供給される第 3、第 4 のアドレス情報を含む一つのメッセージを生成し、上記相互結合ネットワークに送信する回路とを有し、

上記メモリアクセス回路は、上記メッセージの送信後に、上記第 2 のアドレス情報に基づいて、上記送信完了を示す制御情報を上記メモリ内の、上記第 2 のデータが記憶されている記憶位置と異なる記憶位置に書き込み、

該第 1 のアドレス情報は、上記第 1 のデータを保持する、上記メモリ内の領域のアドレス 30  
を指示し、上記データ長情報は、上記第 1 のデータの長さを指定し、該第 2 のアドレス情報は、上記第 2 のデータを保持し、かつ、該第 1 のデータの送信完了を示す制御情報をさらに保持すべき、上記メモリ内の領域のアドレスを指示し、該第 3 のアドレス情報は、第 1 のデータを書き込むべき、受信側の要素プロセッサのメモリ内の領域のアドレスを指示し、該第 4 のアドレス情報は、上記第 1 のデータおよび該第 1 のデータの受信完了を示す制御情報とを格納するための、受信側の要素プロセッサのメモリ内の領域のアドレスを指示する計算機システム。

【請求項 7】

上記ネットワークインタフェース回路は、他の要素プロセッサから送信されたメッセージを上記相互結合ネットワークから受信する回路をさらに有し、 40

上記メモリアクセス回路は、受信されたメッセージ内の上記第 3 のアドレス情報に基づいて、受信されたメッセージ内の上記第 1 のデータを上記メモリに書き込み、受信されたメッセージ内の上記第 4 のアドレス情報に基づいて、受信されたメッセージ内の上記第 2 のデータを上記メモリに書き込み、上記受信されたメッセージ内の上記第 1、第 2 のデータの書き込みの終了後に、上記第 4 のアドレス情報に基づいて、上記第 1 のデータの受信完了を示す制御情報を、上記メモリ内の、受信された上記第 2 のデータが記憶されている記憶位置と異なる記憶位置に書き込む回路を有する請求項 6 記載の計算機システム。

【請求項 8】

上記転送制御情報はモードビットをさらに有し、

上記メモリアクセス回路は、上記モードビットが第 1 の値の時に、上記第 2 のデータの読 50

み出しを実行し、上記モードビットが第2の値の時に、上記第2のデータの読み出しを実行せず、

上記メッセージ送信回路は、上記モードビットが第1の値の時に、上記第2のデータを含むメッセージを生成し、上記モードビットが第2の値の時に、上記第2のデータを含まないメッセージを生成し、

上記メモリアクセス回路は、上記モードビットが第1の値の時には、上記メッセージの送信後に、上記第2のアドレス情報に基づいて、上記送信完了を示す制御情報を、上記メモリ内の、上記第2のデータが記憶されている記憶位置と異なる記憶位置に書き込み、上記モードビットが第2の値の時には、上記メッセージの送信後に、上記送信完了を示す制御情報を、上記第2のアドレス情報に依存する、上記メモリ内の記憶位置に書き込む回路を有する請求項6記載の計算機システム。

10

【請求項9】

上記メッセージ送信が生成するメッセージは、上記モードビットを含み、

上記ネットワークインタフェース回路は、他の要素プロセッサから送信されたメッセージを上記相互結合ネットワークから受信する回路をさらに有し、

上記メモリアクセス回路は、

受信されたメッセージ内のモードビットが第1の値の時には、受信されたメッセージ内の上記第3のアドレス情報に基づいて、受信されたメッセージ内の上記第1のデータを上記メモリに書き込み、上記受信されたメッセージ内の上記第4のアドレス情報に基づいて、受信されたメッセージ内の上記第2のデータを上記メモリに書き込み、上記受信されたメッセージ内の上記第1、第2のデータの書き込みの終了後に、上記第4のアドレス情報に基づいて、上記第1のデータの受信完了を示す制御情報を、上記メモリ内の、受信された上記第2のデータが記憶されている記憶位置と異なる記憶位置に書き込み、

20

受信されたメッセージ内のモードビットが第2の値の時には、受信されたメッセージ内の上記第1のデータを、受信されたメッセージ内の上記第3のアドレス情報に基づいて、上記メモリ内の記憶位置に書き込み、上記受信されたメッセージ内の上記第1のデータの書き込みの終了後に、上記第1のデータの受信完了を示す制御情報を、上記第4のアドレス情報に依存する、上記メモリ内の記憶位置に書き込む回路を有する請求項8記載の計算機システム。

【発明の詳細な説明】

30

【0001】

【発明の属する技術分野】

本発明は、相互結合ネットワークを介して接続された複数の要素プロセッサ間でのデータ転送方法およびそれに適した計算機システムに関する。

【0002】

【従来の技術】

従来、並列計算機は、ローカルメモリと命令プロセッサから構成される複数の要素プロセッサを相互結合ネットワークで結合した構成をとっている。一般に、このような形態の並列計算機は、分散メモリ型並列計算機と呼ばれる。各要素プロセッサは、相互結合ネットワークを介して個々のローカルメモリに格納されているデータの授受を行い、並列に処理を実行する。

40

【0003】

一般に、分散メモリ型の並列計算機では、メッセージパッシングと呼ぶプログラミングモデルを用いてデータ転送を実現する。メッセージパッシングモデルでは、ユーザが並列プログラム中に明示的に送信(S E N D)手続きおよび受信(R E C E I V E)手続きを記して、要素プロセッサ間で必要になるデータの授受をメッセージのやりとりという形で行う。命令プロセッサはこれらの通信手続きを解析して、相互結合ネットワークにデータを送信したり、相互結合ネットワークからデータを受信しながら処理を進める。送信元の要素プロセッサは、転送先の要素プロセッサ番号を指定してメッセージを転送し、転送先の要素プロセッサでメッセージをバッファリングする。メッセージパッシングモデルでは、メ

50

ッセージ通信に伴ってデータのバッファリングやフロー制御が必要となり、送受信オーバーヘッドが大きくなってしまふ。

**【 0 0 0 4 】**

この送受信オーバーヘッドを削減するために、近年複数の並列計算機において、要素プロセッサのローカルメモリの内容を、メッセージの生成あるいは受信を行う送受信回路が直接アクセスするメモリ間直接転送方法が使用されている。この方法を実行する代表例は、PUT / GET通信である。例えば、"情報処理学会並列処理シンポジウム J S P P ' 9 5、" P P . 2 3 3 - 2 4 0 ( 1 9 9 5 年 5 月 ) 参照。メモリ間直接転送方法では、各要素プロセッサのローカルメモリからの送信データの読み出しあるいはそのメモリへの受信データの書き込みをメッセージの生成あるいは受信を行う送受信回路が直接実行するため、これらのデータをOS管理の領域にコピーする必要がなく、このコピーに由来するオーバーヘッドを削減できる。

10

**【 0 0 0 5 】**

しかしながら、このようなPUT / GET通信を実際に行う通信ライブラリは、各並列計算機メーカーや研究機関が独自に開発しているため、それをを用いて作成された並列プログラムを他機種へ移植することは困難であった。この問題点を解決するために、MPI ( Message Passing Interface ) に代表される、メッセージパッシングライブラリの標準化が進みつつある。MPIは、米国各大学および並列計算機メーカーがメッセージパッシングインタフェース標準化団体MPI Forumを組織し、その研究成果をまとめた仕様である。この仕様に基づいて作成されたライブラリ(以下MPIライブラリと呼ぶことがある)は、今後の並列プログラム開発支援ライブラリの主流になると考えられる。MPIライブラリを用いて記述された並列プログラムは、異機種間で変更なしに走らせることができる。各計算機メーカーは、自社の並列計算機上で高性能を達成するようにMPIライブラリを開発している。

20

**【 0 0 0 6 】**

上記MPI仕様は、PUT / GET通信に関する仕様を含んでいない。しかし、データ転送の高速化のためには、PUT / GET通信を併用することが重要である。このために、並列計算機メーカー等は、PUT / GET通信ライブラリを使用可能にしたMPIライブラリを開発している。たとえば、本出願人による、「並列計算機SR22D支援ライブラリ」参照。したがって、MPIライブラリは、各計算機メーカーごとに異なるものであるが、ユーザプログラムから見れば、MPIライブラリとの間のインタフェースは、いずれの計算機メーカーのMPIライブラリに対しても同じである。従って、そのユーザプログラムは、いずれの計算機のメーカーの上でも実行できることになる。

30

**【 0 0 0 7 】**

MPIライブラリを用いるユーザプログラムは、データを送信する時点で、MPIライブラリをコールする。従来は、ユーザプログラムがPUT / GET通信ライブラリを使用するときには、ユーザプログラムは送信すべきユーザデータおよびデータ長をこのコール文の引数をもって指定すればよい。しかし、MPIライブラリを使用するときには、ユーザプログラムは、このコール文の中でこのMPIライブラリにより定められた付加情報を引数としてさらに指定する必要がある。この付加情報は、送信先のプロセスの識別子、プロセスグループ識別子等を含む固定長のデータであり、メッセージの送信先の要素プロセッサにおいて、受信したメッセージがそこで実行中のユーザプロセスが発行する受信要求が要求したメッセージか否かの識別に使用される。以下、この情報をMPI付加情報とも呼ぶ。従来のMPIライブラリとPUT / GETライブラリを併用した通信方法では、ユーザデータおよびMPI付加情報を異なる二つのメッセージにより転送していた。

40

**【 0 0 0 8 】****【 発明が解決しようとする課題 】**

従来のように、同じ転送先プロセッサに対しユーザデータおよびMPI付加情報を2つのメッセージとして転送する場合、各メッセージの転送に異なる転送制御情報が必要になる。その結果、これらの情報の生成も2度行わなければならない。このために、従来のMP

50

ライブラリを使用したデータ転送では、ユーザプログラムがデータの転送を要求してから、実際に転送が開始されるまでの遅延時間（転送レイテンシと呼ばれる）が大きい。

【0009】

さらに、これらの2種のデータに対して別々のメッセージとして送信処理、受信処理を行うと、メッセージ数に比例してローカルメモリに対するアクセス回数（転送制御情報の読み出し、転送データの読み出し、フラグの書き込み）が増加してしまう。

【0010】

本発明の目的は、以上の問題を減少させ、より高速にデータを転送できるデータ転送方法を提供することにある。

【0011】

【発明を解決するための手段】

上記の目的を達成するために、本発明によるデータ転送方法は、送信側の要素プロセッサで走行中のユーザプロセスから発行されたデータ送信要求が要求する、送信データとそれに関連する付加情報とを、送信元のメッセージパッシングライブラリから送信元のメモリ間直接転送ライブラリに通知し、

そのメモリ間直接転送ライブラリにより、上記送信データおよび付加情報と、それらの受信の完了を示す制御情報を書き込むべき領域を指定するための、受信側の要素プロセッサがあらかじめ決定した受信側アドレス情報とを含むメッセージの送信をネットワークインタフェース回路に対して要求し、

上記ネットワークインタフェース回路により、上記メッセージを組立て、受信側の要素プロセッサに宛てて上記相互結合ネットワークに送信し、

上記メッセージの送信後に、そのメモリ間直接転送ライブラリが決定したアドレスを有する、上記メモリ内の記憶位置に送信完了を示す制御情報を上記ネットワークインタフェース回路により書き込む。

【0012】

より具体的には、上記メモリ間直接転送ライブラリは、上記付加情報を上記メモリ内の領域に書き込み、ユーザが指定した送信データの記憶位置を示す第1のアドレスとそのデータの長さ、上記付加情報の記憶位置を示す第2のアドレスと、送信先の要素プロセッサのメモリにおける、送信データの記憶位置を指定する第3のアドレスと、上記付加情報を記憶する位置を示す第4のアドレス等を指定し、これらの情報を含むメッセージの送信を上記ネットワークインタフェース回路に要求する。

【0013】

この回路は、上記第1のアドレスと上記データ長により送信データを読み出し、上記第2のアドレスにより付加情報を読み出し、これらの送信データ、付加情報および上記第3、第4のアドレスを含むメッセージを生成し、送信先にあてて送信する。この送信の完了後に、上記第2のアドレスを使用して、送信完了を示す制御情報をメモリ内の、上記付加情報の書き込み位置と異なる位置、具体的には、付加情報が書き込まれた記憶位置の次の記憶位置に書き込む。

【0014】

さらに、受信側の要素プロセッサにおいても、ネットワークインタフェース回路が、上記メッセージ内の上記送信データと付加情報とをそれぞれ上記第3、第4のメモリアドレスが指定する記憶位置に書き込むとともに、この書き込みの完了後に、受信完了を示す制御情報を上記第4のアドレスの基づいて上記付加情報の書き込み位置と異なる記憶位置に書き込む。

【0015】

本発明のより具体的な態様では、ユーザプロセスが、メッセージパッシングライブラリを介して行う上記の転送とともに、他のユーザプロセスはメッセージパッシングライブラリを介さないでメモリ間直接転送ライブラリにデータの送信要求を発行することができるようになっている。この場合には、メモリ間直接転送ライブラリとネットワークインタフェース回路は、上に述べた処理における、付加情報が存在しない場合の処理と基本的に同じ

10

20

30

40

50

処理を行う。

【0016】

【発明の実施の形態】

以下、本発明に係る計算機システムを図面に示した実施の形態を参照してさらに詳細に説明する。

【0017】

<発明の実施の形態>

図1に、本発明における並列計算機の概略構成を示す。図中、101～104は並列計算機を構成する要素プロセッサ、105は相互結合ネットワークである。要素プロセッサ101～104は、相互結合ネットワーク105に接続し、相互結合ネットワーク105を介して、要素プロセッサ間でデータの授受を行う。相互結合ネットワーク105の構成方法(トポロジ)は、クロスバ結合、格子結合、リング結合、多段結合等多種存在するが、本発明は、これらのいずれにも適用可能であり、特定の相互結合ネットワークトポロジに限定されない。図3に要素プロセッサ101の概略構成を示す。図中、301は命令プロセッサ、302はキャッシュローカルメモリ、303はストレージコントローラ、304はローカルメモリ、305はネットワークインタフェース回路、306はI/Oインタフェース回路である。この並列計算機は、個々の要素プロセッサがそれぞれ固有のローカルメモリ304を有する分散ローカルメモリ型の並列計算機である。

10

【0018】

各要素プロセッサは、他の要素プロセッサとの間でメッセージパッシングによる通信を実行するように構成されている。すなわち、各要素プロセッサは、標準のメッセージパッシングインタフェース、たとえばMPIを有するメッセージパッシングライブラリ(以下、MPIライブラリと呼ぶ)と、このライブラリと交信して自要素プロセッサ内のローカルメモリとの間で直接データの授受を行うメモリ間直接転送を実行するライブラリとして、PUT/GET型通信を実行するためのライブラリ(以下、PUT/GET型ライブラリと呼ぶ)と、PUT/GET型ライブラリからのコマンドにより起動されるネットワークインタフェース回路305を有している。なお、本発明は、この特定のメッセージパッシングライブラリに限定されるのではなく、他のメッセージパッシングライブラリたとえばPVM、PARMACSとして知られているライブラリも適用できる。

20

【0019】

本実施の形態では、各要素プロセッサ内のユーザプロセスがMPIライブラリに対してデータ送信要求を発行したときに、MPIライブラリ、PUT/GET型ライブラリおよびネットワークインタフェース回路305は、協同してユーザデータとMPI付加情報を一つのメッセージにて転送し、さらに他の要素プロセッサからユーザデータとMPI付加情報を含むメッセージを受信したときに、ネットワークインタフェース回路305は、これらのデータを区分してローカルメモリに書き込むところに特徴がある。この付加情報は、MPIライブラリを介した通信のために使用されるもので、MPIライブラリが指定した形式の、データ送信に関連する複数の情報からなり、それぞれの情報は、ユーザプロセスにより指定される。具体的には、既に例示したように、この付加情報は、受信側のユーザプロセスの識別子、プロセスグループの識別子等を含む。

30

40

【0020】

より具体的には、各要素プロセッサで実行中のユーザプロセスからMPIライブラリに対する送信要求が発行されたときに、MPIライブラリは、その送信要求が指定するユーザデータと付加情報の送信をPUT/GET型ライブラリに要求する。PUT/GET型ライブラリにその送信要求が指定する付加情報をローカルメモリに書き込み、そのユーザデータとMPI付加情報の両方を一つのメッセージとして転送するための転送制御情報を生成し、ローカルメモリ304に書き込み、その後転送制御情報によるユーザデータおよび付加情報の送信をネットワークインタフェース回路305に要求する。

【0021】

ネットワークインタフェース回路305は、この送信要求に回答して、転送制御情報に従

50

ってユーザデータとM P I付加情報を含む一つのメッセージを組み立て、受信側の要素プロセッサに相互結合ネットワーク105を介して転送する。受信側の要素プロセッサでは、ネットワークインタフェース回路305は、このメッセージを受信すると、メッセージに含まれたユーザデータおよびM P I付加情報をメッセージのヘッダ内の転送制御情報が指定するローカルメモリ内の二つのアドレスに書き込み、それぞれを受信側のユーザプロセッサおよび受信側のM P Iライブラリに引き渡す。

#### 【0022】

図2に転送制御情報の例を示す。転送制御情報200には、G E TあるいはP U T動作の場合に使用される転送先プロセッサ番号201、P U T動作の場合に送信されるユーザデータが格納されているローカルメモリ領域の先頭アドレスである送信データアドレス203、P U T動作の場合に送信完了フラグを書き込むローカルメモリ領域の先頭アドレスである送信フラグアドレス204、G E TあるいはP U T動作の場合に使用される、転送されるデータの長さである転送データ長205、受信側の要素プロセッサにおいて受信データを書き込むローカルメモリ領域の先頭アドレスである受信データアドレス206、その要素プロセッサにおいて、その受信データに対する受信完了フラグを書き込むローカルメモリ領域の先頭アドレスである受信フラグアドレス207、その他通信処理に必要な情報208等を格納する。

#### 【0023】

さらに、本実施の形態では、モードビット202がセットされていない場合、送信フラグアドレス204および受信フラグアドレス207は、それぞれP U T動作時の送信完了フラグおよびG E T動作時の受信完了フラグを書き込むローカルメモリアドレスを指定する。しかし、モードビット202がセットされている場合、送信フラグアドレスフィールド204は、P U T動作時にM P I付加情報を読み出すべきローカルメモリアドレスを指定するのに使用され、受信フラグアドレスフィールド207は、G E T動作時に受信したM P I付加情報を書き込むべきローカルメモリアドレスを指定するのに使用される。

#### 【0024】

この結果、P U T動作時に送信完了フラグを書き込むべきローカルメモリアドレスが転送制御情報200により指定されなくなるが、本実施の形態では、あらかじめセットされたM P I付加情報サイズを送信フラグアドレスに加算し、その結果得られるアドレスにユーザデータおよびM P I付加情報という二つのデータの送信完了フラグを書き込む。同様に、G E T動作時には、G E T動作時には、あらかじめセットされたM P I付加情報サイズを受信フラグアドレスフィールドに加算し、その結果得られるアドレスに2種類のデータの受信完了フラグを書き込む。これにより、1つの転送制御情報200を用いてユーザデータとM P I付加情報という2つの種類のデータを1つのメッセージで送信または受信し、従来と同様に送信完了フラグあるいは受信完了フラグもローカルメモリに書き込むことができる。

#### 【0025】

命令プロセッサ301は、プログラム処理を行うユニットである。キャッシュローカルメモリ302は、命令プロセッサ301に付随する、高速かつ小容量のローカルメモリである。ネットワークインタフェース回路305は、相互結合ネットワーク105に接続し、命令プロセッサ301からの指示に従って、ローカルメモリ分散型の並列計算機の特徴であるデータ転送処理を命令プロセッサ301のプログラム処理とは独立して行うユニットである。ストレージコントローラ303は、命令プロセッサ301、ネットワークインタフェース回路305およびI/Oインタフェース回路306から発行されるデータアクセス要求に従って、適当な記憶媒体にアクセスする。ローカルメモリ304は、ストレージコントローラ303で制御され、データ等を格納する。命令プロセッサ301およびネットワークインタフェース回路305は独立に動作するため、ストレージコントローラ303は、命令プロセッサ301からローカルメモリ304へのアクセス要求を処理すると同時に、ネットワークインタフェース回路305からのデータ転送に伴うローカルメモリ304へのアクセスも処理する。I/Oインタフェース回路306は、ストレージコント

10

20

30

40

50



ローラ 303 からのアクセス要求に従って、I/O 装置にアクセスする。I/O インタフェース回路 306 は、要素プロセッサの構成によっては、存在しない場合もある。

【0026】

図 4 に示すように、ストレージコントローラ 303 は、命令プロセッサインタフェース回路 401、アドレス解析部 402、メモリアクセスインタフェース回路 403 およびデータ転送インタフェース回路 404 で構成される。命令プロセッサインタフェース回路 401 は、命令プロセッサ 301 からローカルメモリ 304 へのアクセスおよび命令プロセッサ 301 からネットワークインタフェース回路 305 へのコマンド発行というトランザクションを線 401S から受け取る。通常、このコマンドは、ネットワークインタフェース回路 305 内部の制御レジスタへのアクセス要求である。命令プロセッサインタフェース回路 401 は、このトランザクションへの返答、ストレージコントローラ 303 やネットワークインタフェース回路 305 で検出した割り込み要因を線 402S を介して命令プロセッサ 301 へ伝える。このトランザクションの応答は、たとえば、ローカルメモリからの読み出しデータである。

10

【0027】

アドレス解析部 402 は、命令プロセッサ 301 が発行した、ローカルメモリアクセス要求およびネットワークインタフェース回路 305 へのコマンドを線 403S を介して受け取り、そのアクセス要求あるいはコマンドが指定するアクセス先アドレスを解析する。ローカルメモリアクセス要求は線 404S を介してメモリアクセスインタフェース回路 403 に伝えられる。また、ネットワークインタフェース回路 305 へのコマンドは、線 406S を介してデータ転送インタフェース回路 404 に伝達される。

20

【0028】

メモリアクセスインタフェース回路 403 は、アドレス解析部 402 からのローカルメモリアクセス要求を線 404S を介して受け、線 407S を介してローカルメモリ 304 に伝達する。ローカルメモリアクセス要求がローカルメモリからの読み出し要求であった場合、この読み出し要求が指定するデータがローカルメモリ 304 から線 408S を介して伝達される。読み出しデータは、メモリアクセスインタフェース回路 403 から線 409S を介して命令プロセッサインタフェース回路 401 に伝達され、線 402S を介して命令プロセッサ 301 に伝達される。また、ローカルメモリアクセスインタフェース 403 は、データ転送に関わるローカルメモリアクセスも処理する。データ転送処理に関わるローカルメモリアクセスは、データ転送インタフェース回路 404 から線 410S を介して伝達される。メモリアクセスインタフェース回路 403 は、アドレス解析部 402 からローカルメモリアクセス要求が伝達された時と同様に、ローカルメモリアクセス要求をローカルメモリ 304 に対して発行し、読み出しアクセスに対しては、読み出しデータを線 411S を介してデータ転送インタフェース回路 404 に返送する。

30

【0029】

データ転送インタフェース回路 404 は、アドレス解析部 402 から線 406S を介して伝達されるネットワークインタフェース回路 305 へのコマンドを受け取り、線 412S を介してネットワークインタフェース回路 305 に伝達する。ネットワークインタフェース回路 305 からは、線 413S を介してコマンドに対する返答およびデータ転送に関わるローカルメモリアクセス要求が伝達される。データ転送インタフェース回路 404 は、前記コマンドに対する返答を線 414S を介して命令プロセッサインタフェース回路 401 に伝達し、線 402S を介して命令プロセッサ 301 に伝達する。ローカルメモリアクセス要求は、線 410S を介してメモリアクセスインタフェース回路 403 に伝達する。ローカルメモリ読み出しデータは、線 411S を介してデータ転送インタフェース回路 404 に伝達され、データ転送インタフェース回路 404 から線 412S を介してネットワークインタフェース回路 305 に伝達される。データ転送インタフェース回路 404 は、ネットワークインタフェース回路 305 内部で発生した割り込み伝達要求を受ける場合もある。この場合、割り込み伝達要求は線 414S を介して命令プロセッサインタフェース回路 401 に伝達され、さらに命令プロセッサ 301 に伝達される。

40

50

## 【 0 0 3 0 】

ネットワークインタフェース回路 3 0 5 は、コマンド受信部 4 0 5、コマンド処理部 4 0 6、メッセージ生成部 4 0 7、メッセージ送信部 4 0 8、メッセージ受信部 4 0 9、メッセージ分解部 4 1 0 およびコマンド送信部 4 1 1 で構成される。コマンド受信部 4 0 5 は、ストレージコントローラ 3 0 3 から線 4 1 2 S を介して、ネットワークインタフェース回路 3 0 5 内部の制御レジスタへのアクセスあるいはネットワークインタフェース回路 3 0 5 が要求したローカルメモリ 3 0 4 から読み出されたデータ等を受け取る。このデータは、線 4 1 5 S を介してコマンド処理部 4 0 6 に伝達され、転送データとして使用されたり、ネットワークインタフェース回路 3 0 5 の動作制御用データとしてネットワークインタフェース回路内部の制御レジスタに設定されたり、データ送信時にメッセージ生成用データ（転送先プロセッサ番号、送信データアドレス、送信フラグアドレス、転送データ長、受信データアドレス、受信フラグアドレス等）として使用される。

10

## 【 0 0 3 1 】

コマンド処理部 4 0 6 は、線 4 1 5 S を介してコマンド受信部 4 0 5 から伝達されるネットワークインタフェース回路 3 0 5 内部の制御レジスタへのアクセスを行う。制御レジスタ読み出しアクセスを受けた場合、コマンド処理部 4 0 6 は、読み出し結果を線 4 1 7 S を介してコマンド送信部 4 1 1 に伝達する。また、制御レジスタ書き込みアクセスを受けた場合には、コマンド処理部 4 0 6 は、その書き込みを実行する。メッセージの送信処理は、上記の制御レジスタへのアクセスがメッセージ送信起動用レジスタへの書き込み要求である場合に開始する。メッセージの送信処理では、データを宛先の要素プロセッサに転送するのに必要な情報であるヘッダを作成したり、転送データが存在するローカルメモリアドレスを知るために、ローカルメモリ 3 0 4 に格納されている転送制御情報 2 0 0 を読み出すローカルメモリアクセス要求が発生される。このアクセス要求は線 4 1 7 S を介してコマンド送信部 4 1 1 に伝達され、線 4 1 3 S、データ転送インタフェース回路 4 0 4、線 4 1 0 S、メモリアクセスインタフェース回路 4 0 3、線 4 0 7 S を介してローカルメモリ 3 0 4 から読み出される。ローカルメモリ 3 0 4 からの読み出し結果は、線 4 0 8 S、メモリアクセスインタフェース回路 4 0 3、線 4 1 1 S、データ転送インタフェース回路 4 0 4、線 4 1 2 S、コマンド受信部 4 0 5 を介してコマンド処理部 4 0 6 に伝達され、メッセージ生成部 4 0 7 に伝達される。

20

## 【 0 0 3 2 】

メッセージ生成部 4 0 7 は、本実施の形態での特徴的な回路の一つであり、4 1 8 S を介して伝達された転送データと転送制御情報 2 0 0 を含むヘッダからメッセージを生成し、線 4 1 9 S を介してメッセージ送信部 4 0 8 に送出する。転送制御情報 2 0 0 内のモードビット 2 0 2 が 1 である場合、メッセージ生成部 4 0 7 は、転送制御情報 2 0 0 から生成されるヘッダと、転送制御情報 2 0 0 内の送信データアドレス 2 0 3 に従ってローカルメモリ 3 0 4 から読み出した送信データと、送信フラグアドレス 2 0 4 に従ってローカルメモリ 3 0 4 から読み出した M P I 付加情報からメッセージを組み立て、メッセージ送信部 4 0 8 に送出する。

30

## 【 0 0 3 3 】

図 5 に、メッセージ生成部 4 0 7 の内部構成を示す。メッセージ生成部 4 0 7 は、ローカルメモリ 3 0 4 から読み出した転送制御情報 2 0 0 内の、送信データアドレス、送信フラグアドレスフィールド内のアドレス、転送データ長をそれぞれ保持するレジスタ 5 0 1、5 0 2、5 0 3 を有する。メッセージ生成部 4 0 7 は、さらに、ローカルメモリ 3 0 4 から生成中のメッセージのためにローカルメモリからすでに読みだされたデータの量を保持するレジスタ 5 0 4 と、本実施の形態に特徴的な回路として、M P I 付加情報のサイズを保持するレジスタ 5 0 5 と、転送制御情報 2 0 0 内モードビットを保持するレジスタ 5 0 6 を有する。レジスタ 5 0 4 に保持された読み出し済みのデータの量はデータがローカルメモリ 3 0 4 から読み込まれるたびにカウントアップされ、送信すべきすべてのデータがメッセージ送信部 4 0 8 に伝達された後リセットされる。従って、この読み出し済みのデータの量は、送信済みのデータの総量と考えることができる。M P I 付加情報のサイズは

40

50

、ユーザプロセスからの初期化要求によりM P IライブラリとP U T / G E Tライブラリが初期化される時にレジスタ5 0 5にあらかじめセットされる。レジスタ5 0 6には、転送制御情報2 0 0の読み出し時にその情報内のモードビットがセットされ、すべてのデータがメッセージ送信部4 0 8に伝達された後このレジスタ内のモードビットがリセットされる。メッセージ生成部4 0 7は、さらに、ローカルメモリ読み出し要求発行部5 1 1と、ローカルメモリ書込み要求発行部5 1 2と、メッセージ組み立て部5 1 3の他に、本実施の形態に特徴的な回路として、アンドゲート5 0 7と、加算器5 0 8および5 0 9と比較回路5 1 0とを有する。

#### 【0034】

アンドゲート5 0 7は、モードビットが1である場合には、レジスタ5 0 5内のM P I付加情報サイズを出力し、0である場合にはモードビットの値0を出力する。加算器5 0 9は、アンドゲート5 0 7の出力とレジスタ5 0 2に保持されている送信フラグアドレスフィールドの値を加算する。P U T動作時にモードビットが1にセットされている場合、レジスタ5 0 2に保持されている送信フラグアドレスフィールドには、ローカルメモリ3 0 4内の、M P I付加情報7 1 4の先頭アドレスが含まれているので、この加算の結果アドレスは、そのM P I付加情報の次のアドレスを指すことになり、P U T動作時の送信完了フラグの書込みアドレスとして使用される。加算器5 0 8は、アンドゲート5 0 7の出力と、レジスタ5 0 3に保持された転送データ長を加算する。この加算の結果は、モードビットに1がセットされた場合、P U T動作時に送信すべきメッセージに含まれるべき、ローカルメモリ3 0 4から読み出すべきデータの総量を示す。

#### 【0035】

比較回路5 1 0は、レジスタ5 0 4に保持された読み出し済みのデータ量を、レジスタ5 0 3に保持された転送データ量と加算器5 0 8から出力される転送データの総量とを比較する。比較の結果として、レジスタ5 0 4に保持された読み出し済みデータ量がレジスタ5 0 3内の転送データ長を越えていない場合には、そのことを示す比較結果信号をローカルメモリ読み出し要求発行部5 1 1に出力する。レジスタ5 0 4に保持された読み出し済みデータ量がレジスタ5 0 3内の転送データ長を越えているが、加算器5 0 8から与えられる、ローカルメモリ3 0 4から読み出すべきデータの総量を越えていない場合には、そのことを示す比較結果信号をローカルメモリ読み出し要求発行部5 1 1に出力する。レジスタ5 0 4に保持された読み出し済みデータ量が読み出すべきデータの総量に達した場合、そのことを示す比較結果信号をローカルメモリ書込み要求発行部5 1 2に出力する。

#### 【0036】

ローカルメモリ読み出し要求発行部5 1 1は、比較回路5 1 0からの比較結果信号と、レジスタ5 0 1内の送信データアドレスおよびレジスタ5 0 2内の送信フラグアドレスとからローカルメモリ読み出し要求を生成し、コマンド送信部4 1 1に送信する。すなわち、比較結果信号が、読み出し済みのデータが転送データ長を超えていないことを示すときには、送信データアドレスを元に後続の未読み出しのユーザデータを読み出すためのローカルメモリ読み出し要求を生成し、読み出し済みのデータが転送データ長を越えているが、読み出すべきデータの総量を越えていないときには、送信フラグアドレスを元にして未読み出しのM P I付加情報を読み出すためのローカルメモリ読み出し要求を生成する。

#### 【0037】

ローカルメモリ書込み要求発行部5 1 2は、比較回路5 1 0からの比較結果信号と、加算器5 0 9から出力される送信完了フラグアドレスとを元にローカルメモリ書込み要求をコマンド送信部4 1 1に送信する。すなわち、比較結果信号が、読み出し済みデータ量が読み出すべきデータの総量に達したことを示す場合、加算器5 0 8より与えられる送信完了フラグアドレスに送信完了フラグを書き込むことを要求する書き込み要求を生成する。

#### 【0038】

メッセージ組み立て部5 1 3では、本実施の形態に特徴的なセレクタ5 1 4がモードビットの値に従ってメッセージを組み立て、メッセージ送信部4 0 8にそのメッセージの送信要求を送付する。モードビットが1である場合、ヘッダとデータとM P I付加情報を含む

10

20

30

40

50

メッセージ515を組立て、モードビットが0である場合には、ヘッダとデータのみからなるメッセージ516を組立る。

【0039】

図4において、メッセージ送信部408は、メッセージ生成部407からのメッセージ送信要求を受けて線420Sを介して相互結合ネットワーク105にメッセージを送出する。送出されたメッセージは、相互結合ネットワーク105を介してそのヘッダ情報に従って宛先に転送される。メッセージ生成部407における送信処理は、ネットワークインタフェース回路305へ送信される複数のメッセージ送信要求に対してそれらの送信要求の到着順に順次行われる。

【0040】

次に、メッセージ受信部409について説明する。メッセージ受信部409は、線421Sを介して、相互結合ネットワーク105からメッセージを順次受け取り、メッセージ分解部410に線422Sを介して転送する。メッセージ分解部410はこのメッセージのヘッダ部に含まれる転送制御情報に従ってこのメッセージをデータ部とヘッダ部に分解し、ローカルメモリ304へこのデータ部や受信完了フラグの書き込みを要求する書き込み要求を線423Sを介してコマンド送信部411に伝達する。

【0041】

図6にメッセージ分解部410の内部構成を示す。メッセージ分解部410には、メッセージヘッダ内の受信データアドレス、受信フラグアドレス、転送データ長、モードビットをそれぞれ保持するレジスタ601、602、603、606が設けられている。レジスタ606は本実施の形態で特徴的なレジスタであり、レジスタ606には、ヘッダ受信時にヘッダ内のモードビットがセットされ、メッセージ内の全データを受信したときにそのモードビットがリセットされる。さらに、受信したメッセージ内のデータの内、ローカルメモリに書き込み済みのデータの総量を保持するレジスタ604と、MPI付加情報のサイズを保持する、本実施の形態に特徴的なレジスタ605が設けられている。レジスタ605には、MPI初期化時あるいはジョブ起動時にあらかじめ定められたMPI付加情報サイズがセットされ、このサイズ情報はメッセージ内の全データが受信されたときにリセットされる。レジスタ604に保持された書き込み済みのデータの量は受信されたデータがローカルメモリ304に書き込まれるたびにカウントアップされ、受信すべきすべてのデータがローカルメモリ304に書き込まれた後リセットされる。従って、この書き込み済みのデータの総量は、受信済みのデータの総量であるとも考えることができる。

【0042】

メッセージ分解部410には、ローカルメモリ書き込み要求発行部611の他に、本実施の形態で特徴的な、アンドゲート607、加算器608、609と比較回路610とがさらに設けられている。アンドゲート607は、レジスタ606内のモードビットが1である場合にレジスタ605内のMPI付加情報サイズを出力し、モードビットが0である場合には0を出力する。加算器608および609の動作は、加算器508および509(図5)と同様である。比較回路610は、レジスタ504内の受信済みデータの総量を、転送データ長、アンドゲート607の出力結果データと比較し、比較結果信号をローカルメモリ書き込み要求発行部611に出力する。すなわち、この比較回路は、受信済みのデータが転送データ長よりも短いか、受信済みのデータの総量が転送データ長より大きいか、転送データ長とアンドゲート607から出力されるMPI付加情報のサイズの和より小さいか、あるいは受信済みのデータの総量が転送データ長とMPI付加情報サイズの和より大きいかを判別する。

【0043】

メモリ書き込み要求発行部611は、比較回路610の比較結果信号と、レジスタ601内の受信データアドレスと、加算器609の出力とレジスタ602内の受信フラグアドレスフィールドの値とから、ローカルメモリ書き込み要求を生成し、コマンド送信部411に線425Sを介して伝達する。加算器609は、モードビットが1である場合に、レジスタ602内の受信フラグアドレスフィールドの値にMPI付加情報サイズを加算したアドレス

10

20

30

40

50

を受信完了フラグを書き込むべきローカルメモリアドレスとして出力する。

【 0 0 4 4 】

比較回路 6 1 0 の出力が、受信済みのデータの総量が転送データ長よりも大きいことを示す場合、メモリ書き込み要求発行部 6 1 1 は、受信されたデータをそれまでに受信したデータの書き込み位置に続けて書き込むためのローカルメモリ書き込み要求をレジスタ 6 0 1 内の受信データアドレスに基づいて生成する。

【 0 0 4 5 】

比較回路 6 1 0 の出力が、受信済みのデータの総量が転送データ長よりも大きく転送データ長とアンドゲート 6 0 7 から出力される M P I 付加情報のサイズの和より小さいことを示す場合、新たに受信されたデータは、M P I 付加情報である。従って、この場合には、メモリ書き込み要求発行部 6 1 1 は、受信されたデータをそれまでに受信した M P I 付加情報の書き込み位置に続けて書き込むためのローカルメモリ書き込み要求を、レジスタ 6 0 2 内の受信フラグアドレスフィールドに含まれる、M P I 付加情報の書き込みアドレスに基づいて生成する。

10

【 0 0 4 6 】

比較回路 6 1 0 の出力が、受信済みのデータの総量が転送データ長とアンドゲート 6 0 7 から出力される M P I 付加情報のサイズの和より大きいことを示す場合、すべてのデータが受信されたこととなる。従って、この場合には、メモリ書き込み要求発行部 6 1 1 は、受信完了フラグを M P I 付加情報の書き込み位置に続けて書き込むためのローカルメモリ書き込み要求を、加算器 6 0 9 の出力に基づいて生成する。

20

【 0 0 4 7 】

コマンド送信部 4 1 1 は、これらの書き込み要求に従って、受信されたデータあるいは受信完了フラグをローカルメモリ 3 0 4 に書き込む。コマンド送信部 4 1 1 は、線 4 1 7 S を介して伝達される、ネットワークインタフェース回路 3 0 5 内部の制御レジスタからの読み出しデータ、線 4 2 4 S を介して伝達される、メッセージ送信処理において使用される送信データのローカルメモリ 3 0 4 からの読み出し要求、線 4 2 4 S を介して伝達される、メッセージ送信処理の完了に伴う送信完了フラグのローカルメモリ 3 0 4 への書き込み要求、線 4 2 3 S を介して伝達される、メッセージ受信処理に伴う受信データあるいは受信完了フラグのローカルメモリ 3 0 4 への書き込み要求およびネットワークインタフェース回路 3 0 5 内で発生した割り込み要求を、線 4 1 3 S を介してストレージコントローラ 3 0 3 内のデータ転送インタフェース回路 4 0 4 に伝達する。また、コマンド送信部 4 1 1 は、ネットワークインタフェース回路 3 0 5 の動作制御に関わる情報のローカルメモリ 3 0 4 からの読み出し要求、メッセージ送信処理に使用する転送制御情報 2 0 0 のローカルメモリ 3 0 4 からの読み出し要求を線 4 1 3 S を介してデータ転送インタフェース回路 4 0 4 に伝達する。

30

【 0 0 4 8 】

次に本実施の形態におけるメッセージ転送の流れを説明する。最初に自プロセッサのローカルメモリに格納されているユーザデータおよび M P I 付加情報を転送先の要素プロセッサのローカルメモリに直接書き込む P U T 処理について図 4、5、7 を用いて説明する。まず、本実施の形態に係わる通信方式では、送信側のユーザプロセスおよび受信側のユーザプロセスは、M P I ライブラリを使用する前に、M P I ライブラリ内の初期化ルーチンたとえば M P I \_ i n i t をコールするコマンド発行する。この初期化ルーチンのコールを受けると、M P I ライブラリは、P U T / G E T ライブラリ内のいくつかの通信準備手続きをコールする。これらの通信準備手続きは、使用する P U T / G E T ライブラリにより予め定められているが、以下では、後の説明に関連する部分および本実施の形態で新規に行われる処理のみを説明する。

40

【 0 0 4 9 】

本実施の形態では、送信側の P U T / G E T ライブラリおよび受信側の P U T / G E T ライブラリは、いずれもこれらの通信準備手続きにおいて以下の処理をすると仮定する。すなわち、ローカルメモリ 7 0 1 をユーザ空間にあらかじめマップし、さらに、通信領域 7 0

50

3、704(図7)を確保する。さらに、それぞれの通信領域内に送信データ領域およびそれに対応する送信完了フラグ領域を確保する。図7では、713は送信データ領域の例を示す。図では、受信側の通信領域704内の送信データ領域は図示していない。本実施の形態では、送信完了フラグ領域として、MPI付加情報714および送信完了フラグ715の両方を格納する連続した領域を確保する点で従来と異なる。同様に、各通信領域内に、受信データ領域および受信完了フラグ領域を確保する。本実施の形態では図7では、721は受信データ領域の例を示す。図では、送信側の通信領域703内の受信データ領域は図示していない。本実施の形態では、受信完了フラグ領域として、MPI付加情報722および受信完了フラグ723の両方を格納する連続した領域を確保する点で従来と異なる。なお、MPI付加情報714の長さは例えば64バイト程度である。なお、送信完了フラグあるいは受信完了フラグを書き込む領域715,723は例えば8バイトである。

10

#### 【0050】

PUT/GET型ライブラリとして、送信すべきユーザデータを、ローカルメモリ701に常駐させることを前提とする場合とそうでない場合とがある。前者の場合には、送信データ領域713は、この常駐されたユーザデータの領域と一致するように、送信データ領域713が決定される。一方、後者の場合には、送信データ領域713は、送信側のユーザプロセスが使用するユーザデータに割り当てられたローカルメモリ内の領域とは独立に決定される。本発明はいずれの構造のPUT/GET型ライブラリにも適用可能である。しかし、後者の場合には、後に述べるように、ユーザプロセスが使用しているユーザデータに割り当てられたローカルメモリ内の領域のデータを送信データ領域713にコピーする処理が必要となる。しかし、前者の場合にはこのコピー動作が必要でなく、それだけデータ転送動作が高速化される。

20

#### 【0051】

以上のようにして、通信準備手続きが実行された後に、送信側のユーザプロセスの処理が進むと、そのユーザプロセスは、送信すべきデータを送信側の通信領域703内のユーザデータ領域713に書き込んだ後に、データ送信要求コマンド、たとえばMPI\_sendを送信側のMPIライブラリに対して発行する。このコマンドの名称は、使用するメッセージパッシングライブラリにより定まり、それが指定する引数も同様にそのライブラリにより定められた複数の種類の情報からなる。ここで仮定するMPIライブラリの場合には、このコマンドの引数は、送信すべきユーザデータの先頭アドレス、ユーザデータ長と付加情報からなり、この先頭アドレスは、ユーザプロセスに割り当てられた仮想メモリ空間内での、そのユーザデータに対する仮想アドレスである。この付加情報は、受信側のユーザプロセスの識別子、プロセスグループの識別子等を含む。最初の二つの引数は、MPIライブラリを介さないでデータ転送をユーザプロセスがPUT/GET型ライブラリに直接要求するためのデータ送信要求が指定する引数と同じであり、付加情報がMPIライブラリに対するデータ転送要求が新たに指定する引数である。

30

#### 【0052】

送信側のMPIライブラリは、このデータ送信要求コマンドMPI\_sendに応答して、送信側のPUT/GET型ライブラリに、この送信要求が指定するデータの送信を要求する。送信側のMPIライブラリは、この要求を、MPIライブラリとPUT/GET型ライブラリにより予め定められた一つまたは複数のコマンドの形で発行する。以下では、それらのコマンドの内、本実施の形態で使用すると仮定する主なコマンドのみを説明する。

40

#### 【0053】

まず、送信側のMPIライブラリは、送信権の取得を要求するコマンドを発行する。送信側のPUT/GET型ライブラリは、このコマンドに응答して、受信側のユーザプロセスと交信してそのプロセスに対するデータの送信権を得る。受信側のユーザプロセスおよびそのプロセスを実行している要素プロセッサの番号は、上記データ送信要求コマンドMPI\_sendが指定する付加情報中の、受信側のプロセス識別番号とプロセスグループ識

50

別番号とにより決定される。

【 0 0 5 4 】

送信側の M P I ライブラリは、さらに、受信側のユーザプロセスの受信データ領域および受信フラグ領域のそれぞれの先頭位置を示す受信データアドレスおよび受信フラグアドレスを受信する。但し、この後に再度同じデータ送信要求コマンド M P I \_ s e n d を送信側のユーザプロセスが発行したときには、このコマンドを実行する必要はない。

【 0 0 5 5 】

既に述べたように、P U T / G E T 型ライブラリが、送信すべきユーザデータがローカルメモリ 7 0 1 に常駐されることを前提としない場合には、送信側の M P I ライブラリは、送信すべきユーザデータを、ユーザデータ領域 7 1 3 にコピーすることを要求するコマンドを発行し、送信側の P U T / G E T 型ライブラリにより、このコマンドの引数で指定されるユーザデータのアドレスとデータ長で指定されるユーザデータに割り当てられた、ローカルメモリ内の領域のデータを、先に決定された送信データ領域 7 0 3 にコピーする。P U T / G E T 型ライブラリが、送信すべきユーザデータをローカルメモリ 7 0 1 に常駐させることを前提とする場合には、このコピー動作は不要である。

10

【 0 0 5 6 】

次に、送信側の M P I ライブラリは、付加情報のローカルメモリへの書き込みを要求するコマンドを発行する。送信側の P U T / G E T 型ライブラリは、このコマンドに回答して、このコマンドの引数で指定される M P I 付加情報を、先に決定された送信データ領域 7 0 3 に対応して決定された送信フラグ領域の先頭の領域 7 1 4 に書き込む。

20

【 0 0 5 7 】

送信側の M P I ライブラリは、付加情報のローカルメモリへの書き込みを要求するコマンドを発行する。送信側の P U T / G E T 型ライブラリは、このコマンドに回答して、このコマンドの引数で指定される M P I 付加情報を、先に決定された送信データ領域 7 0 3 に対応して決定された送信フラグ領域の先頭の領域 7 1 4 に書き込む。

【 0 0 5 8 】

送信側の M P I ライブラリは、転送制御情報を生成することを要求するコマンドを発行する。送信側の P U T / G E T 型ライブラリは、このコマンドに回答して、転送制御情報 7 0 0 を生成して、ローカルメモリ 7 0 0 内の適当な領域に書き込む。この転送制御情報 7 0 0 に含まれた情報は以下の通りである。転送先プロセッサ番号 2 0 1 は、受信側のユーザプロセスが実行されているプロセッサの番号であり、この番号は、すでに述べたように、M P I 初期化ルーチンにおいて決定されている。モードビット 2 0 2 は、送信フラグアドレスフィールド 2 0 4 を本実施の形態に従って拡張して使用するか否かを示すビットである。送信側の P U T / G E T 型ライブラリは、送信側の M P I ライブラリから、データ転送を要求されたときに、P U T / G E T 型ライブラリはこのモードビットを 1 にセットする。モードビット 2 0 2 が 1 であることは、ユーザデータと M P I 付加情報とを一つのメッセージで送信することを指示する。なお、本実施例の形態では、送信側のユーザプロセスが、M P I ライブラリに対してでなく、P U T / G E T 型ライブラリに対して直接データ送信要求コマンドを発行した場合には、P U T / G E T 型ライブラリは、そのデータ要求に対して、図 7 に示す転送制御情報 7 0 0 と同じ構造を有し、モードビット 2 0 2 の値が 0 である転送制御情報を生成する。

30

40

【 0 0 5 9 】

送信データアドレス 2 0 3 は、送信側のユーザプロセスにより転送が要求されたユーザデータまたはそのコピーを保持する送信データ領域 7 1 3 の先頭アドレスである。転送データ長 2 0 5 は、データ送信要求コマンド M P I \_ s e n d が指定した、ユーザデータのデータ長であり、そのコマンドを受けた M P I ライブラリが P U T / G E T 型ライブラリに通知する。

【 0 0 6 0 】

送信フラグアドレスフィールド 2 0 4 には、従来では送信完了フラグを書き込むためのローカルメモリアドレスが格納されるが、本実施の形態では、モードビット 2 0 2 が 1 であ

50

る場合には、ローカルメモリ701内のMPI付加情報714の先頭アドレスをこの送信フラグアドレスフィールド204に格納する。なお、モードビット202が1の場合には、送信完了フラグを格納するローカルメモリ内の領域715のアドレスは送信フラグアドレスフィールド204によっては明には指定されないことになる。本実施の形態では、ネットワークインタフェース回路305が、ユーザプロセスとMPI付加情報とに対する共通の送信完了フラグを、ローカルメモリ701内のMPI付加情報714の最終のアドレスの次のアドレスの領域715に格納するようになっている。このため、送信型のPUT/GET型ライブラリは、ユーザプロセスからの先のデータ送信要求コマンドに対する応答として、この送信完了フラグが書き込まれた時点で、送信完了を送信側のユーザプロセスに通知するようになっている。なお、モードビット202が0である場合には、PUT/GET型ライブラリは、送信フラグアドレスフィールド204に、そのライブラリが決定した送信完了フラグを書き込むアドレスをセットする。

10

#### 【0061】

受信データアドレス206は、受信側の要素プロセッサにおいて、受信したデータを格納するためのローカルメモリ領域721(図7)のアドレスである。受信フラグアドレスフィールド207には、この受信データとともに受信した付加情報と受信完了を示す受信完了フラグを格納する領域の先頭アドレスである。本実施の形態では、付加情報を記憶する領域722の後続の領域723に受信完了フラグを書き込む。したがって、受信フラグアドレスフィールド207には、この領域722の先頭アドレスが書き込まれる。受信データアドレス206と受信フラグアドレス207は、いずれもMPI初期化ルーチンにて送信先のPUT/GET型ライブラリにより通知される。これらのアドレスは、GET動作時に送信側のPUT/GET型ライブラリにより使用される。

20

#### 【0062】

その他制御情報208は、PUTメッセージあるいはGETメッセージあるいはGET要求メッセージの種別を示したり、1対1通信あるいは1対多通信等の通信形態を示すといった、その他通信処理に必要な情報を含む。

#### 【0063】

なお、転送制御情報200の中に、MPI付加情報を読み出すべきローカルメモリアドレスを指定するフィールドおよび受信したMPI付加情報を書き込むべきローカルメモリアドレスを指定するフィールドを別に設けることも可能であるが、本実施の形態のようにモードビット202を用いて、送信フラグアドレス204および受信フラグアドレス207が指定する二つのアドレスを切り替えることにより転送制御情報200の構造と大きさを、ユーザプロセスがPUT/GET型ライブラリのみを使用して送信あるいは受信する場合と同じとすることができる。

30

#### 【0064】

こうして、転送制御情報700が生成されると、送信側のMPIライブラリは、生成された転送制御情報に従って、データの送信を行うことをネットワークインタフェース回路305に要求するコマンドを発行する。PUT/GET型ライブラリは、このコマンドに回答して、ネットワークインタフェース回路305内部の送信起動用レジスタ(図示せず)にその転送制御情報700のアドレスを書き込むことを要求する書き込みコマンドを発行する。コマンド処理部406は、この書き込みコマンドを実行して、ネットワークインタフェース回路305内のメッセージ送信起動用レジスタ(図示せず)へ転送制御情報700の先頭アドレスを書き込む。この書き込みによりネットワークインタフェース回路305はメッセージ送信処理を開始する。

40

#### 【0065】

図4に示すように、コマンド処理部406は、メッセージのヘッダを作成するためにローカルメモリ304に格納されている転送制御情報700を読み出すローカルメモリアクセス要求を線417Sを介してコマンド送信部411に伝達する。線413S、データ転送インタフェース回路404、線410S、メモリアクセスインタフェース回路403を介して、送信起動用レジスタ(図示せず)に書き込まれたアドレスをもとにローカルメモリ

50



304から読み出された転送制御情報700は、メモリアクセスインタフェース回路403、線411S、データ転送インタフェース回路404、線412S、コマンド受信部405を介してコマンド処理部406に伝達され、メッセージ生成部407に伝達される。

【0066】

図5において、メッセージ生成部407は、転送制御情報700内の送信データアドレス203、送信フラグアドレスフィールド204の値、転送データ長205、モードビット202をそれぞれレジスタ501、502、503、506にセットする。その他の情報は図示しないレジスタに保持される。レジスタ506内のモードビットは今の場合には1である。レジスタ506内のモードビットが1であるため、アンドゲート507がONになり、レジスタ505内のMPI付加情報サイズを出力する。加算器508は、レジスタ503内の転送データ長とアンドゲート507から出力されたMPI付加情報を足しあわせ、ローカルメモリ701から読み出すべきデータの総量を出力する。比較回路510がレジスタ504に保持された、読み出し済みデータの総量がレジスタ503内の転送データ長よりも小さいと判断した場合には、ローカルメモリ701内の、レジスタ501に保持された送信データアドレスの記憶位置からユーザデータ713を読み出すために、コマンド送信部411にローカルメモリ読み出し要求を伝達する。このユーザデータは、コマンド送信部411からストレージコントローラ303内のメモリアクセスインタフェース回路403を介してローカルメモリ304から読み出される。読み出されたデータは、ストレージコントローラ303内のメモリアクセスインタフェース回路403、データ転送インタフェース回路404およびコマンド受信部405を介してコマンド処理部406に

10

20

【0067】

メッセージ生成部407では、比較回路510が、レジスタ504内の読み出し済みのデータの量が、レジスタ503内の転送データ長より大きい、(転送データ長+レジスタ505内のMPI付加情報サイズ)以下であると判断したときには、メモリ読み出し要求発行部511は、ローカルメモリ701内の、レジスタ502に保持された送信フラグアドレスを有する記憶位置からMPI付加情報714を読み出すためのローカルメモリ読み出し要求をコマンド送信部411に伝達する。この読み出し要求は、コマンド送信部411からストレージコントローラ303内の、データ転送インタフェース回路404およびメモリアクセスインタフェース回路403を介してローカルメモリ304に送られ、MPI付加情報714がそこから読み出される。読み出されたMPI付加情報714は、ストレージコントローラ303内のメモリアクセスインタフェース回路403、データ転送インタフェース回路404およびコマンド受信部405を介してコマンド処理部406に伝達され、線418Sを介してメッセージ生成部407に伝達される。

30

【0068】

メッセージ生成部407では、メッセージ組み立て部513内のセレクタ514は、モードビットが1であることから、ローカルメモリ304から読み出された転送制御情報700内の送信データアドレス、送信フラグアドレス以外の部分をレジスタ515内のヘッダ部の格納する。同様に、ユーザデータ713およびMPI付加情報714をレジスタ515のデータ部に格納する。

40

【0069】

比較回路610が、読み出したユーザデータの総量が、転送データ長とMPI付加情報サイズの和に等しくなったことを検出したとき、メモリ書き込み要求発行部512は、送信完了フラグを書き込むことを要求するローカルメモリ書き込み要求をコマンド送信部411に伝達する。このコマンドは、加算器509により与えられる、レジスタ602内の送信フラグアドレスフィールドに保持されたMPI付加情報の先頭アドレスと、レジスタ505内のMPI付加情報サイズとの和に等しいアドレスにこのフラグを書き込むことを要求する。送信完了フラグ715は、コマンド送信部411からストレージコントローラ303内のメモリアクセスインタフェース回路403を介してローカルメモリ304に書き込まれる。本実施の形態では、メッセージに含まれるべきユーザデータとMPI付加情報の

50

読み出しが完了した時点で、メッセージの送信が完了したと見なして、送信完了フラグ 715 を書き込む。しかし、このメッセージが実際に相互結合ネットワーク 105 に送信された時点でこのフラグを書き込むようにしてもよい。

**【0070】**

こうして、レジスタ 515 内にユーザデータとそれに関連する M P I 付加情報を含む一つのメッセージ 705 が生成される。このメッセージ 705 には、転送制御情報 700 に含まれていたのと同じ転送先プロセッサ番号 716、モードビット 717、転送データ長 718、受信データアドレス 719、受信フラグアドレス 720、その他の制御情報 208 をそのまま含み、転送制御情報 700 に含まれていた送信データアドレス 203 と送信フラグアドレス 204 に代えて、ユーザデータ 726、付加情報 727 を含むことになる。送信データアドレスメッセージ生成部 407 はそのメッセージ 705 をメッセージ送信部 408 に送信する。メッセージ生成部 407 は、メッセージ 705 の生成に使用される上記 3 つの情報をメッセージ送信部 408 へ送出し終ると、メッセージ送信部 408 はレジスタ 504 と 506 をリセットする。メッセージ送信部 408 はそのメッセージ 705 を相互結合ネットワーク 105 に送出する。相互結合ネットワーク 105 はそのメッセージ内の転送先プロセッサ番号 201 により指定されるプロセッサにそのメッセージを転送する。

10

**【0071】**

次に受信処理について説明する。相互結合ネットワーク 105 から転送されたメッセージ 705 は、まずメッセージ受信部 409 で受け取られ、メッセージ分解部 410 に転送される。メッセージ分解部 410 は、メッセージヘッダ内の受信データアドレス 719、受信フラグアドレス 720、転送データ長 718、モードビット 717 をそれぞれレジスタ 601、602、603、606 (図 6) に書き込む。レジスタ 604 はあらかじめ 0 にリセットされ、レジスタ 605 にはあらかじめ M P I 付加情報サイズがセットされている。メッセージ分解部 410 では、メモリ書き込み要求発行部 611 は、受信されたユーザデータをヘッダ内の受信データアドレス 719 が示すローカルメモリ領域に書き込むためのローカルメモリアクセス要求を生成し、線 423 S を介してコマンド送信部 411 に伝達する。受信されたユーザデータは、メッセージヘッダ内の受信データアドレス 719 に従って、ストレージコントローラ 303 内のメモリアクセスインタフェース回路 403 を介して、ローカルメモリ 304 の通信領域 704 内の領域 721 に書き込まれる。レジスタ 606 内のモードビットは 1 であるため、アンドゲート 607 が ON になり、レジスタ 605 内の M P I 付加情報サイズを加算器 608 に供給する。加算器 608 では、レジスタ 603 内の転送データ長とアンドゲート 607 から与えられる M P I 付加情報サイズを足しあわせ、ローカルメモリに書込むべきデータの総量を得る。レジスタ 604 内の受信データ量は、メッセージ受信部 409 が相互結合ネットワーク 105 からメッセージ内のデータの異なる部分を受信することに更新される。

20

30

**【0072】**

メッセージ分解部 410 では、ユーザデータの異なる部分がメッセージ受信部 409 により受信されるごとに、比較回路 610 が、レジスタ 604 内の受信されたデータの総量がレジスタ 603 内の転送データ長よりも小さいか否かを判断し、前者が後者より小さいと判断したときには、メモリ書き込み要求発行部 611 は、レジスタ 601 に保持された受信データアドレスにしたがって受信されたデータを受信側のローカルメモリ 702 に書き込むことを要求するコマンド送信部 411 にローカルメモリ書き込み要求を伝達する。

40

**【0073】**

その後比較回路 610 が、レジスタ 604 内の、受信されたデータの総量がレジスタ 603 内の転送データ長よりも大きいと判断したならば、メモリ書き込み要求発行部 611 は、レジスタ 602 に保持された受信フラグアドレスフィールドの値のアドレスに、受信された M P I 付加情報を書き込むことを要求するローカルメモリ書き込み要求を伝達する。

**【0074】**

50

加算器 609 は、レジスタ 602 内の受信フラグアドレスフィールドの値とアンドゲート 607 から与えられる M P I 付加情報サイズを加算し、受信完了フラグを書き込むべきメモリアドレスを決定する。比較回路 610 が、レジスタ 604 内の受信データ数が、加算器 608 より与えられる、転送データ長 + M P I 付加情報サイズに等しくなったことを検出すると、メモリ書き込み要求発行部 611 は、加算器 609 により与えられるアドレスに受信完了フラグを書き込むことを要求するローカルメモリ書き込み要求をコマンド送信部 411 に伝達する。その結果、受信完了フラグ 723 は、コマンド送信部 411 から、ストレージコントローラ 303 内のメモリアクセスインタフェース回路 403 を介してローカルメモリ 304 に書き込まれる。こうして、データ受信処理が完了する。また、受信の完了でもってデータ転送処理が終了する。

10

**【0075】**

なお、受信されたデータおよび付加情報は、以下のようにして受信側の要素プロセッサ 702 で使用される。受信側のユーザプロセスが、他の要素プロセッサから送信されたデータの受信を要求するコマンド、例えば、M P I \_ r e c v を発行する。このコマンドは、受信すべきユーザデータを指定するアドレスと、そのデータの最大長、および M P I ライブラリにより定められた付加情報とからなる引数を指定する。この付加情報は、送信元ユーザプロセスの識別子その他の情報からなる。このコマンドで指定される上記アドレスは、受信側のユーザプロセスに割り当てられたアドレス空間に属する仮想アドレスである。受信側の M P I ライブラリは、この受信コマンドの引数で指定されるユーザデータがローカルメモリ 704 に書き込み済みであるか否かを、ローカルメモリ 704 に書き込まれたユーザデータ 721、付加情報 722、受信完了フラグ 723 に基づいて判別する。もし、この要求されたデータがローカルメモリ 704 に書き込み済みであるときには、受信側の M P I ライブラリは、受信側のユーザプロセスに受信完了を通知する。この要求されたデータがローカルメモリ 704 に書き込み済みでないときには、M P I ライブラリは、上記判別が成功するまでその判別を繰り返す。

20

**【0076】**

なお、受信側のユーザプロセスが指定した仮想アドレスを有するデータ領域が、ローカルメモリ 704 に常駐していない場合には、受信側の M P I ライブラリは、上記通知を行う前に、受信されたユーザデータ 721 を、ユーザプロセスが指定するアドレスに割り当てられたローカルメモリ領域にコピーする。もし、受信側のユーザプロセスが指定した仮想アドレスを有するデータ領域が、ローカルメモリ 704 に常駐している場合には、このコピーは不要である。

30

**【0077】**

受信側のユーザプロセスは、この受信完了の通知を受けると、受信データを読み出す命令を実行する。したがって、この受信側でのユーザデータの受信判別処理では、ユーザデータ 721 と付加情報 722 が書き込み済みであるか否かを検出するのに、共通の受信完了フラグ 723 を使用するところが従来と異なる。

**【0078】**

以上は、ユーザプロセスがデータの送信要求コマンドを M P I ライブラリに対して発行した場合である。本実施の形態では、他のユーザプロセスは、P U T / G E T ライブラリに対してデータ送信要求を発行することもできるようになっている。この場合には、このデータ送信要求コマンドは、送信すべきユーザデータを示す仮想アドレスと、データ長を指定する。送信側の M P I ライブラリは、このコマンドに回答して、先に述べたと同じようにして、転送制御情報 700 を生成する。但し、この情報の中のモードビット 202 の値は 0 である。さらに、送信フラグアドレスフィールド 204 は、送信完了フラグの書き込み領域 715 のアドレスを指定する。受信フラグアドレスフィールド 723 についても同じである。上記データ送信要求コマンドの場合には、付加情報記憶領域 714、722 は不要である。

40

**【0079】**

転送制御情報 700 内のモードビット 202 が 0 である場合、転送制御情報 700 内のモ

50

ードビット202が1である場合と比べると、メッセージ生成部407およびメッセージ分解部410の動作が異なる。すなわち、レジスタ506内のモードビットが0であるため、加算器508の出力は、レジスタ503内の転送データ長に等しく、加算器509の出力は、レジスタ502内の送信完了フラグアドレスに等しい。したがって、メッセージ生成部407では、比較回路510が、レジスタ504内の読み出し済みのデータの総量がアンドゲート507より与えられる、レジスタ503内の転送データ長に等しくなったことを検出したときに、メモリ読み出し要求発行部511は、ローカルメモリ701からの送信すべきデータの読み出しを終了する。この読み出されたデータを含むメッセージがすべて相互結合ネットワーク105へ送出されると、メッセージ生成部407内のメモ書き込み要求発行部512は、加算器509より与えられる、レジスタ502内の送信フラ 10  
グアドレスに示されるローカルメモリ領域に送信完了フラグを書き込むことを要求する書き込み要求をコマンド送信部411に伝達する。

#### 【0080】

受信側の要素プロセッサでは、相互結合ネットワーク105から転送されたメッセージ内のヘッダ内のモードビット717は0であるため、レジスタ606(図6)には0がセットされる。レジスタ606内のモードビットが0であるため、加算器608の出力はレジスタ603内の転送データ長に等しく、加算器609の出力は、レジスタ602にセットされる受信フラグアドレスに等しい。したがって、比較回路610が、レジスタ604内の受信済みのデータの総量が加算器608から出力される、転送データ長に等しくなったことを検出するまで、メモリ書き込み要求発行部611は、受信データのローカルメモリ 20  
への書き込みを要求する書き込み要求を発行する。すべての受信データがローカルメモリに書き込まれた後、メモリ書き込み要求発行部611は加算器609が出力する、レジスタ602内の受信フラグアドレスに従ってローカルメモリ304に受信完了フラグを書き込むための書き込み要求を発行する。こうして、MPIライブラリが要求したデータ転送が終了する。

#### 【0081】

なお、比較のために、従来のPUT処理では、MPI\_\_sendに示すように、送信側のPUT/GETライブラリは、ユーザデータとそれに対する付加情報をそれぞれ転送するための転送制御情報A、B(800、801)を作成し、通信領域802内にユーザデータとそれに対する送信完了フラグA、付加情報とそれに対する送信完了フラグBを記憶す 30  
るように構成され、送信側のネットワークインタフェース回路はこれらの転送制御情報に基づいて、二つのメッセージA、B(804、805)を送信する。受信側のネットワークインタフェース回路は、通信領域803内にユーザデータおよびそれに対する受信完了フラグ、付加情報とそれに対する受信完了フラグを書き込むように構成される。なお、転送制御情報A、Bには、本実施の形態で言うモードビットが存在しない。

#### 【0082】

これに対して、本実施の形態では、ユーザデータおよびその付加情報という異なる2つのデータをローカルメモリ間直接転送にしたがって1回のメッセージ転送で行なえる。したがって、2回のメッセージ転送が必要であった従来よりも転送レイテンシおよびローカルメモリアクセス回数を削減して並列処理効率を向上できる。さらに、MPIライブラリを 40  
介さない従来メモリ間直接転送も実行できる。

#### 【0083】

さらに、送信側のネットワークインタフェース回路は、送信完了フラグアドレスを付加情報と送信完了を示す制御情報の書き込みの両方に使用しているので、PUT/GETライブラリがネットワークインタフェース回路に対して指定すべき情報量が少なくて済む。また、受信側においても、ネットワークインタフェース回路は、受信完了アドレスを、付加情報と受信完了を示す制御情報の書き込みの両方に使用しているので、PUTメッセージに含まれる情報量が少なくて済み、それだけネットワークの混雑を防ぐことができ、さらに、メッセージの、送信元の要素プロセッサでの送信時間、ネットワーク上の転送時間および送信先の要素プロセッサにおける受信時間が短くでき、全体としてメッセージの転送 50

時間が短くなる。

【 0 0 8 4 】

以上のように、従来の方式では、異なる2種類のデータ送信処理を行う場合、ネットワークインタフェース回路305がローカルメモリを6回アクセス(転送制御情報の読み出し×2と転送データの読み出し×2とフラグの書き込み×2)し、受信処理の場合も送信と同様に、ネットワークインタフェース回路305がローカルメモリを4回アクセス(転送データの書き込み×2とフラグの書き込み×2)することになり、本実施の形態に比べて処理オーバーヘッドが大きい。本実施の形態では、転送制御情報にモードビットを持つことで、転送制御情報の読み出し回数、およびフラグのローカルメモリへの書き込み回数を送信側4回、受信側3回に削減してローカルメモリアクセスに関わる処理オーバーヘッドを小さくできる。また、従来は2回にわけて転送していたメッセージを1回で転送できるため、データ転送のレイテンシを削減できる。

10

【 0 0 8 5 】

以上ではMPIライブラリを介したPUT処理について説明したが、他の通信として、送信元の要素プロセッサが宛先の要素プロセッサのローカルメモリに格納されているユーザデータを取ってくるGET処理がある。本実施の形態は、MPIライブラリを介したGET処理にも同様に適用できる。すなわち、MPIライブラリとPUT/GET処理を併用して、ユーザデータとそれに対する付加情報を同時に一つのメッセージでGET(転送)できる。

【 0 0 8 6 】

まず、要求元のユーザプロセスは、GET要求コマンドをMPIライブラリに対して発行する。このコマンドは、要求するデータのアドレスとデータ長およびMPIライブラリが定めた付加情報等からなる点で、先に説明した送信要求コマンドと同様の引数を指定する。このコマンドを受けて、MPIライブラリとPUT/GETライブラリとネットワークインタフェース回路は、PUT動作の場合と同様にしてGET要求メッセージを宛先要素プロセッサに対して転送する。このメッセージには、PUTのときのメッセージにおける受信データアドレス、受信完了フラグアドレスの代わりに、GETすべきデータに対する送信データアドレスおよび送信完了フラグが含まれ、ユーザデータと付加情報は含まれない。このメッセージが従来のGET要求メッセージと異なる点は、モードビットを有することである。

20

【 0 0 8 7 】

宛先プロセッサは、GET要求メッセージに含まれたデータアドレスで指定されるローカルメモリ領域からユーザデータを読み出し、送信完了フラグアドレスで指定されるローカルメモリ領域から付加情報を読み出し、送信元プロセッサにそのユーザデータと付加情報を含むGETメッセージを送り返す。このメッセージは、受信データアドレスと受信完了フラグも含む。宛先プロセッサではそのメッセージを全て相互結合ネットワークに送出した後、送信完了フラグをローカルメモリ内の、上記付加情報の記憶領域の次のアドレス位置に書き込む。

30

【 0 0 8 8 】

送信元要素プロセッサは、GETメッセージを受信し、メッセージ内の受信データアドレスに従ってローカルメモリにそのメッセージ内のデータを書き込み、そのメッセージ内の受信完了フラグアドレスに従って、メッセージ内の付加情報をローカルメモリに書き込む。これらのユーザデータと付加情報をすべて受信した後に受信完了フラグをローカルメモリに書き込む。

40

【 0 0 8 9 】

PUT処理と同様に、GET処理でも、転送制御情報の読み出し回数およびフラグのローカルメモリへの書き込み回数を削減してローカルメモリアクセスに関わる処理オーバーヘッドを小さくできる。また、従来はユーザデータとそれに対する付加情報とを2回にわけて転送していたが、本実施の形態では、一つのGETメッセージでこれらの二つのデータを転送できるため、データ転送のレイテンシを削減できる。

50

## 【 0 0 9 0 】

## &lt; 変形例 &gt;

本発明は以上の実施の形態に限定されるのではなく、以下に示す変形例を含むいろいろの実施の形態により実施可能である。たとえば、ローカルメモリ304に格納する転送制御情報に拡張サイズフィールドを設ける。拡張サイズフィールドには、拡張したいフラグ領域のサイズをセットする。モードビットがセットされた場合、送信処理において、転送制御情報を含むヘッダと、送信データアドレスに従ってローカルメモリ領域からよみだされる転送データ長分のデータと、送信フラグアドレスに従ってローカルメモリ領域から読み出される拡張サイズフィールドに設定されたサイズ分の別のデータからメッセージを生成し、相互結合ネットワーク105に送出する。受信処理においては、メッセージをヘッダ部とデータ部に分解し、ヘッダ内の転送制御情報に含まれる受信データアドレスで示されるローカルメモリ領域にデータを転送データ長分書き込み、さらに、ヘッダ内の転送制御情報に含まれる受信フラグアドレスで示されるローカルメモリ領域に対し、別のデータを拡張サイズフィールドに設定されたサイズ分書き込む。モードビットがセットされていない場合、拡張サイズフィールドにセットされた値は無視される。

10

## 【 0 0 9 1 】

## 【 発明の効果 】

本実施の形態によれば、MPIのようなメッセージパッシングライブラリを介してメモリ間直接転送を行う場合におけるデータ転送処理をより高速に行うことができる。

## 【 図面の簡単な説明 】

20

【 図 1 】 本実施の形態が対象とする並列計算機の概略構成を示す図である。

【 図 2 】 本実施の形態で使用する転送制御情報の例を示す図である。

【 図 3 】 図 1 の装置に使用する要素プロセッサの概略構成を示す図である。

【 図 4 】 図 3 の要素プロセッサに使用するストレージコントローラおよびネットワークインタフェース回路の構成を示す図である。

【 図 5 】 図 4 のネットワークインタフェース回路内のメッセージ生成部の内部構成を示す図である。

【 図 6 】 図 4 のネットワークインタフェース回路内のメッセージ分解部の内部構成を示す図である。

【 図 7 】 本実施の形態におけるデータ転送処理の概要を説明する図である。

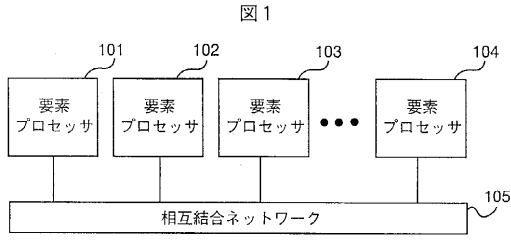
30

【 図 8 】 従来例のデータ転送処理の概要を説明する図である。

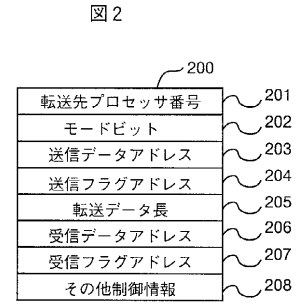
## 【 符号の説明 】

105 ... 相互結合ネットワーク

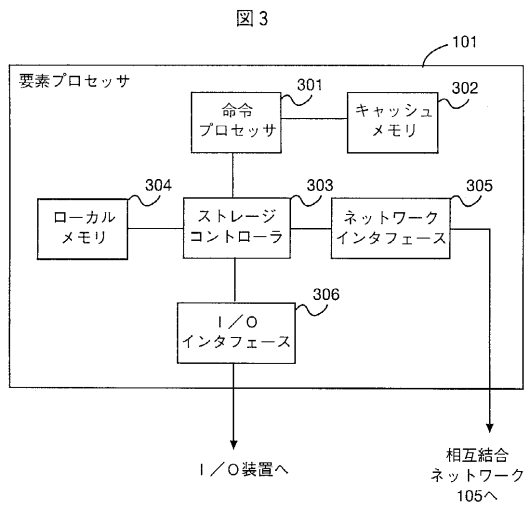
【 図 1 】



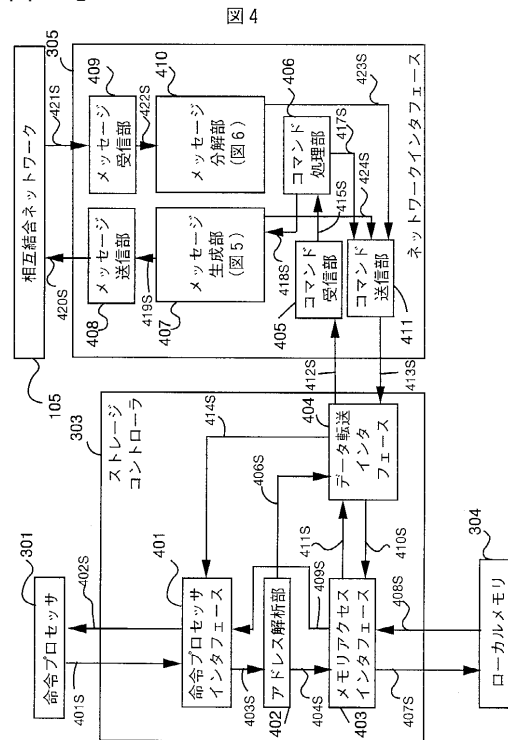
【 図 2 】



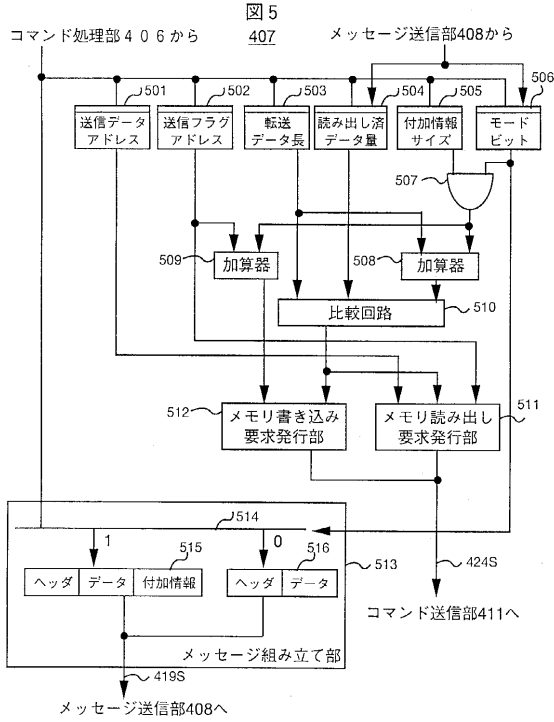
【 図 3 】



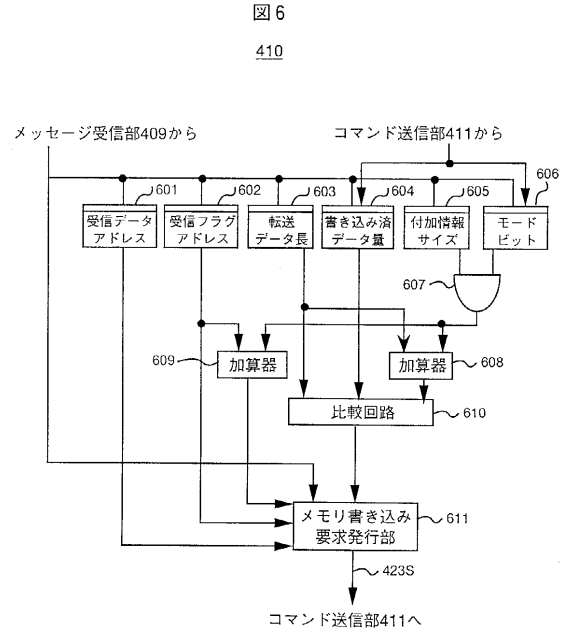
【 図 4 】



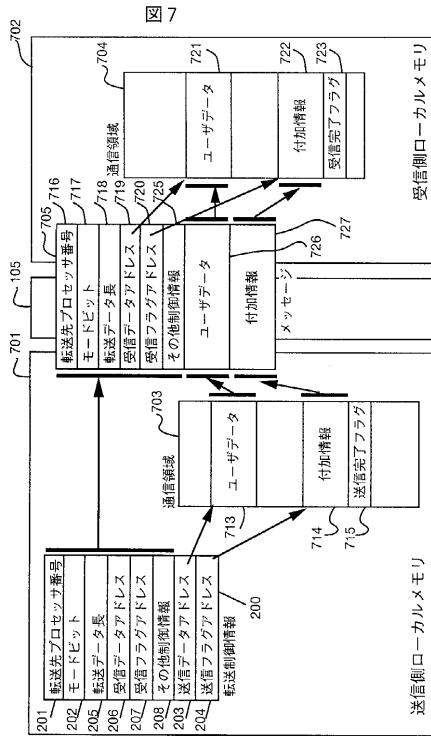
【 図 5 】



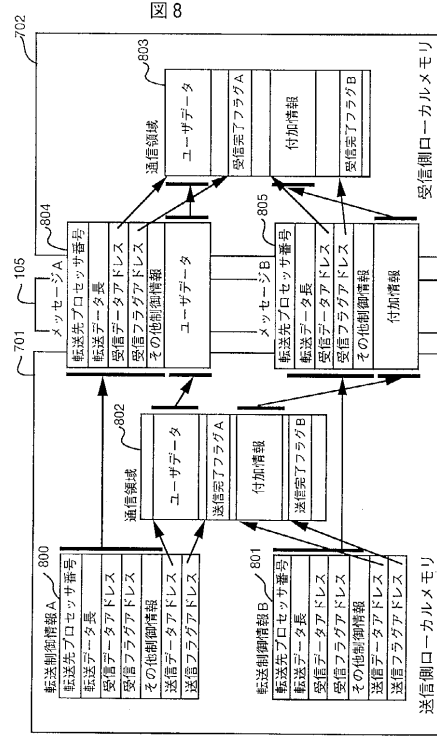
【 図 6 】



【 図 7 】



【 図 8 】





---

フロントページの続き

審査官 鳥居 稔

(56)参考文献 特開平7 - 311750 (JP, A)

特開平8 - 030566 (JP, A)

白木長武,小柳洋一,今村信貴,林憲一,清水俊幸,堀江健志,石畑宏明,高並列計算機AP1000  
+のメッセージハンドリング機構,情報処理学会論文誌 第37巻 第7号,日本,社団法人情  
報処理学会,1996年 7月15日,第37巻第7号,1388-1398

(58)調査した分野(Int.Cl.<sup>7</sup>, DB名)

G06F 15/16-173

G06F 13/00

G06F 9/46