

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号  
特許第7645495号  
(P7645495)

(45)発行日 令和7年3月14日(2025.3.14)

(24)登録日 令和7年3月6日(2025.3.6)

(51)国際特許分類 F I  
G 0 6 F 15/173(2006.01) G 0 6 F 15/173 6 6 5 D

請求項の数 3 (全15頁)

(21)出願番号	特願2022-530631(P2022-530631)	(73)特許権者	520211960 株式会社情報科学研究所 新潟県新潟市西区真砂一丁目15番36号
(86)(22)出願日	令和3年6月11日(2021.6.11)	(73)特許権者	311009158 株式会社エヌアイエスプラス 東京都文京区小石川2-4-10
(86)国際出願番号	PCT/JP2021/022268	(74)代理人	100121441 弁理士 西村 竜平
(87)国際公開番号	WO2021/251479	(74)代理人	100154704 弁理士 齊藤 真大
(87)国際公開日	令和3年12月16日(2021.12.16)	(74)代理人	100206151 弁理士 中村 惇志
審査請求日	令和5年12月27日(2023.12.27)	(74)代理人	100218187 弁理士 前田 治子
(31)優先権主張番号	特願2020-102484(P2020-102484)		
(32)優先日	令和2年6月12日(2020.6.12)		
(33)優先権主張国・地域又は機関	日本国(JP)		

最終頁に続く

(54)【発明の名称】 並列分散計算システム

(57)【特許請求の範囲】

【請求項1】

アドレス変換バッファ(TLB)を有するプロセッサ、物理メモリ、及び、当該物理メモリに直接アクセス可能なネットワークインターフェースコントローラ(NIC)を備えた複数の計算機を、複数のデータリンクを介して相互接続した並列分散計算システムであって、

送信元計算機のプロセス(以下、送信側プロセス)が、送信先計算機のプロセスを規定する操作対象プロセス(以下、受信側プロセス)の識別子と、当該受信側プロセスのメモリ領域を規定する操作対象アドレスと、書き込まれるデータサイズと、データ列とを含む操作要求パケットを送信し、

前記送信先計算機が、前記操作要求パケットを受信し、前記受信側プロセス及び前記操作対象アドレスにより規定されるメモリ領域に前記データ列を格納するものであり、

前記送信側プロセスは、前記受信側プロセス毎に前記操作要求パケットの送信数をカウントし、前記操作要求パケットに当該操作要求パケットのカウントアップ前又は後のカウント数を添付する送信数カウント部を有し、

前記受信側プロセスは、前記送信側プロセス毎に前記操作要求パケットの受信数をカウントする受信数カウント部を有し、

前記送信先計算機は、

前記受信数カウント部のカウント数と、前記操作要求パケットに添付されたカウント数とが連続している場合に、前記操作要求パケットの処理を行い、前記操作要求パケットに添付されたカウント数によって前記受信数カウント部のカウント数を更新し、

10

20

前記受信側プロセスへの到着順序が乱れて未着の操作要求パケットが存在する場合には、当該未着の操作要求パケットの前までしか前記受信数カウント部によるカウント数の更新を行わずに、前記未着の操作要求パケット以降の操作要求パケットの到着に関しては、操作要求パケットの処理を行い、前記受信数カウント部とは別に操作要求パケットに添付されたカウント数を既着カウント記録部に記録し、

未着だった操作要求パケットの到着によって前記既着カウント記録部を用いて前記受信側プロセスへの到着記録の記録内容が前記受信数カウント部から連続した場合に、その連続している範囲内で前記受信数カウント部を更新し、前記既着カウント記録部から更新済みの記録を削除し、

前記送信側プロセスは、自分が送ろうとする操作要求パケットがそれに先行する操作要求パケットの処理が済むまで操作されることを禁止したい場合は、確実に操作が終了して欲しい操作要求パケットの送信数カウントを複数リンク同期情報として前記操作要求パケットに付加する、並列分散計算システム。

10

#### 【請求項 2】

前記送信先計算機は、前記複数リンク同期情報の付加された操作要求パケットを受信した場合に、前記受信側プロセスの前記受信数カウント部を参照して、前記複数リンク同期情報に記述されたカウント数が前記受信数カウント部のカウント数以下であれば、操作要求パケットの処理を実行し、前記受信数カウント部及び前記既着カウント記録部の更新作業を行う、請求項 1 記載の並列分散計算システム。

#### 【請求項 3】

20

前記送信先計算機は、前記複数リンク同期情報の付加された操作要求パケットを受信した場合に、前記受信側プロセスの前記受信数カウント部を参照して、前記複数リンク同期情報に記述されたカウント数が前記受信数カウント部のカウント数より大きければ、前記受信数カウント部のカウント数が増えて、前記複数リンク同期情報に記述されたカウント数以上になるまで、操作要求パケットの処理を遅延し、前記受信数カウント部及び前記既着カウント記録部の更新作業も遅延させる、請求項 2 記載の並列分散計算システム。

#### 【発明の詳細な説明】

#### 【技術分野】

#### 【0001】

本発明は、アドレス変換バッファ (TLB) を有するプロセッサ、物理メモリ、及び、当該物理メモリに直接アクセス可能なネットワークインターフェースコントローラ (NIC) を備えた複数の計算機を、複数のデータリンクを介して相互接続した並列分散計算システムに関するものである。

30

#### 【背景技術】

#### 【0002】

本願発明者は、非特許文献 1 に示すように、メモリベース通信同期方式に基づく通信同期機構のうちのメモリベース通信ファシリティ (MBCF; Memory Based Communication Facility) の開発を進めている。この MBCF は、特殊な通信同期ハードウェアを一切必要とせず、一般のネットワークインタフェースカード (NIC) を使用して、ソフトウェアのみで遠隔メモリ操作による高速高機能通信同期を実現する機構である。

40

#### 【先行技術文献】

#### 【非特許文献】

#### 【0003】

【文献】 Matsumoto, T.: A Study on Memory-Based Communications and Synchronization in Distributed-Memory Systems. Dissertation Thesis, Graduate School of Science, Univ. of Tokyo (February 2001).

#### 【発明の概要】

#### 【発明が解決しようとする課題】

#### 【0004】

具体的に上記の MBCF は、アドレス変換バッファ (TLB) を有するプロセッサ、物理メ

50

メモリ、及び、当該物理メモリに直接アクセス可能なネットワークインターフェースコントローラ（NIC）を備えた計算機を用いて構成される。例えば、送信元計算機のプロセス（以下、送信側プロセス）が、送信先計算機のプロセスを規定する操作対象プロセス（以下、受信側プロセス）の識別子と、当該受信側プロセスのメモリ領域を規定する操作対象アドレスと、書き込まれるデータサイズと、データ列を含む操作要求パケットを送信する。そして、送信先計算機が、送信側プロセスが送信した操作要求パケットを受信し、受信側プロセス及び操作対象アドレスにより規定されるメモリ領域にデータ列を格納する。

【0005】

ところで、本願発明者は、複数の計算機を、複数のデータリンクを介して相互接続した（リンクアグリゲーション結合された）並列分散計算システムを構成することを考えている。このとき、各ノード（計算機）には複数のデータリンクが存在することが想定されるが、複数のデータリンクが存在するがゆえに、複数のリンクにパケットを負荷分散することにより、遠隔メモリ操作要求パケットの到着順序が乱れてメモリコンシステンシが破れてしまうという問題が考えられる。

10

【0006】

そこで本発明は、上記問題点を解決すべくなされたものであり、アドレス変換バッファ（TLB）を有するプロセッサ、物理メモリ、及び、当該物理メモリに直接アクセス可能なネットワークインターフェースコントローラ（NIC）を備えた複数の計算機を、複数のデータリンクを介して相互接続した並列分散計算システムにおいて、パケットの到着順序の乱れが存在する状況においてメモリコンシステンシを保持することを主たる課題とするものである。

20

【課題を解決するための手段】

【0007】

すなわち本発明に係る並列分散計算システムは、アドレス変換バッファ（TLB）を有するプロセッサ、物理メモリ、及び、当該物理メモリに直接アクセス可能なネットワークインターフェースコントローラ（NIC）を備えた複数の計算機を、複数のデータリンクを介して相互接続した並列分散計算システムであって、送信元計算機のプロセス（以下、送信側プロセス）が、送信先計算機のプロセスを規定する操作対象プロセス（以下、受信側プロセス）の識別子と、当該受信側プロセスのメモリ領域を規定する操作対象アドレスと、書き込まれるデータサイズと、データ列を含む操作要求パケットを送信し、前記送信先計算機が、前記操作要求パケットを受信し、前記受信側プロセス及び前記操作対象アドレスにより規定されるメモリ領域に前記データ列を格納するものであり、前記送信側プロセスは、前記受信側プロセス毎に前記操作要求パケットの送信数をカウントし、前記操作要求パケットに当該操作要求パケットのカウントアップ前又は後のカウント数を添付する送信数カウント部を有し、前記受信側プロセスは、前記送信側プロセス毎に前記操作要求パケットの受信数をカウントする受信数カウント部を有し、前記送信先計算機は、前記受信数カウント部のカウント数と、前記操作要求パケットに添付されたカウント数とが連続している場合に、前記操作要求パケットの処理を行い、前記操作要求パケットに添付されたカウント数によって前記受信数カウント部のカウント数を更新することを特徴とする。

30

【0008】

このような並列分散計算システムであれば、複数の計算機を複数のデータリンクを介して相互接続し、複数の計算機が相互にMBCFで通信同期を行う並列分散計算システムにおいて、送信側プロセスが、操作要求パケットに当該操作要求パケットのカウントアップ前又は後のカウント数を添付する送信数カウント部を有し、受信側プロセスが、送信側プロセス毎に操作要求パケットの受信数をカウントする受信数カウント部を有しており、送信先計算機が、受信数カウント部のカウント数と、操作要求パケットに添付されたカウント数とが連続している場合に、操作要求パケットの処理を行い、操作要求パケットに添付されたカウント数によって前記受信数カウント部のカウント数を更新するので、操作要求パケットの到着順序の乱れを必要に応じて解消することができる。

40

【0009】

50

未着の操作要求パケットが存在する場合の具体的な処理態様としては、前記送信先計算機は、前記受信側プロセスへの到着順序が乱れて未着の操作要求パケットが存在する場合には、当該未着の操作要求パケットの前までしか前記受信数カウント部によるカウント数の更新を行わずに、前記未着の操作要求パケット以降の操作要求パケットの到着に関しては、操作要求パケットの処理を行い、前記受信数カウント部とは別に操作要求パケットに添付されたカウント数を既着カウント記録部に記録する。

【0010】

未着の操作要求パケットが到着した後の具体的な処理態様としては、前記送信先計算機は、未着だった操作要求パケットの到着によって前記既着カウント記録部を用いて前記受信側プロセスへの到着記録の記録内容が前記受信数カウント部から連続した場合に、その連続している範囲内で前記受信数カウント部を更新し、前記既着カウント記録部から更新済みの記録を削除する。

10

【0011】

操作要求パケットの処理順序を保障するためには、前記送信側プロセスは、自分が送ろうとする操作要求パケットがそれに先行する操作要求パケットの処理が済むまで操作されることを禁止する意図を伝えることとし、確実に操作が終了して欲しい先行する操作要求パケットの送信数カウントを複数リンク同期情報として前記操作要求パケットに付加する。

【0012】

操作要求パケットの処理順序を保障するための具体的な操作要求パケットの処理態様としては、前記送信先計算機は、前記複数リンク同期情報の付加された操作要求パケットを受信した場合に、前記受信側プロセスの前記受信数カウント部を参照して、前記複数リンク同期情報に記述されたカウント数が前記受信数カウント部のカウント数以下であれば、操作要求パケットの処理を実行し、前記受信数カウント部及び前記既着カウント記録部の更新作業を行う。

20

【0013】

操作要求パケットの処理順序を保障するための具体的な操作要求パケットの処理態様としては、前記送信先計算機は、前記複数リンク同期情報の付加された操作要求パケットを受信した場合に、前記受信側プロセスの前記受信数カウント部を参照して、前記複数リンク同期情報に記述されたカウント数が前記受信数カウント部のカウント数より大きければ、前記受信数カウント部のカウント数が増えて、前記複数リンク同期情報に記述されたカウント数以上になるまで、操作要求パケットの処理を遅延し、前記受信数カウント部及び前記既着カウント記録部の更新作業も遅延させる。

30

【発明の効果】

【0014】

このように構成した本発明によれば、複数の計算機を複数のデータリンクを介して相互接続し、複数の計算機が相互にMBCFで通信同期を行う並列分散計算システムにおいて、パケットの到着順序の乱れが存在する状況においてメモリコンシステンシを保持することが低コストで達成できる。

【図面の簡単な説明】

40

【0015】

【図1】本発明の一実施形態における並列分散計算システムの全体構成を模式的に示す図である。

【図2】同実施形態の送信元計算機と送信先計算機の物理構成を示す模式図である。

【図3】同実施形態のMBCF\_WRITEの使用環境を示す模式図である。

【図4】同実施形態の操作要求側におけるパケット生成過程を含む、操作要求パケット送信の動作手順を示す模式図である。

【図5】同実施形態の受信側ノードに操作要求パケットが到着した状況を示す模式図である。

【図6】同実施形態の受信側ノードにおけるパケット受信割込みルーチン内のMBCF関連

50

処理を示す模式図である。

【図 7】操作要求パケットの到着順序が乱れた状態を示す模式図である。

【図 8】同実施形態の計算機の機能ブロック図である。

【図 9】ノード Na、Nb、Nc において、ノード Na の送信数カウント部及び受信数カウント部を代表させて示す模式図である。

【図 10】同実施形態の同期情報を挿入した遠隔メモリ操作の使用例 1 である。

【図 11】同実施形態の同期情報を挿入した遠隔メモリ操作の使用例 2 である。

【符号の説明】

【0016】

100・・・並列分散計算システム

10

2・・・計算機

2X・・・送信元計算機

2Y・・・送信先計算機

21・・・プロセッサ

22・・・物理メモリ

23・・・ネットワークインターフェースコントローラ (NIC)

3・・・データリンク

201・・・送信数カウント部

202・・・受信数カウント部

203・・・既着カウント記録部

20

【発明を実施するための形態】

【0017】

以下、本発明の一実施形態に係る並列分散計算システム 100 について、図面を参照して説明する。

【0018】

本実施形態の並列分散計算システム 100 は、図 1 に示すように、複数の計算機 2 を複数のデータリンク 3 を介して相互接続したものである。複数の計算機 2 を複数のデータリンク 3 を介して相互接続することによりリンクアグリゲーション結合された並列分散計算システムが構成される。

【0019】

30

各計算機 2 は、図 2 に示すように、アドレス変換バッファ (TLB; Translation Look-aside Buffer) を有するプロセッサ 21、物理メモリ 22、及び、当該物理メモリ 22 に直接アクセス可能なネットワークインターフェースコントローラ (NIC; Network Interface Card) 23 を備えたものである。

【0020】

そして、並列分散計算システム 100 は、特殊な通信同期ハードウェアを一切必要とせず、一般のネットワークインタフェースカード (NIC) 23 を使用して、ソフトウェアのみで遠隔メモリ操作による高速高機能通信同期を実現するメモリベース通信ファシリティ (MBCF; Memory Based Communication Facility) を構築するものである。具体的に並列分散計算システム 100 は、各計算機 2 のカーネル空間内に格納されたオペレーティングシステム (OS) によって、メモリベース通信ファシリティ (MBCF) を構築する。

40

【0021】

並列分散計算システム 100 は、例えば、以下に示す遠隔メモリ書き込みを行う write コマンド (MBCF\_WRITE) や遠隔メモリ読み出しを行う read コマンド (MBCF\_READ) などの各種の操作コマンドのバリエーションを有する。

【0022】

例えば、送信元計算機 2 (2X) のプロセス (以下、送信側プロセス) が、送信先計算機 2 (2Y) のプロセスを規定する操作対象プロセス (以下、受信側プロセス) の識別子と、当該受信側プロセスのメモリ領域を規定する操作対象アドレスと、書き込まれるデータサイズと、データ列とを含む操作要求パケットを送信し、送信先計算機 2 が、操作要求

50

パケットを受信し、受信側プロセス及び操作対象アドレスにより規定されるメモリ領域にデータ列を格納する (MBCF\_WRITE)。

【 0 0 2 3 】

また、送信側プロセスが、送信先計算機 2 ( 2 Y ) のプロセスを規定する操作対象プロセス (以下、受信側プロセス) の識別子と、当該受信側プロセスのメモリ領域を規定する操作対象アドレスと、読み出すデータサイズと、送信側プロセスのデータ格納領域アドレスとを含む操作要求パケットを送信し、送信先計算機 2 ( 2 Y ) が、操作要求パケットを受信し、受信側プロセス及び操作対象アドレスにより規定されるメモリ領域からデータ列を読み出して、送信側プロセスのデータ格納領域に返送する (MBCF\_READ)。

【 0 0 2 4 】

ここで、MBCF\_WRITEの手順について、図 3 ~ 図 6 を参照して詳述する。なお、その他のコマンドについても基本となる手順は共通している。

【 0 0 2 5 】

図 3 はMBCF\_WRITEの使用環境を示しており、Pnode1が操作要求パケットの送信元ノード (送信元計算機) であり、Pnode 2 が受信側ノード (送信先計算機) である。送信元ノードのプロセッサは、自分のメモリのNIC DMA領域内に配送情報の書かれたヘッダとペイロードを含むパケットイメージを作成する。NIC DMA領域メモリには送信又は受信のためにNICが直接アクセスすることが可能である。送信用のパケットイメージの生成が終わると、プロセッサはNICにDMA読み出しによる送信動作を開始させるための指示 (kick動作) を行う。受信側ノード (Pnode2) は、自分のメモリのNIC DMA領域内に到達するパケットのためのリングバッファを持っている。受信側ノードのNICは、パケットが自分のノード宛てのパケットであるか判断して (通常はMACアドレスで判断する)、自分宛てのパケットに関しては、リングバッファにそのコピーを生成する。その後、NICは、受信側ノードのプロセッサにパケットの到着を知らせるためにハードウェア割り込みを発生させる。

【 0 0 2 6 】

次に、図 4 を参照して、操作要求側におけるパケット生成過程を含む、操作要求パケット送信の動作手順を示す。

【 0 0 2 7 】

送信側プロセス (要求側タスク) では、受信側プロセス (要求先タスク) の識別子[Ltask1]、受信側プロセスの操作対象メモリアドレス[Laddr1]、受信側プロセスのメモリ空間操作のためのアクセスキー[AccessKey]、MBCFのコマンド種別[MBCF\_WRITE]、遠隔書き込みを行うデータサイズ[n]、書き込むデータが格納された領域の先頭へのポイント[Laddr0]から成るパラメータを用意する。そして、これらのパラメータを伴ってMBCF要求送信用システムコールを呼び出す。システムコールを受けてOSは、送信側プロセスのタスク表を参照して、受信側プロセスを示す論理タスクIDを物理タスクID[(Pnode2, Ptask5)]に変換する。物理タスクIDは物理的なノードIDであるPnode2を含んでいるため、この情報から受信側ノードへの経路情報 (配送先情報) が設定できる。使用するネットワークがEthernetであれば、MACアドレスが配送先情報として使用される。この配送先情報によってNICは操作要求パケットを受信側ノードまで配送することが可能になる。そして、OSはNICに操作要求パケットを送信させる。

【 0 0 2 8 】

次に、図 5 及び図 6 を参照して、操作要求先における操作要求パケットの受信手順を示す。

図 5 は、受信側ノードに操作要求パケットが到着した状況を示す。操作要求パケットは受信側ノード[Pnode2]までネットワークによって運ばれてくる。パケット到着時に、受信側ノードのNICはDMAによって操作要求パケットのデータイメージをリングバッファにコピーし、その後受信側ノードのプロセッサにパケットが到着したことを知らせるために、割り込み信号を発生させる。

【 0 0 2 9 】

10

20

30

40

50

図 6 に受信側ノードにおけるパケット受信割込みルーチン内のMBCF関連処理を示す。NICからの受信割込み発生により、受信側ノードのプロセッサはパケット受信ルーチンに制御が切り替わり、NICによって要求される低レベルのパケット受信手順を最初に行う。パケット受信ルーチンでは、まず物理タスクID ( 図中のPtask5 ) から受信側プロセスを特定する。具体的には、そのプロセスのプロセス構造体へのポインタを得る。次に、操作要求パケット内のAccessKeyが受信側プロセスのものと同じかどうかチェックして、一致している場合のみ受信側プロセスのメモリ空間内のメモリ操作を許可する。一致した場合は、メモリ空間のコンテキストを受信側プロセスのものに切り替えて、操作要求パケットで運ばれてきたnバイトのデータを操作対象論理アドレス ( Laddr1 ) から、特権レベルとしてではなくユーザ権限のstore命令によって書き込む。そして、メモリ空間のコンテキストを割込み発生時のものに戻しておく。ここまでで割込みルーチン内の処理が終了する。

#### 【 0 0 3 0 】

そして、本実施形態の並列分散計算システム 1 0 0 は、複数のデータリンク 3 を有することから、図 7 に示すように、操作要求パケットの到着順序が乱れてしまうという問題が考えられる。このため、本実施形態の並列分散計算システム 1 0 0 は、受信側プロセスへの操作要求パケットの到達順序をユーザの意図する範囲内において保障する機能を有している。

#### 【 0 0 3 1 】

ユーザの意図する範囲内でのみの保障というのは、具体的には、受信側プロセスのメモリ上にA, B, C, D, Fという変数があり、送信側プロセスがそのA, B, C, DにそれぞれMBCF\_WRITEで遠隔書き込みを行い、最後に、それらの書き込みが済んだことを示す値 ( 例えば 1 ) をFにMBCF\_WRITEで書き込むことにすると、A, B, C, DのMBCF\_WRITEの間の順序保障は不要であるが、それら 4 つのMBCF\_WRITEとFに対するMBCF\_WRITEの間の順序保障は不可欠である。逆に、全部のMBCF\_WRITE間で順序保障を行うことは過剰の保障であり、余分なオーバーヘッドを発生させる可能性がある。よって、本実施形態では、Dの書き込みまでのMBCF操作要求パケットがすべて処理された後に、Fの書き込みを行うという同期情報を、FのMBCF\_WRITE要求パケットにのみ付加する。

#### 【 0 0 3 2 】

操作要求パケットの到達順序をユーザの意図する範囲内において保障する機能は、具体的には、図 8 に示すように、送信元計算機 2 ( 2 X ) の各送信側プロセスは、受信側プロセス毎に操作要求パケットの送信数をカウントし、操作要求パケットに当該操作要求パケットのカウントアップ前又は後のカウント数を添付する送信数カウント部 2 0 1 を有している。また、送信先計算機 2 ( 2 Y ) の各受信側プロセスは、送信側プロセス毎に操作要求パケットの受信数をカウントする受信数カウント部 2 0 2 を有している。

#### 【 0 0 3 3 】

そして、送信先計算機 2 ( 2 Y ) は、受信数カウント部 2 0 2 のカウント数と、操作要求パケットに添付されたカウント数とが連続している場合に、操作要求パケットの処理を行い、操作要求パケットに添付されたカウント数によって受信数カウント部 2 0 2 のカウント数を更新する。

#### 【 0 0 3 4 】

ここで、送信先計算機 2 ( 2 Y ) は、受信側プロセスへの到着順序が乱れて未着の操作要求パケットが存在する場合には、当該未着の操作要求パケットの前までしか受信数カウント部 2 0 2 によるカウント数の更新を行わない。また、未着の操作要求パケット以降の操作要求パケットの到着に関しては、以下に述べる複数リンク同期情報が付加されていないければ、操作要求パケットの処理を行い、受信数カウント部 2 0 2 とは別に操作要求パケットに添付されたカウント数を既着カウント記録部 2 0 3 に記録する。

#### 【 0 0 3 5 】

そして、送信先計算機 2 ( 2 Y ) は、未着だった操作要求パケットの到着によって既着カウント記録部 2 0 3 を用いて受信側プロセスへの到着記録の記録内容が受信数カウント

10

20

30

40

50

部 2 0 2 から連続した場合に、その連続している範囲内で受信数カウント部 2 0 2 を更新し、既着カウント記録部 2 0 3 から更新済みの記録を削除する。

【 0 0 3 6 】

このとき、送信側プロセスは、自分が送ろうとする操作要求パケットがそれに先行する操作要求パケットの処理が済むまで操作されることを禁止したい場合、つまり到着順序をユーザ（送信側プロセス）が保障したい場合は、確実に操作が終了して欲しい操作要求パケットの送信数カウントを複数リンク同期情報として操作要求パケットに付加する。

【 0 0 3 7 】

送信先計算機 2 ( 2 Y ) は、複数リンク同期情報の付加された操作要求パケットを受信した場合に、受信側プロセスの受信数カウント部 2 0 2 を参照して、複数リンク同期情報に記述されたカウント数が受信数カウント部 2 0 2 のカウント数以下であれば、操作要求パケットの処理を実行し、受信数カウント部 2 0 2 及び既着カウント記録部 2 0 3 の更新作業を行う。

10

【 0 0 3 8 】

また、送信先計算機 2 ( 2 Y ) は、複数リンク同期情報の付加された操作要求パケットを受信した場合に、受信側プロセスの受信数カウント部 2 0 2 を参照して、複数リンク同期情報に記述されたカウント数が受信数カウント部 2 0 2 のカウント数より大きければ、受信数カウント部 2 0 2 のカウント数が増えて、複数リンク同期情報に記述されたカウント数以上になるまで、操作要求パケットの処理を遅延し、受信数カウント部 2 0 2 及び既着カウント記録部 2 0 3 の更新作業も遅延させる。

20

【 0 0 3 9 】

次に、受信数カウント部 2 0 2 のカウント数の更新方法について、既着カウント記録部 2 0 3 にビットベクタを用いた一例を説明する。

ビットベクタを使ってパケット番号に抜けが無い状態でカウンタをカウントアップが可能になったら、カウンタとビットベクタを更新する。具体的な更新方法は以下である。

【 0 0 4 0 】

ビットベクタ長 8bit で 101 番のパケットが欠落して、102 番、103 番が到着している場合は、カウンタの値は 100 で、ビットベクタは二進数で 00000110 となる。つまり、最下位より (102-100)、(103-100) 番目の bit が 1 にセットされる。

【 0 0 4 1 】

ここに 105 番目のパケットが到着すると、カウンタの値は 100 のままで、ビットベクタは二進数で 00010110 となる。つまり、(105-100) 番目の bit が 1 にセットされる。

30

【 0 0 4 2 】

この状態で、101 番目のパケットが到着すると、(101-100) 番目の bit が 1 にセットされ、ビットベクタは二進数で 00010111 となる。一番右端の bit に 1 がセットされているので、このことはカウンタの値に連続したパケットの到着を意味し、一番右端の bit がゼロになるまでこのビットベクタを右シフトして、そのシフト数をカウンタに加える。

【 0 0 4 3 】

つまり、シフト数は 3 になるので、カウンタの値は 103 で、ビットベクタは二進数で 0000010 となる。ここに 104 番目のパケットが到着すると、(104-103) 番目の bit が 1 にセットされ、ビットベクタは二進数で 0000011 となり、先ほどと同じく、1 が続く 2 をカウンタに加えてカウンタの値は 105 で、ビットベクタは二進数で 0000000 となる。ビットベクタがオールゼロの状態は、パケットに抜けが無くパケットがカウンタの値まで到着している状態である。

40

【 0 0 4 4 】

また、以下の方法であっても良い。

今、カウンタの値は 105 で、ビットベクタは二進数で 00000000 で 106 番のパケットが到着した（順序の乱れがない）とすると、一旦、ビットベクタの (106-105) 番目の bit が 1 にセットされる。このとき、カウンタの値は 105 で、ビットベクタは二進数で 00000001 となり、ビットベクタの一番右端の bit に 1 がセットされているので、1bit 分右シフトして

50

、カウンタの値は+1で106となり、ビットベクタは二進数で00000000となる。

【0045】

この方法では、順序の乱れがないことを確認する手順を入れずに、抜けがある場合の手順と同一で処理することが可能となる。

【0046】

なお、上記では、8bitのビットベクタを用いているが、データリンクが8本までのリンクアグリゲーションを仮定すると、ビットベクタは32bitぐらいあれば、実用上問題は無い。万が一、ビットベクタのビット幅を超える順序の混乱が生じた場合は、ビットベクタのビット幅を超える番号を持つパケットは受け取りを拒否して再送対象にする。

【0047】

次に、図9に示すように、3つの計算機2において、それぞれをノードNa、Nb、Ncとし、

ノードNaに、プロセスTa1、Ta2、Ta3、

ノードNbに、プロセスTb1、Tb2、Tb3、

ノードNcに、プロセスTc1、Tc2、Tc3、

の各3つのプロセス(タスクもいう。)が存在して、相互にMBCFで通信同期を行っている場合を考える。なお、送信側プロセスは、MBCFを使って同一計算機の中の自らを含む他プロセスへのメモリ操作要求を遠隔メモリ操作の一環として行うことができる。

【0048】

この場合、プロセスTa1には、送信数カウンタ部201(送信数カウンタ)が、プロセスTa1、Ta2、Ta3、Tb1、Tb2、Tb3、Tc1、Tc2、Tc3に対応して、送信数カウンタXa1\_Ta1、Xa2\_Ta1、Xa3\_Ta1、Xb1\_Ta1、Xb2\_Ta1、Xb3\_Ta1、Xc1\_Ta1、Xc2\_Ta1、Xc3\_Ta1の9個が存在する。ここでは、自分自身や同じノードのプロセスに対しても順序制御すると仮定している。なお、これらではパケット到着順序の乱れが想定されないで順序制御を省略可能であるが、通信相手の受信側プロセスが同一計算機内に存在するかどうかをユーザやプログラマが意識しなければならないことは、プログラム作成の困難度を増す。このため、受信側プロセスがどの計算機に存在するかに関わらず、同一仕様の送信数カウンタと複数リンク同期情報を含むMBCF操作要求パケットを使用するようにシステムを構成する。

【0049】

そして、例えば、プロセスTa1がプロセスTb2に向かって操作要求パケットを発行した場合には、送信数カウンタXb2\_Ta1がカウントアップされる。同じように、他のプロセスにも9個ずつ送信数カウンタが存在する。

【0050】

また、プロセスTa1には、受信数カウンタ部202(受信数カウンタ)が、プロセスTa1、Ta2、Ta3、Tb1、Tb2、Tb3、Tc1、Tc2、Tc3に対応して、受信数カウンタRa1\_Ta1、Ra2\_Ta1、Ra3\_Ta1、Rb1\_Ta1、Rb2\_Ta1、Rb3\_Ta1、Rc1\_Ta1、Rc2\_Ta1、Rc3\_Ta1の9個が存在する。

【0051】

そして、例えば、プロセスTa1宛の操作要求パケットをノードNaがプロセスTc2から受信し、その操作要求パケット内のカウント値が(Rc2\_Ta1+1)であれば、メモリ操作後に受信数カウンタRc2\_Ta1をカウントアップする。操作要求パケット内のカウント値が(Rc2\_Ta1+1)より大きければ、メモリ操作後に受信数カウンタRc2\_Ta1に対応した既存カウント記録部203(ビットベクタを用いて構成)に記録して、追い越された操作要求パケットのメモリ操作終了後まで受信数カウンタRc2\_Ta1のカウントアップを遅延させる。

【0052】

プロセスTc2からのプロセスTa1への操作要求パケットに複数リンク同期情報(同期オプション)が付加されている場合は、受信数カウンタRc2\_Ta1の値が同期オプションで記述されたカウント値以上であった場合は、同期オプションなしの操作要求パケットと同じ

10

20

30

40

50

処理を行う。

【 0 0 5 3 】

一方で、受信数カウンタRc2\_Ta1の値が同期オプションに記述されたカウント値より小さい場合は、受信数カウンタRc2\_Ta1の値が同期オプションに記述されたカウント値以上になるまでメモリ操作や既着カウント記録部（ビットベクタ）への登録も含めて、実行を遅延させる。

【 0 0 5 4 】

< 2 . 本実施形態の効果 >

このように構成した本実施形態の並列分散計算システム 1 0 0 によれば、複数の計算機 2 を複数のデータリンク 3 を介して相互接続し、複数の計算機 2 が相互にMBCFで通信同期を行うものにおいて、送信側プロセスが、操作要求パケットに当該操作要求パケットのカウントアップ前又は後のカウント数を添付する送信数カウント部 2 0 1 を有し、受信側プロセスが、送信側プロセス毎に操作要求パケットの受信数をカウントする受信数カウント部 2 0 2 を有しており、送信先計算機 2 ( 2 Y ) が、受信数カウント部 2 0 2 のカウント数と、操作要求パケットに添付されたカウント数とが連続している場合に、操作要求パケットの処理を行い、操作要求パケットに添付されたカウント数によって受信数カウント部 2 0 2 のカウント数を更新するので、操作要求パケットの到着順序の乱れが発生する環境下において、ユーザの意図する範囲内でメモリ操作順序を保障する、つまり、メモリコンシステンシを維持することができる。

【 0 0 5 5 】

従来、分散配置されたメモリに対して共有メモリアクセスを許す分散共有メモリ型並列計算機において、メモリの操作順序を保障するために、すべての遠隔メモリ操作に対して確認応答 ( Acknowledgement, 以下ではAck ) が要求元に返送され、順序を保障しないとならない場合は、それまでの遠隔メモリ操作すべてに対するAckが戻って来たことを確認後に、順序を保障すべき遠隔メモリ操作の要求を行っていた。この順序保障の仕組みはメモリバリアと呼ばれている。それに対して、本発明においては、順序保障に対して遠隔メモリ操作のAckを必要とせず、直近に要求した遠隔メモリ操作に続けて、順序を保障すべき遠隔メモリ操作の要求を行うことが可能である。このことは先行する遠隔メモリ操作のAckが返ることを待つ必要がないことを意味し、大幅な性能向上を可能にする。また、送信側プロセスは自分が保持する送信数カウント部 2 0 1 の値を使って複数リンク同期情報を付加することができ、送信先計算機 2 ( 2 Y ) は同じ計算機内の受信側プロセスに属する受信数カウント部 2 0 2 の値を使って、同期成立・非成立の判定が可能である。つまり、同期を行うために必要な情報を、送信側プロセスも送信先計算機 2 ( 2 Y ) も極めて低コストで入手することが可能である。

【 0 0 5 6 】

本発明を使うことによって、Ackを待つ必要がないこと以上に同期によるオーバーヘッドコストを削減できる可能性がある。その可能性について、具体例によって説明する。受信側プロセスに、A1, B1, C1というメモリ領域があり、それらに値の格納を終えたことを示すF1というフラグ変数がある。同様に、同じ受信側プロセスにA2, B2, C2, A3, B3, C3というメモリ領域があり、A2, B2, C2に値の格納を終えたことを示すF2というフラグ変数、A3, B3, C3に値の格納を終えたことを示すF3というフラグ変数がある。すべてのメモリ領域とフラグへの値の格納をMBCF\_WRITEによって実現するものとする。A1, B1, C1とF1、A2, B2, C2とF2、A3, B3, C3とF3の遠隔書き込みの間にはそれぞれ守るべき順序関係があるが、AnとBnとCnの間やEm, En ( ただし、m < n ) の間には順序関係はない。これを単純に図 1 0 のように、MBCF\_WRITE要求パケットを発行すると、F1, F2, F3への書き込みは、それぞれ直前のC1, C2, C3への書き込みが行われた後であることを保障せねばならない。複数のリンクが存在する下では、要求を出した順番で届くことは保障されないため、送信先計算機 2 ( 2 Y ) において、操作要求パケット処理を同期のために遅延させなければならない可能性がかなりある。これに対して、図 1 1 のように、無関係なMBCF\_WRITEを間に挿む形で、順序を保障すべきメモリ領域への書き込みとフラグ変

数への書き込みの距離を離すことにより、付加するリンク同期情報のカウント値が、直前に発行したMBCF要求カウントの送信数カウント部201の値ではなく、3パケット前の送信数カウント部201の値で済むことになる。このことは送信先計算機2(2Y)への要求パケットの到着は直前の2つのパケットより早くても構わないことを意味し、図11のケースは図10のケースよりも明らかに、フラグ変数への書き込みが同期待ちによって、送信先計算機2(2Y)において遅延される可能性が低い。つまり、本発明を使用し、同期の必要な操作要求と不要な操作要求、ならびに順序関係がない操作要求の発行順序を最適化することにより、同期のオーバーヘッドコストをより削減できる可能性がある。

【0057】

その他、本発明は前記実施形態に限られず、その趣旨を逸脱しない範囲で種々の変形が可能であるのは言うまでもない。

10

【産業上の利用可能性】

【0058】

本発明は、複数の計算機を複数のデータリンクを介して相互接続し、複数の計算機が相互にMBCFで通信同期を行う並列分散計算システムにおいて、パケットの到着順序の乱れが存在する状況においてメモリコンシステンシを保持することを低コストで達成することができる。

20

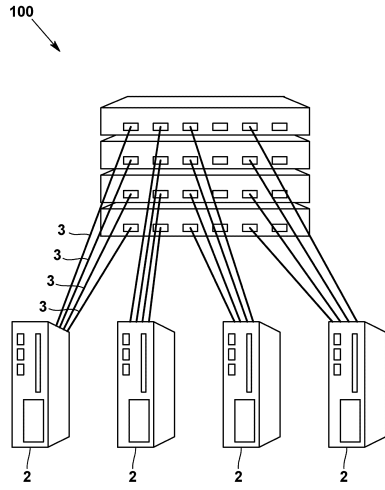
30

40

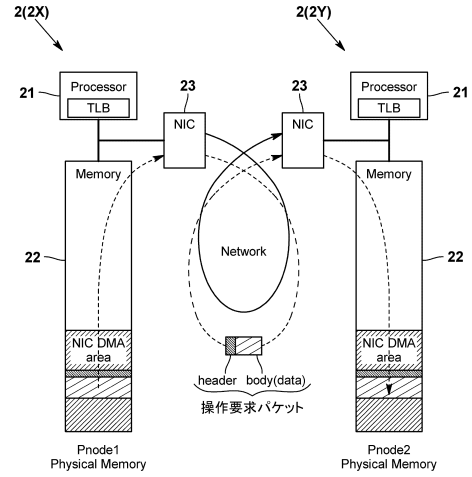
50

【 図面 】

【 図 1 】



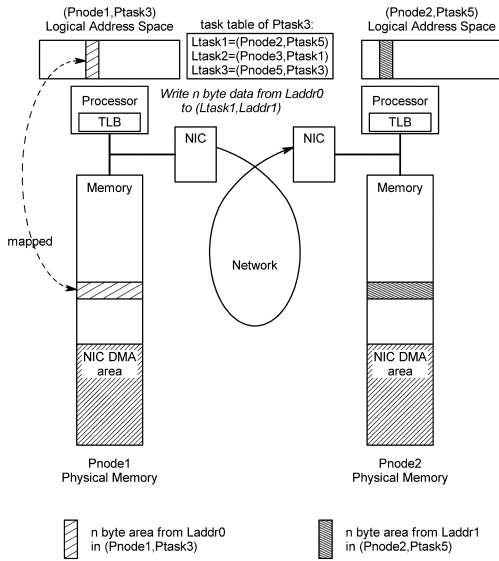
【 図 2 】



10

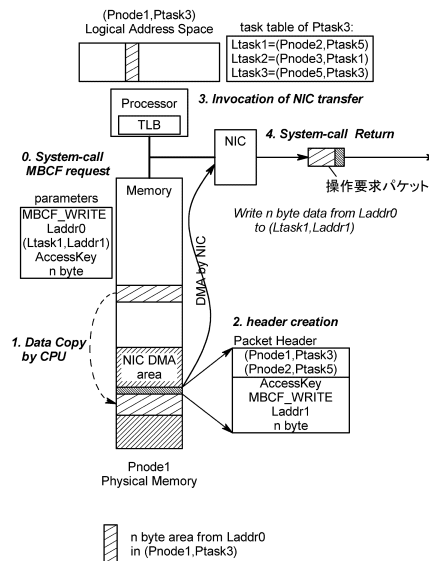
20

【 図 3 】



Ptask : Physical Process (=task) ID  
 Pnode : Physical node ID  
 Ltask : Logical Process (=task) ID  
 Laddr : Logical address

【 図 4 】



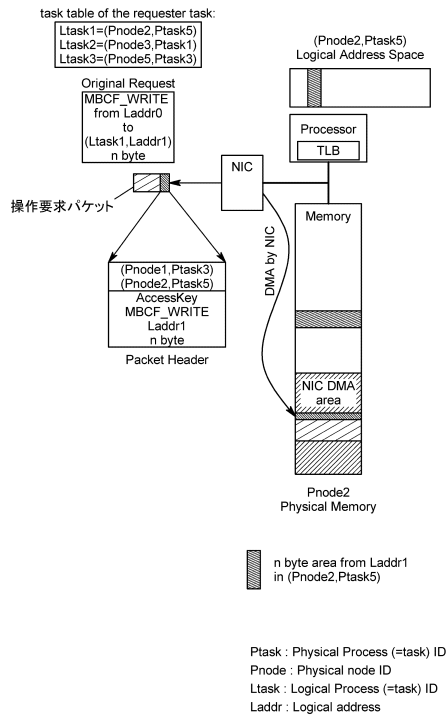
Ptask : Physical Process (=task) ID  
 Pnode : Physical node ID  
 Ltask : Logical Process (=task) ID  
 Laddr : Logical address

30

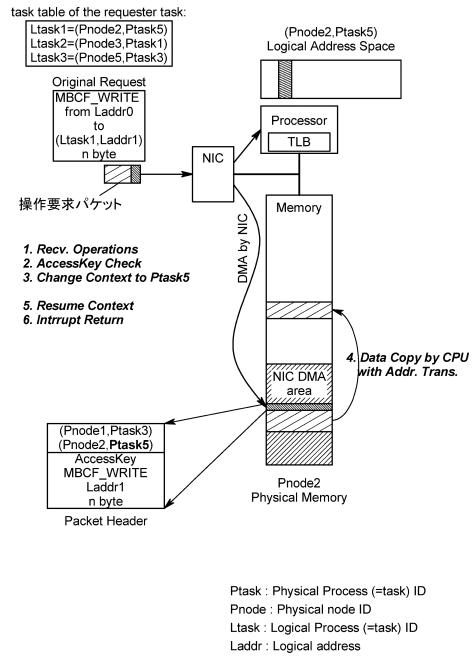
40

50

【 図 5 】



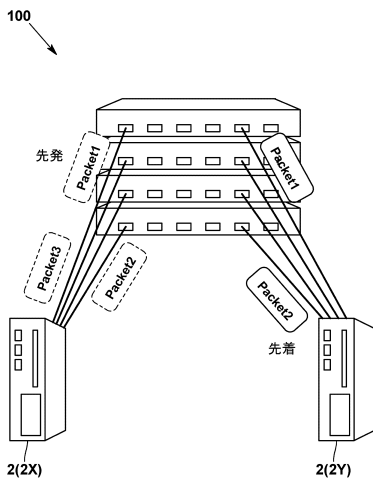
【 図 6 】



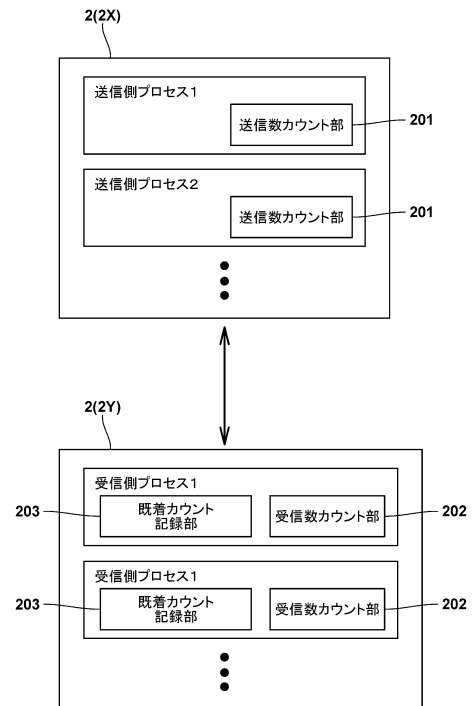
10

20

【 図 7 】



【 図 8 】

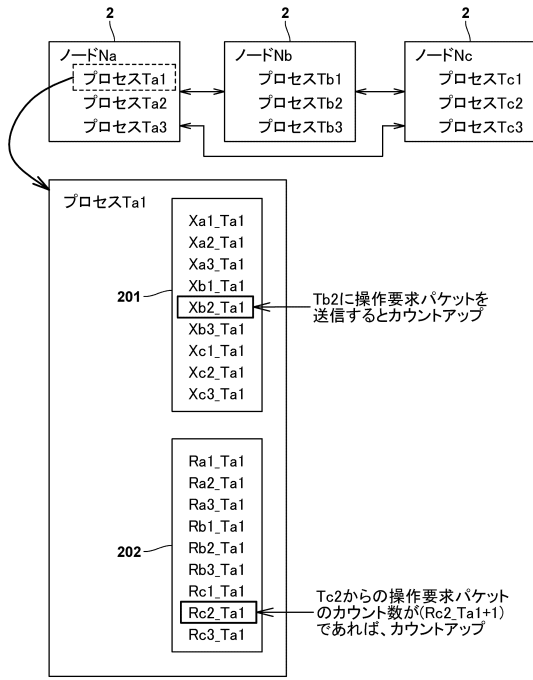


30

40

50

【図 9】



【図 10】

送信数カウント	種数リンク同期情報	コマンド	対象アドレス	サイズ(バイト数)	データ列
1	—	MBCF_WRITE	A1のアドレス	N	a <sub>1</sub> d <sub>1</sub> ...
2	—	MBCF_WRITE	B1のアドレス	N	b <sub>1</sub> d <sub>1</sub> ...
3	—	MBCF_WRITE	C1のアドレス	N	c <sub>1</sub> d <sub>1</sub> ...
4	3	MBCF_WRITE	F1のアドレス	4	1
5	—	MBCF_WRITE	A2のアドレス	N	a <sub>1</sub> d <sub>1</sub> ...
6	—	MBCF_WRITE	B2のアドレス	N	b <sub>1</sub> d <sub>1</sub> ...
7	—	MBCF_WRITE	C2のアドレス	N	c <sub>1</sub> d <sub>1</sub> ...
8	7	MBCF_WRITE	F2のアドレス	4	1
9	—	MBCF_WRITE	A3のアドレス	N	a <sub>1</sub> d <sub>1</sub> ...
10	—	MBCF_WRITE	B3のアドレス	N	b <sub>1</sub> d <sub>1</sub> ...
11	—	MBCF_WRITE	C3のアドレス	N	c <sub>1</sub> d <sub>1</sub> ...
12	11	MBCF_WRITE	F3のアドレス	4	1

10

20

【図 11】

送信数カウント	種数リンク同期情報	コマンド	対象アドレス	サイズ(バイト数)	データ列
1	—	MBCF_WRITE	A1のアドレス	N	a <sub>1</sub> d <sub>1</sub> ...
2	—	MBCF_WRITE	A2のアドレス	N	a <sub>1</sub> d <sub>1</sub> ...
3	—	MBCF_WRITE	A3のアドレス	N	a <sub>1</sub> d <sub>1</sub> ...
4	—	MBCF_WRITE	B1のアドレス	N	b <sub>1</sub> d <sub>1</sub> ...
5	—	MBCF_WRITE	B2のアドレス	N	b <sub>1</sub> d <sub>1</sub> ...
6	—	MBCF_WRITE	B3のアドレス	N	b <sub>1</sub> d <sub>1</sub> ...
7	—	MBCF_WRITE	C1のアドレス	N	c <sub>1</sub> d <sub>1</sub> ...
8	—	MBCF_WRITE	C2のアドレス	N	c <sub>1</sub> d <sub>1</sub> ...
9	—	MBCF_WRITE	C3のアドレス	N	c <sub>1</sub> d <sub>1</sub> ...
10	7	MBCF_WRITE	F10のアドレス	4	1
11	8	MBCF_WRITE	F20のアドレス	4	1
12	9	MBCF_WRITE	F30のアドレス	4	1

30

40

50

## フロントページの続き

(72)発明者 松本 尚

新潟県新潟市西区真砂一丁目15番36号 株式会社情報科学研究所内

審査官 三坂 敏夫

(56)参考文献

特開2014-191497(JP, A)

特表2010-510590(JP, A)

松本 尚, メモリーベース通信ファシリティの評価, 電子情報通信学会技術研究報告, 社団法人電子情報通信学会, 2001年09月27日, Vol. 101 No. 329, 第31頁 - 第40頁, ISSN:0913-5685

松本 尚 他, 汎用超並列オペレーティングシステム: SSS-CORE, 情報処理学会研究報告, 社団法人情報処理学会, 1996年08月27日, Vol. 96 No. 79, 第115頁 - 第120頁, ISSN:0919-6072

丹羽 純平 他, 汎用超並列オペレーティングシステムSSS-CORE上の非対称分散共有メモリにおけるコンパイル技法, コンピュータソフトウェア, 日本ソフトウェア科学会, 1998年05月15日, Vol. 15 No. 3, 第242頁 - 第246頁, ISSN:0289-6540

(58)調査した分野 (Int.Cl., DB名)

G06F 15/16 - 15/177