

(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2009-282917

(P2009-282917A)

(43) 公開日 平成21年12月3日(2009.12.3)

(51) Int.Cl.
G06F 15/173 (2006.01)

F I
G06F 15/173 650S

テーマコード (参考)
5B045

審査請求 未請求 請求項の数 6 O L (全 12 頁)

(21) 出願番号 特願2008-136943 (P2008-136943)
(22) 出願日 平成20年5月26日 (2008.5.26)

(71) 出願人 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
(74) 代理人 110000442
特許業務法人 武和国際特許事務所
(72) 発明者 ▲高▼瀬 亮
神奈川県秦野市堀山下1番地 株式会社日立製作所エンタープライズサーバ事業部内
(72) 発明者 情野 雄太郎
神奈川県秦野市堀山下1番地 株式会社日立製作所エンタープライズサーバ事業部内
(72) 発明者 藤原 至誠
神奈川県秦野市堀山下1番地 株式会社日立製作所エンタープライズサーバ事業部内
Fターム(参考) 5B045 BB15 BB28 BB34

(54) 【発明の名称】 サーバ間通信機構及びコンピュータシステム

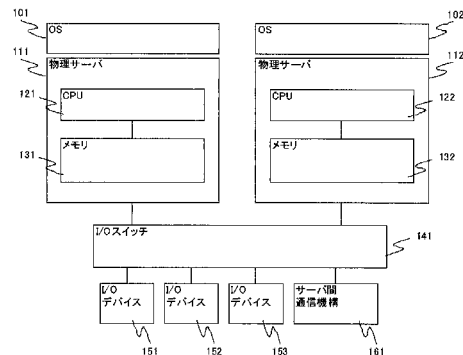
(57) 【要約】

【課題】物理サーバ毎にサーバ間通信を行うための外付けのI/Oデバイスを用意する必要がなく、プロトコル変換によるオーバーヘッドも生じることのないようにする。

【解決手段】サーバ間通信機構161には、I/Oリンク201、図示しないI/Oスイッチを介して複数の物理サーバに接続されている。サーバ間通信機構161は、各物理サーバのデータをアクセスするための命令を発行する読み出し命令生成部203、読み出したデータを他のサーバに送信する書き込み命令生成部204を有し、サーバ間通信機構161の内部でデータの送信元からデータを読み出し、読み出したデータをそのままデータの送信先に書き出して、サーバ間通信機構161でデータを直接折り返すことにより物理サーバ相互間でのデータ転送を行う。

【選択図】 図1

図1



【特許請求の範囲】**【請求項 1】**

メモリの内容を読み出すための命令を生成する読み出し命令生成手段と、前記読み出し命令の結果として返却されるメモリデータ返却命令を受信する返却命令受信手段と、メモリデータ返却命令と共に返却されるメモリデータをバッファリングするデータバッファと、バッファリングされたメモリデータを書き込むための命令を生成する書き込み命令生成手段と、前記読み出し命令及び前記書き込み命令の宛先情報を付加する宛先情報付加手段とを備え、I/Oスイッチを介して接続されている複数の物理サーバにおける、データ送信元の物理サーバのメモリ上にあるデータを送信先の物理サーバのメモリ上に転送することを特徴とするサーバ間通信機構。

10

【請求項 2】

請求項 1 記載のサーバ間通信機構において、前記サーバ間通信機構は、制御レジスタを内蔵しており、前記制御レジスタは、複数の物理サーバから共通に読み出し、あるいは、書き込みが行われることを特徴とするサーバ間通信機構。

【請求項 3】

請求項 1 記載のサーバ間通信機構であって、該サーバ間通信機構は、複数のアップストリームポートと 1 つ以上のダウンストリームポートとを有する I/O スイッチに内蔵されていることを特徴とするサーバ間通信機構。

【請求項 4】

それぞれが 1 つ以上の CPU とメモリとを有する複数の物理サーバと、複数のアップストリームポートと 1 つ以上のダウンストリームポートを持つ I/O スイッチとからなり、前記複数の物理サーバのそれぞれが前記 I/O スイッチのアップストリームポートを経由して接続されたコンピュータシステムにおいて、

20

前記 I/O スイッチのダウンストリームポートには、請求項 1 記載のサーバ間通信機構が接続されていることを特徴とするコンピュータシステム。

【請求項 5】

それぞれが 1 つ以上の CPU とメモリとを有する複数の物理サーバと、複数のアップストリームポートと 1 つ以上のダウンストリームポートを持つ I/O スイッチとからなり、前記複数の物理サーバのそれぞれが前記 I/O スイッチのアップストリームポートを経由して接続されたコンピュータシステムにおいて、

30

前記 I/O スイッチには、請求項 1 記載のサーバ間通信機構が内蔵されており、前記ダウンストリームポートには I/O デバイスが接続されていることを特徴とするコンピュータシステム。

【請求項 6】

それぞれが 1 つ以上の CPU とメモリとを有する複数の物理サーバと、複数のアップストリームポートと 1 つ以上のダウンストリームポートを持つ I/O スイッチとからなり、前記複数の物理サーバのそれぞれが前記 I/O スイッチのアップストリームポートを経由して接続されたコンピュータシステムにおいて、

前記 I/O スイッチが多段に接続され、最上位の前記 I/O スイッチのアップストリームポートに前記物理サーバが接続され、多段に接続された前記 I/O スイッチのそれぞれには、請求項 1 記載のサーバ間通信機構が内蔵されていることを特徴とするコンピュータシステム。

40

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、サーバ間通信機構及びコンピュータシステムに係り、特に、2 台以上の物理サーバが I/O スイッチを介して接続されて構成されたコンピュータシステムにおける物理サーバ間の通信を行うためのサーバ間通信機構及びコンピュータシステムに関する。

【背景技術】**【0002】**

50

近年のコンピュータシステムは、CPU単体の高性能化及びCPUのマルチコア化等のCPUの処理性能の向上に伴って、1台の物理サーバとしてのコンピュータ上に複数の仮想サーバを稼働させるサーバ統合のニーズが高まってきている。このようなサーバ統合を行うことにより、1台の物理サーバで稼働するOSやアプリケーションの数を増大させることが可能になり、コンピュータシステムの高性能化を図ることができる。その結果、そのようなコンピュータシステムは、1台の物理サーバとしてのコンピュータに接続しなければならないI/Oデバイスが増加していくことが予測される。そして、より多くのI/Oデバイスを搭載するために、この種のコンピュータシステムは、サーバとI/Oデバイスとの間をPCI-Express(R)スイッチ等のI/Oスイッチを介して接続して構成する必要性が高まっている。

10

【0003】

前述したようなI/Oスイッチを用いてサーバに接続する物理的なI/Oデバイスの数を増加させるというアプローチの一方で、物理サーバ相互間あるいは仮想サーバ相互間でI/Oデバイスを共有するI/O仮想化の普及も予測される。I/O仮想化とは、物理的なI/Oデバイスの上に複数の仮想I/Oデバイスを構築し、各物理サーバあるいは各仮想サーバのそれぞれに仮想I/Oデバイスを割り当てることにより、I/Oデバイスを物理サーバ相互間あるいは仮想サーバ相互間で共有するという方法である。

【0004】

I/Oデバイスを複数の物理サーバで共有する場合、そのコンピュータシステムは、複数のアップストリームポートと複数のダウンストリームポートとを有するI/Oスイッチを用意し、このI/Oスイッチのアップストリームポートのそれぞれに物理サーバを接続し、ダウンストリームポートにI/Oデバイスを接続する構成とされる。これにより、I/Oデバイスを複数の物理サーバ相互間で共有することができる。このようなI/O仮想化を用いることにより、I/Oデバイスの物理的な数を増やすことなく、統合されたサーバ上で稼働するOSやアプリケーションのそれぞれから多くのI/Oデバイスを使用することが可能になる。

20

【0005】

また、サーバ統合を行って1台の物理サーバ上に複数の仮想サーバを稼働させるようにしたコンピュータシステムは、ある物理サーバ上で稼働する仮想サーバを、稼働中に別の物理サーバに集約したり再構成したりしたいという要求がある。仮想サーバをある物理サーバから別の物理サーバに移動させるためには、その仮想サーバが使用しているメモリの内容や稼働状態を、そのまま別の物理サーバに引き継ぐということを意味する。すなわち、仮想サーバをある物理サーバから別の物理サーバに移動させるためには、このような稼働に関するメモリ上の大量のデータを高速に転送することが必要となり、物理サーバ間での高速な通信手段が必要となる。

30

【0006】

コンピュータ相互間での高速な通信を行うことを可能とした通信手段に関する従来技術として、例えば、特許文献1等に記載された技術が知られている。この従来技術は、マルチノードコンピュータシステムにおいて、ノード間通信に汎用I/Oインタフェースを用い、各ノードに存在する通信制御装置がデータ転送のための転送コマンドを解釈して、汎用I/Oインタフェースの制御を行ってノード間のデータ転送を高速に行うことを可能にするというものである。

40

【特許文献1】特開2006-58956号公報

【発明の開示】

【発明が解決しようとする課題】

【0007】

前述した従来技術は、物理サーバ相互間の通信に適用した場合、各物理サーバにサーバ間通信用のI/Oデバイスを接続し、各物理サーバの通信用I/Oデバイス相互間で通信を行うことになり、これにより、サーバ間での高速通信を実現することが可能になる。このような手法を用いて、複数の物理サーバの全てがその相互間で通信を行うためには、各

50

物理サーバのそれぞれに通信用の I / O デバイスを接続する必要がある。

【 0 0 0 8 】

このため、前述した従来技術は、ブレードサーバに代表される物理サーバの集積度が高いサーバに適用した場合に、各物理サーバ間の通信を実現するために必要な通信用の I / O デバイスの数が非常に多くなってしまおうという問題点を生じる。

【 0 0 0 9 】

また、前述した従来技術は、それを適用し、通信用の I / O デバイスを用いてサーバ間通信を実現した場合、物理サーバと通信用 I / O デバイスとの間のインタフェースと、通信用 I / O デバイス相互間のインタフェースとで用いられる通信プロトコルが異なるため、サーバ間通信を行うためにプロトコル変換のためのオーバーヘッドが生じ、通信スループットが低下したり、通信レイテンシが大きくなったりするという問題点が生じてしまう。

10

【 0 0 1 0 】

前述した通信用の I / O デバイスの数が非常に多くなるという問題点は、I / O 仮想化の技術を利用してサーバ間通信用の I / O デバイスを物理サーバ間で共有するようにすることにより、I / O デバイスの数を低減することができ、問題をある程度解決することが可能であるが、プロトコル変換によるオーバーヘッドが生じるという問題を解決することができない。

【 0 0 1 1 】

本発明の目的は、前述した従来技術の問題点を解決し、物理サーバ毎にサーバ間通信を行うための外付けの I / O デバイスを用意する必要がなく、プロトコル変換によるオーバーヘッドが生じることのないサーバ間通信機構及び該サーバ間通信機構を用いたコンピュータシステムを提供することにある。

20

【課題を解決するための手段】

【 0 0 1 2 】

本発明によれば前記目的は、メモリの内容を読み出すための命令を生成する読み出し命令生成手段と、前記読み出し命令の結果として返却されるメモリデータ返却命令を受信する返却命令受信手段と、メモリデータ返却命令と共に返却されるメモリデータをバッファリングするデータバッファと、バッファリングされたメモリデータを書き込むための命令を生成する書き込み命令生成手段と、前記読み出し命令及び前記書き込み命令の宛先情報を付加する宛先情報付加手段とを備え、I / O スイッチを介して接続されている複数の物理サーバにおける、データ送信元の物理サーバのメモリ上にあるデータを送信先の物理サーバのメモリ上に転送することにより達成される。

30

【発明の効果】

【 0 0 1 3 】

本発明によれば、物理サーバ毎にサーバ間通信を行うための外付けの I / O デバイスを用意する必要がなく、プロトコル変換によるオーバーヘッドが生じさせることを防止して、通信のスループットを増大させることができ、通信によるレイテンシの増大を防止することができる。

【発明を実施するための最良の形態】

【 0 0 1 4 】

以下、本発明によるサーバ間通信機構及びコンピュータシステムの実施形態を図面により詳細に説明する。

40

【 0 0 1 5 】

図 1 は本発明の第 1 の実施形態によるコンピュータシステムの構成を示すブロック図である。

【 0 0 1 6 】

本発明の第 1 の実施形態によるコンピュータシステムは、複数の物理サーバ 1 1 1、1 1 2 が、複数のアップストリームポートと複数のダウンストリームポートとを有する I / O スイッチ 1 4 1 のアップストリームポートに接続され、I / O スイッチ 1 4 1 のダウンストリームポートに複数の I / O デバイス 1 5 1 ~ 1 5 3 とサーバ間通信機構 1 6 1 とが

50

接続され、物理サーバ111、112の上でOS101、102が稼働するように構成されている。物理サーバ111は、CPU121及びメモリ131を有して構成され、また、物理サーバ112も、同様に、CPU122及びメモリ132を有して構成されている。

【0017】

すなわち、本発明の第1の実施形態によるコンピュータシステムは、I/Oスイッチ141を経由して物理サーバ111、112と、I/Oデバイス151～153とが接続されて構成されている。I/Oデバイス151～153のそれぞれは、物理サーバ111と112との両方で共有して使用されるI/Oデバイスであっても、物理サーバ111または112のどちらかで占有して使用されるI/Oデバイスであってもよい。また、I/Oスイッチ141のダウンストリームポートには、本発明によるサーバ相互間の通信に利用するサーバ間通信機構161が接続されており、I/Oスイッチ141を経由して物理サーバ111及び112と接続されている。

10

【0018】

前述において、本発明の実施形態では、OSが稼働している物理サーバを2台、サーバ間通信機構を1台備えるとして示しているが、本発明は、物理サーバがさらに多数設けられてもよく、また、サーバ間通信機構を2台以上設けてもよい。これにより、2台のサーバの組でサーバ相互間の通信を行っている場合にも、他の2台のサーバの組でサーバ相互間の通信を並行して行うことが可能となる。

【0019】

図2はサーバ間通信機構161の構成を示すブロック図である。サーバ間通信機構161は、I/Oリンク201によってI/Oスイッチ141と接続されている。I/Oリンク201とサーバ間通信機構161との内部は、I/Oインタフェース202により接続されている。

20

【0020】

サーバ間通信機構161は、物理サーバ内のメモリデータを取り込むためのメモリ読み出し命令生成部203と、メモリデータを他の物理サーバに送り出すためのメモリ書き込み命令生成部204と、割り込みを発生させるための割り込み命令生成部205とを有している。メモリ読み出し命令生成部203、メモリ書き込み命令生成部204、割り込み命令生成部205には、命令を正しい宛先の物理サーバに送出するための宛先情報付加機構として、それぞれ送信側サーバ宛先情報付加部206と受信側サーバ宛先情報付加部207とが接続されている。また、サーバ間通信機構161は、メモリ読み出し命令生成部203、メモリ書き込み命令生成部204及び割り込み命令生成部205の動作を制御するために命令発行用シーケンサ208を備えている。また、サーバ間通信機構161は、メモリ読み出し命令により、物理サーバから読み出されたデータを取り込むためのメモリデータ返却命令受信部209を備えると共に、ここで受信したメモリデータを格納するためのメモリデータバッファ210を備えている。さらに、サーバ間通信機構161は、ソフトウェアからサーバ間通信機構を制御するための機構として、サーバ間通信機構レジスタ211を備えている。このサーバ間通信機構レジスタ211には、送信メモリアドレスレジスタ212、受信メモリアドレスレジスタ213、送信メモリ領域長レジスタ214、起動レジスタ215が含まれている。

30

40

【0021】

図3はサーバ間通信機構が物理サーバ相互間でメモリデータの転送を行う制御を説明するシーケンスチャートであり、次に、これについて説明する。ここでのメモリデータの転送は、図1に示すコンピュータシステムの例における物理サーバ111内のメモリ131に存在するデータを、物理サーバ112内のメモリ132に転送する場合を例としている。

【0022】

(1)まず、データ送信側となる物理サーバ111上で稼働しているOS101が、送信するメモリ領域の先頭アドレスを設定する。この設定は、データ送信側となる物理サー

50

バ 1 1 1 からサーバ間通信機構 1 6 1 の送信メモリアドレスレジスタ 2 1 2 に対して書き込み命令を送信することにより行われる (ステップ 3 0 1)。

【 0 0 2 3 】

(2) 同様に、データ送信側となる物理サーバ 1 1 1 上で稼働している OS 1 0 1 は、サーバ間通信機構 1 6 1 の送信メモリ領域長レジスタ 2 1 4 に対して書き込みを行い、送信するメモリ領域の大きさを設定する (ステップ 3 0 2)。

【 0 0 2 4 】

(3) また、同様に、データ受信側となる物理サーバ 1 1 2 上で稼働している OS 1 0 2 が、サーバ間通信機構 1 6 1 の受信メモリアドレスレジスタ 2 1 3 に対して書き込みを行い、受信するメモリ領域の先頭アドレスを設定する (ステップ 3 0 3)。

10

【 0 0 2 5 】

前述までの処理でサーバ間通信を行うための初期設定処理が完了する。

【 0 0 2 6 】

(4) 次に、データ送信側となる物理サーバ 1 1 1 上で稼働している OS 1 0 1 が、サーバ間通信機構 1 6 1 の起動レジスタ 2 1 5 に書き込みを行うことにより、サーバ間通信機構 1 6 1 を起動させる (ステップ 3 0 4)。

【 0 0 2 7 】

(5) サーバ間通信機構 1 6 1 は、起動されると命令発行用シーケンサ 2 0 8 が動作を開始し、メモリ読み出し命令生成部 2 0 3 から送信側サーバ宛先情報付加部 2 0 6 を経由してデータ送信側の物理サーバ 1 1 1 の宛先情報を付加し、メモリ読み出し命令を発行する。このメモリ読み出し命令は、I/Oスイッチ 1 4 1 を通過する際に前記の送信側サーバ宛先情報を用いてデータ送信側の物理サーバ 1 1 1 へと正しく送信される (ステップ 3 0 5 a)。

20

【 0 0 2 8 】

(6) メモリ読み出し命令を受け取ったデータ送信側の物理サーバ 1 1 1 は、サーバ間通信機構 1 6 1 に対して、送信するメモリデータを含むデータ返却命令を送信することにより、ステップ 3 0 5 a でのメモリ読み出し命令に対するメモリデータ返却を行う。なお、1 回のデータ返却命令で送信するメモリデータの量は、予め定められた量であり、送信すべきメモリ領域内のデータ量が大きい場合には、複数回に分けて送信される (ステップ 3 0 6 a)。

30

【 0 0 2 9 】

(7) サーバ間通信機構 1 6 1 は、データ返却命令をメモリデータ返却命令受信部 2 0 9 で受信し、メモリデータ部分をメモリデータバッファ 2 1 0 に格納する。そして、サーバ間通信機構 1 6 1 は、メモリデータを受け取ると、メモリ書き込み命令生成部 2 0 4 がメモリデータバッファ 2 1 0 からメモリデータを取り出し、受信側サーバ宛先情報付加部 2 0 7 を経由してデータ受信側の物理サーバ 1 1 2 の宛先情報を付加し、転送すべきメモリデータを含むメモリ書き込み命令を発行する。このメモリ書き込み命令は、I/Oスイッチ 1 4 1 を通過する際に前記の受信側サーバ宛先情報を用いてデータ受信側の物理サーバ 1 1 2 へと正しく送信される (ステップ 3 0 7 a)。

【 0 0 3 0 】

(8) 命令発行用シーケンサ 2 0 8 は、前述したステップ 3 0 5 a でのメモリ読み出し命令の送信、ステップ 3 0 6 a でのメモリ返却命令の受信、ステップ 3 0 7 a でのメモリ書き込み命令の送信の一連の動作を、送信メモリ領域長レジスタ 2 1 4 で指定されたデータ長を転送し終えるまで繰り返すことにより、データ送信側の物理サーバ 1 1 1 からデータ受信側の物理サーバ 1 1 2 へデータの転送を行う (ステップ 3 0 5 b、3 0 6 b、3 0 7 b)。

40

【 0 0 3 1 】

前述したステップ 3 0 5 b、3 0 6 b、3 0 7 b の処理を指定されたデータ長を転送し終えるまで繰り返すことにより、データ送受信の処理が完了する。

【 0 0 3 2 】

50

(9) データの転送が終了すると、サーバ間通信機構 161 の命令発行用シーケンサ 208 は、割り込み命令生成部 205 を起動する。割り込み命令生成部 205 は、データ送信側の物理サーバ 111 及びデータ受信側の物理サーバ 112 のそれぞれに対して、データ送信完了を通知する割り込み命令を発行する。割り込み命令生成部 205 で生成された割り込み命令は、送信側サーバ宛先情報付加部 206 あるいは受信側サーバ宛先情報付加部 207 を経由して正しい宛先情報が付加され、それぞれデータ送信側の物理サーバ 111、データ受信側の物理サーバ 112 へ送信される (ステップ 308、309)。

【0033】

前述したステップ 308、309 の処理を行うことにより、サーバ間でのデータ転送処理の全てが完了する。

【0034】

前述した本発明の第 1 の実施形態では、サーバ間通信機構 161 に設定指示をしたり、サーバ間通信機構を起動したりする動作主体を OS であるとして説明したが、動作主体は、サーバ間通信機構用のデバイスドライバであっても、アプリケーションプログラム等であっても、仮想サーバを管理するハイパバイザ等であってもよい。

【0035】

また、図 3 に示す処理を開始する契機となるのは、図 1 に示していないが、物理サーバ 111、112 に接続されている管理用のコンピュータからの指示であっても、また、物理サーバ 111、112 の中に構築されているサービスプロセッサからの指示であってもよい。また、この指示は、管理用のコンピュータ、サービスプロセッサが、物理サーバの状態から自動的に行われるものであっても、システム管理者等から行われるものであってもよい。

【0036】

図 4 は本発明の第 2 の実施形態によるコンピュータシステムの構成を示すブロック図である。図 4 に示す本発明の第 2 の実施形態は、I/O スイッチ内にサーバ間通信機構を設けるようにした例である。

【0037】

図 1 ~ 図 3 により説明した本発明の第 1 の実施形態は、サーバ間通信機構 161 を外付けの I/O デバイスとして、I/O スイッチ 141 のダウンストリームポートに接続するものとして説明したが、本発明の第 2 の実施形態は、図 4 に示すように、I/O スイッチにサーバ間通信機構 421 を内蔵させて通信機構内蔵 I/O スイッチ 411 として構成し、物理サーバ 401 及び 402 相互間のデータ転送を行うようにしたものである。

【0038】

本発明の第 2 の実施形態は、このように構成することにより、通信機構内蔵 I/O スイッチ 411 内のサーバ間通信機構 421 により、図 3 により説明した場合と同様に、物理サーバ 401、402 相互間の通信を行うことが可能となる。

【0039】

前述した本発明の第 2 の実施形態は、サーバ間通信機構 421 を接続するためにのために、I/O スイッチのダウンストリームスロットを占有する必要がなくなり、他のデバイスのために I/O スイッチのダウンストリームスロットを開放することができるという利点を得ることができる。また、サーバ間通信機構 421 を外付けのデバイスとして用意する必要がないため、サーバ間通信を導入するコストをより低減することができるという利点を得ることができる。

【0040】

図 5 は本発明の第 3 の実施形態によるコンピュータシステムの構成を示すブロック図である。図 5 に示す本発明の第 3 の実施形態は、物理サーバに接続する I/O デバイスの数を増加するために、通信機構内蔵 I/O スイッチを多段に設けて構成した場合のコンピュータシステムの例である。

【0041】

図 5 に示す本発明の第 3 の実施形態によるコンピュータシステムは、複数の物理サーバ

10

20

30

40

50

501、502が1段目の通信機構内蔵I/Oスイッチのアップストリームポートに接続され、複数の2段目の通信機構内蔵I/Oスイッチ512、513が1段目の通信機構内蔵I/Oスイッチのダウンストリームポートに接続され、2段目の通信機構内蔵I/Oスイッチ512、513のそれぞれの複数のダウンストリームポートにI/Oデバイスが接続されて構成されている。そして、物理サーバ501と502との間のサーバ間通信は、初段の通信機構内蔵I/Oスイッチ511に内蔵したサーバ間通信機構521を用いて行うことができる。

【0042】

前述の本発明の第3の実施形態によれば、初段のI/Oスイッチ内のサーバ間通信機構を用いてサーバ間通信を行うことができるので、図6により後述するI/Oスイッチにサーバ間通信機構を内蔵していないI/Oスイッチを多段に設けて構成した例に比較して、サーバ間通信に要する通信レイテンシを低減することができる。

10

【0043】

図5に示す本発明の第3の実施形態によるコンピュータシステムは、初段の通信機構内蔵I/Oスイッチ511を1台だけ備えて構成されているが、この例では、初段の通信機構内蔵I/Oスイッチ511を複数台設けることも可能であり、また、2段目の通信機構内蔵I/Oスイッチ512、513を3台以上設けて構成することができる。このように構成された本発明の第3の実施形態は、2段目の通信機構内蔵I/Oスイッチの1つに複数の初段の通信機構内蔵I/Oスイッチを接続するような構成をとることにより、異なる初段の通信機構内蔵I/Oスイッチ内のサーバ間通信機構を用いては通信を行うことができない物理サーバ間の通信を、2段目の通信機構内蔵I/Oスイッチに内蔵したサーバ間通信機構を用いて行うことができる。

20

【0044】

図6は本発明の第4の実施形態によるコンピュータシステムの構成を示すブロック図である。図6に示す本発明の第4の実施形態は、物理サーバに接続するI/Oデバイスの数を増加させるために、I/Oスイッチとして、サーバ間通信機構を内蔵していないI/Oスイッチを多段(図示例では2段)に設け、2段目のI/Oスイッチにサーバ間通信機構を接続して構成した例である。

【0045】

すなわち、図6に示す例は、複数の物理サーバ501、502を1段目のI/Oスイッチ611に接続し、このI/Oスイッチ611に複数の2段目のI/Oスイッチ612、613を接続し、2段目のI/Oスイッチ612、613にI/Oデバイスとサーバ間通信機構を接続して構成されている。

30

【0046】

図6に示す本発明の第4の実施形態によるコンピュータシステムは、I/Oスイッチにサーバ間通信機構を内蔵していないI/Oスイッチを多段に設けて構成しているため、物理サーバ501と502との間のサーバ間通信は、2段目のI/Oスイッチ612に接続されたサーバ間通信機構622で行う必要がある。なお、図6には示していないが、2段目のI/Oスイッチ613にもサーバ間通信機構を接続することができる。

【0047】

前述した本発明の各実施形態は、サーバ間通信機構によるデータ送受信の処理の完了を、割り込み命令生成部からデータ送信側の物理サーバ及びデータ受信側の物理サーバに割り込み命令を発行することにより、各物理サーバに通知するものであったが、本発明は、データ送信処理の完了を他の方法により、各物理サーバに通知するようにすることもできる。

40

【0048】

図7はサーバ間通信機構が物理サーバ相互間でメモリデータの転送を制御する他の例を説明するシーケンスチャートであり、次に、これについて説明する。ここでのメモリデータの転送は、図3に示して説明した場合と同様に、図1に示すコンピュータシステムの場合における物理サーバ111内のメモリ131に存在するデータを、物理サーバ112内の

50

メモリ 132 に転送する場合を例としている。

【0049】

そして、この例においては、サーバ間通信機構 161 内のサーバ間通信機構レジスタ 211 に、データ送信側の物理サーバからもデータ受信側の物理サーバからも共通して読み出すことが可能な完了ステータスレジスタを設け、データ送受信処理を終了すると、サーバ間通信機構 161 が、サーバ間通信機構レジスタ 211 内に設けた完了ステータスレジスタに完了を登録するようにし、物理サーバが、完了ステータスレジスタへのポーリングにより、完了ステータスレジスタを読み出すことにより、送信側及び受信側の物理サーバがデータの転送の終了を知るようにコンピュータシステムが構成される。

【0050】

図 7 に示すシーケンスにおけるステップ 301 ~ 307 b までの処理は、図 3 に示して説明したと同一であるので、ここまでの説明については省略する。

【0051】

(1) ステップ 301 ~ 307 b までの処理でデータの送受信処理が完了すると、サーバ間通信機構 161 が、サーバ間通信機構レジスタ 211 内に設けた完了ステータスレジスタに完了を登録する。その後、サーバ間通信機構 161 は、データ送信側物理サーバ 111 から完了ステータスレジスタの読み出し命令を受け取って、完了ステータスレジスタの内容をデータ送信側物理サーバ 111 に返却する(ステップ 701、702)。

【0052】

(2) 同様に、サーバ間通信機構 161 は、データ受信側物理サーバ 112 から完了ステータスレジスタの読み出し命令を受け取って、完了ステータスレジスタの内容をデータ受信側物理サーバ 112 に返却する(ステップ 703、704)。

【0053】

前述では、完了ステータスレジスタを用いる例を説明したが、本発明は、サーバ通信機構 161 内の単一のレジスタに対して、データ送信側物理サーバ 111 からデータ受信側物理サーバ 112 から同様に読み出したりは書き込みといったアクセスを可能とすることにより、サーバ間通信機構で単一のレジスタの状態を変更することにより、データ送信側物理サーバ及びデータ受信側物理サーバの両者にその状態を通知することができる。

【0054】

前述した本発明の実施形態での各処理は、プログラムにより構成し、本発明が備えるサーバ間通信機構に実行させることができ、また、それらのプログラムは、FD、CDROM、DVD等の記録媒体に格納して提供することができ、また、ネットワークを介してデジタル情報により提供することができる。

【0055】

前述した本発明の実施形態によれば、物理サーバ間のデータ転送を行う際に、従来は物理サーバ毎に用意していたサーバ間通信用の外付け I/O デバイスの数を削減することができ、また、サーバ間通信機構を I/O スイッチに内蔵することにより、サーバ間通信用の外付け I/O デバイスを不要とすることができる。

【0056】

また、本発明の実施形態によれば、物理サーバ間のデータ通信は、サーバ間通信機構の内部で行うことができるため、プロトコル変換のオーバーヘッドが生じることがなく、通信スループットやレイテンシの面で有利となる。

【0057】

さらに、本発明の実施形態によれば、サーバ間通信機構を I/O スイッチに内蔵することにより、物理サーバ毎にサーバ間通信を行うための外付けの I/O デバイスを用意する必要がなく、サーバ間通信を行うための導入コストを抑えることができる。また、I/O スイッチのスロットがサーバ間通信用の I/O デバイスで占有されることがなく、I/O スロットを有効に活用することが可能となる。また、I/O スイッチを多段の構成にして I/O デバイスの数を増加させたシステムにおいて、中継段の I/O スイッチにもサーバ

10

20

30

40

50

間通信機構を内蔵させることができるため、中継段の I / O スイッチに接続された物理サーバ間のサーバ間通信は、中継段の I / O スイッチ内蔵のサーバ間通信機構で折り返すことが可能であり、通信レイテンシをより低減することができる。

【 0 0 5 8 】

また、本発明の実施形態によれば、アドレス変換を必要とせず大量のメモリデータをサーバ間で転送するような用途で特に効果がある。例えば、仮想サーバのメモリイメージは、大容量であり、かつ、物理アドレスで連続領域であることがある。このような場合に、ハイパバイザが本発明によるサーバ間通信機構を使用して、仮想サーバのマイグレーションを物理サーバ間で行うといった用途に本発明を適用して大きな効果を得ることができる。

10

【 図面の簡単な説明 】

【 0 0 5 9 】

【 図 1 】 本発明の第 1 の実施形態によるコンピュータシステムの構成を示すブロック図である。

【 図 2 】 サーバ間通信機構の構成を示すブロック図である。

【 図 3 】 サーバ間通信機構が物理サーバ相互間でメモリデータの転送を行う制御を説明するシーケンスチャートである。

【 図 4 】 本発明の第 2 の実施形態によるコンピュータシステムの構成を示すブロック図である。

【 図 5 】 本発明の第 3 の実施形態によるコンピュータシステムの構成を示すブロック図である。

20

【 図 6 】 本発明の第 4 の実施形態によるコンピュータシステムの構成を示すブロック図である。

【 図 7 】 サーバ間通信機構が物理サーバ相互間でメモリデータの転送を制御する他の例を説明するシーケンスチャートである。

【 符号の説明 】

【 0 0 6 0 】

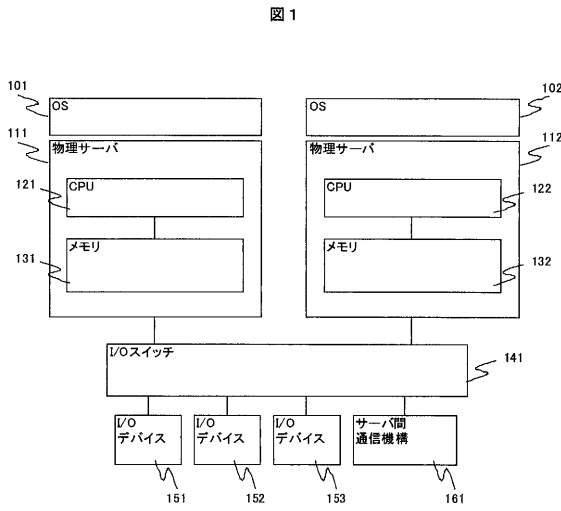
- 1 0 1、1 0 2 O S
- 1 1 1、1 1 2、4 0 1、4 0 2、5 0 1、5 0 2 物理サーバ
- 1 2 1、1 2 2 C P U
- 1 3 1、1 3 2 メモリ
- 1 4 1、6 1 1、6 1 2、6 1 3 I / O スイッチ
- 1 5 1 ~ 1 5 3 I / O デバイス
- 1 6 1、4 2 1、5 1 1、5 2 2、5 2 3、6 2 2 サーバ間通信機構
- 2 0 1 I / O リンク
- 2 0 2 I / O インタフェース
- 2 0 3 メモリ読み出し命令生成部
- 2 0 4 メモリ書き込み命令生成部
- 2 0 5 割り込み命令生成部
- 2 0 6 送信側サーバ宛先情報付加部
- 2 0 7 受信側サーバ宛先情報付加部
- 2 0 8 命令発行用シーケンス
- 2 0 9 メモリデータ返却命令受信部
- 2 1 0 メモリデータバッファ
- 2 1 1 サーバ間通信機構レジスタ
- 2 1 2 送信メモリアドレスレジスタ
- 2 1 3 受信メモリアドレスレジスタ
- 2 1 4 送信メモリ領域長レジスタ
- 2 1 5 起動レジスタ
- 4 1 1、5 1 1、5 1 2、5 1 3 通信機構内蔵 I / O スイッチ

30

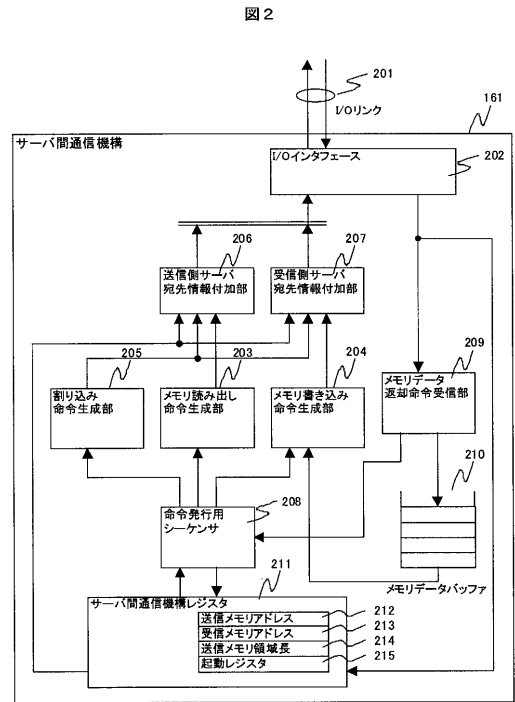
40

50

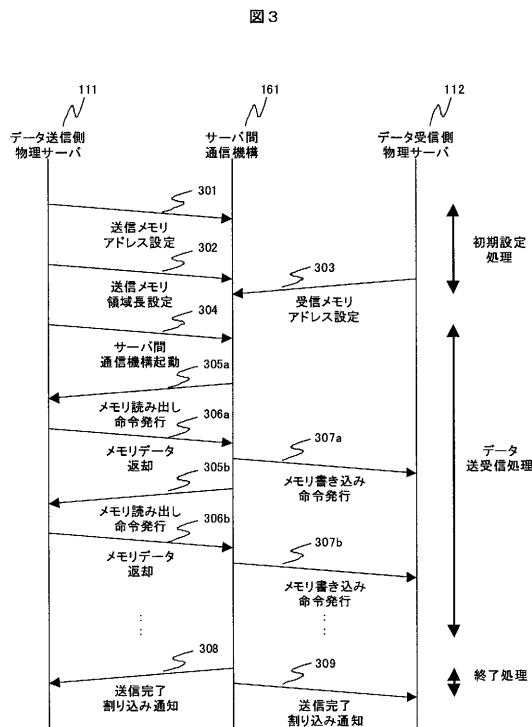
【 図 1 】



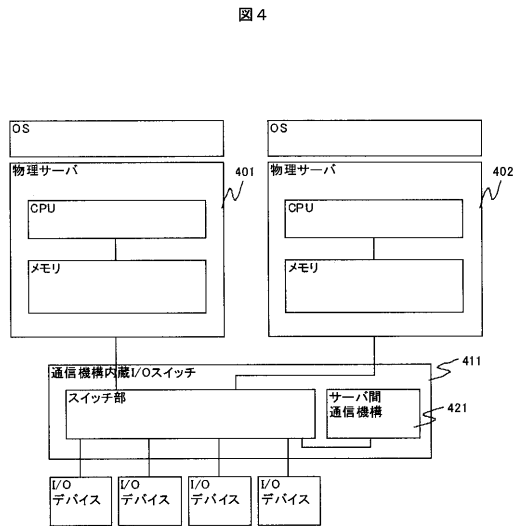
【 図 2 】



【 図 3 】

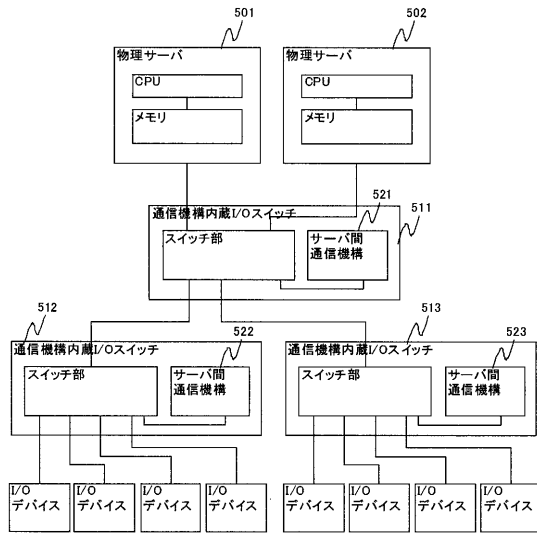


【 図 4 】



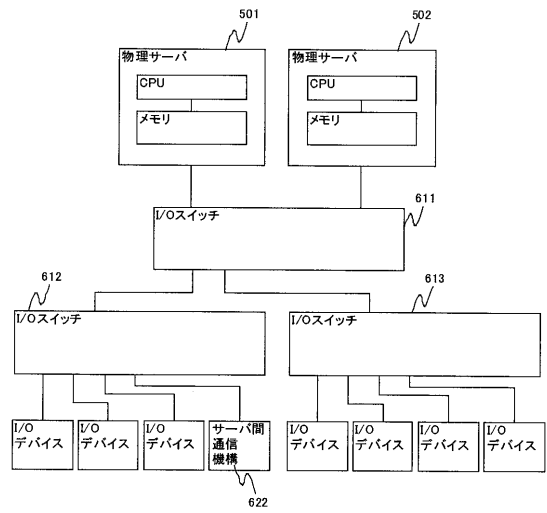
【 図 5 】

図 5



【 図 6 】

図 6



【 図 7 】

図 7

