

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4001889号

(P4001889)

(45) 発行日 平成19年10月31日(2007.10.31)

(24) 登録日 平成19年8月24日(2007.8.24)

(51) Int. Cl.		F I		
HO 4 L	29/08	(2006.01)	HO 4 L	13/00 3 O 7 Z
HO 4 L	12/56	(2006.01)	HO 4 L	12/56 1 O O A

請求項の数 10 (全 13 頁)

(21) 出願番号	特願2004-567360 (P2004-567360)	(73) 特許権者	390009531
(86) (22) 出願日	平成15年10月28日(2003.10.28)		インターナショナル・ビジネス・マシー ズ・コーポレーション
(65) 公表番号	特表2006-514454 (P2006-514454A)		I N T E R N A T I O N A L B U S I N E S S M A S C H I N E S C O R P O R A T I O N
(43) 公表日	平成18年4月27日(2006.4.27)		アメリカ合衆国10504 ニューヨーク 州 アーモンク ニュー オーチャード ロード
(86) 国際出願番号	PCT/GB2003/004645	(74) 代理人	100086243
(87) 国際公開番号	W02004/068801		弁理士 坂口 博
(87) 国際公開日	平成16年8月12日(2004.8.12)	(74) 代理人	100091568
審査請求日	平成18年7月28日(2006.7.28)		弁理士 市位 嘉宏
(31) 優先権主張番号	0302117.7	(74) 代理人	100108501
(32) 優先日	平成15年1月30日(2003.1.30)		弁理士 上野 剛史
(33) 優先権主張国	英国 (GB)		

最終頁に続く

(54) 【発明の名称】 ネットワークにおけるバッファ・データのプリエンプティブな再送

(57) 【特許請求の範囲】

【請求項1】

受信機に伝送されるデータ項目を格納するためのバッファと、
前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイ
マ機構であって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔
よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するも
のである、タイマ機構と、

第1の伝送用にデータにアクセスするため、および前記タイマ機構内の前記間隔に設定
された第1のタイムアウト・クロックを開始するための、バッファの第1のアクセス機構
と、

再送用に前記データにアクセスするため、および前記タイマ機構内の前記間隔に設定さ
れた少なくとも第2のタイムアウト・クロックを開始するための、前記タイマ機構による
前記信号発信にตอบสนองする前記バッファの第2のアクセス機構であって、前記第1のアクセ
ス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよ
う試行する、第2のアクセス機構と、

を有する、複数のパスを有するネットワークを介してデータを伝送するための装置。

【請求項2】

前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへの参照の
カウントを維持するための参照カウンタであって、前記第1のアクセス機構および前記第
2のアクセス機構によるそれぞれの参照時に増分され、各データ伝送の完了時に減分され

るものであり、前記参照カウンタは前記カウントがゼロに達すると信号を発信する、参照カウンタと、

前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへのアクセスを読み取ることが可能なように適合され、前記カウントがゼロに達した際の前記参照カウンタの信号発信に応答して前記バッファをフリー・バッファ・プールに戻すように適合された、メモリ・マネージャと、

をさらに有する、請求項1に記載の装置。

【請求項3】

前記データの再送用に前記最適な間隔を決定するための分析機構と、

前記データの再送用の前記最適な間隔に前記タイマ機構を調整するための調整機構と、
をさらに有する、請求項1に記載の装置。

10

【請求項4】

前記分析機構は、前記データの再送用に前記最適な間隔を決定するためにネットワーク監視データを使用するように動作可能である、請求項3に記載の装置。

【請求項5】

受信機に伝送されるデータ項目を格納するためのバッファと、

前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構であって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構と、

20

第1の伝送用にデータにアクセスするため、および前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するための、バッファの第1のアクセス機構と、

再送用に前記データにアクセスするため、および前記タイマ機構内の前記間隔に設定された少なくとも第2のタイムアウト・クロックを開始するための、前記タイマ機構による前記信号発信に응答する前記バッファの第2のアクセス機構であって、前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよう試行する、第2のアクセス機構と、

を有する、複数のパスを有するネットワークを介してデータを伝送するための装置を含むストレージ・コントローラ。

30

【請求項6】

前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへの参照のカウントを維持するための参照カウンタであって、前記第1のアクセス機構および前記第2のアクセス機構によるそれぞれの参照時に増分され、各データ伝送の完了時に減分されるものであり、前記カウントがゼロに達すると信号を発信する、参照カウンタと、

前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへのアクセスを読み取ることが可能なように適合され、前記カウントがゼロに達した際の前記参照カウンタの信号発信に응答して前記バッファをフリー・バッファ・プールに戻すように適合された、メモリ・マネージャと、

をさらに有する、請求項5に記載のストレージ・コントローラ。

40

【請求項7】

前記データの再送用に前記最適な間隔を決定するための分析機構と、

前記データの再送用の前記最適な間隔に前記タイマ機構を調整するための調整機構と、
をさらに有する、請求項5に記載のストレージ・コントローラ。

【請求項8】

受信機に伝送されるデータ項目を格納するためのバッファと、

前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構であって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構と、

50

第1の伝送用にデータにアクセスするため、および前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するための、バッファの第1のアクセス機構と、

再送用に前記データにアクセスするため、および前記タイマ機構内の前記間隔に設定された少なくとも第2のタイムアウト・クロックを開始するための、前記タイマ機構による前記信号発信に応答する前記バッファの第2のアクセス機構であって、前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよう試行する、第2のアクセス機構と、

を有する、複数のパスを有するネットワークを介してデータを伝送するための装置を含むネットワーク・アプライアンス。

10

【請求項9】

受信機に伝送されるデータ項目を格納するためのバッファを提供するステップと、

前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構を提供するステップであって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構を提供するステップと、

第1の伝送用に第1のアクセス機構によってバッファ内のデータにアクセスするステップと、

前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するステップと、

20

前記タイマ機構による前記信号発信に応答して第2のアクセス機構によってバッファ内のデータにアクセスするステップと、

前記第2のアクセス機構によって前記データを再送するステップと、

前記タイマ機構内の前記間隔に設定された第2のタイムアウト・クロックを開始するステップと、

前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するように前記第2のアクセス機構によって試行するステップと、

を有する、複数のパスを有するネットワークを介してデータを伝送するための方法。

【請求項10】

コンピュータ読取り可能メディア内で明白に具体化され、コンピュータ・システムにロードされ実行された場合に、

30

受信機に伝送されるデータ項目を格納するためのバッファを提供するステップと、

前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構を提供するステップであって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信する、タイマ機構を提供するステップと、

第1の伝送用に第1のアクセス機構によってバッファ内のデータにアクセスするステップと、

前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するステップと、

40

前記タイマ機構による前記信号発信に応答して第2のアクセス機構によってバッファ内のデータにアクセスするステップと、

前記第2のアクセス機構によって前記データを再送するステップと、

前記タイマ機構内の前記間隔に設定された第2のタイムアウト・クロックを開始するステップと、

前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するように前記第2のアクセス機構によって試行するステップと、

を実行するためのコンピュータ・プログラム・コード手段を有するコンピュータ・プログラム。

【発明の詳細な説明】

50

【技術分野】

【0001】

本発明は、分散ネットワーク型の耐障害 (fault tolerant) システムの分野に関し、とりわけ、データ伝送の速度および信頼性が重要なシステムに関する。

【背景技術】

【0002】

分散システムは、ストレージ・コントローラ・タイプの機能を含む様々な機能に関するプラットフォームを提供する手段として、ますます普及しつつある。その人気は、こうしたシステムが提供する柔軟性および拡張容易性に由来している。耐障害性は、冗長ネットワーク・インフラストラクチャまたは冗長ストレージ接続機構の提供などの、いくつかの相互にサポートし合う方法でインプリメントされる。分散アプリケーションは、それらの機能を実行するためのネットワークの接続性および通信機能に依存する。これらの耐障害機能がシステムの可用性を向上させる。多くのアプリケーションにとっては、可用性の向上もますます重要になってきている。

10

【0003】

多くのシステムでは、以下のように動作するそれらのネットワーク・インターフェース用の再試行アルゴリズムをインプリメントする。

1. パケットがドロップするなどのエラーが発生する。
2. タイムアウト間隔が満了する。
3. ネットワーク・ハードウェアまたはプロトコル・スタックがエラーを検出する。
4. 要求発行者にエラーが報告される。
5. 発行者が、代替ハードウェアを使用して2度目の要求を試行する。

20

【0004】

こうした方式は簡単であり、オリジナル要求の失敗まで待機することが受け入れ可能な場合は適切である。しかしながら、いくつかの重要な環境では、最小限の可用性 (サービスへのアクセスを試行している間に実際の障害がない) では不十分である。これらの環境では、一定時間内に応答を受け取ることが重要である。その時間内での応答に失敗すると、サービスにアクセスしている間の本格的な障害に匹敵するペナルティとなり、たとえば他のアプリケーションがタイムアウトしてエラー状態に戻るか、またはインターネット・ユーザがイライラしてクリック1つで競合相手のウェブ・サイトに行ってしまう可能性がある。

30

【0005】

したがって、所望の時間制限内で再試行が実行できる方法を提供することが有利となるが、既存のシステムでは本来これが実行できない。対処すべきいくつかの問題がある。第1に、システムはよりタイムリーな様式でエラーを検出しようとする可能性がある。インターフェース・アダプタまたはハードウェア内のタイムアウト間隔を短くすることは可能であるかもしれないが、これらの間隔をどこまで短くできるかについては、しばしばアーキテクチャ上の制限がある。たとえばファイバ・チャネルでは、正常に完了しなかった交換はスイッチによって定義されたエラー・タイムアウト間隔が満了するまで再使用できず、これはしばしば10秒間である。さらに、多くのネットワーク・インプリメンテーションは、ネットワークのエラー・タイムアウト間隔があまりに短すぎると、確実に動作しない。

40

【0006】

失敗したインターフェース・ソフトウェアまたはハードウェアに関連付けられていない何らかの他のタイムアウト機構を使用して、要求の再駆動を試行することが可能な場合があるが、オリジナルの要求が依然としてアクティブであるため、これでは問題の解決にならない。冗長ハードウェアによって提供された代替パスを使用してオリジナル要求の再駆動を試行しようとする、オリジナル要求に関連付けられた依然として使用中のリソースがあるため、オリジナル要求が完了するまでブロックされることになる。

【0007】

50

具体的な例として、ファイバ・チャネル・アダプタを使用して伝送インターフェースをインプリメントし、マルチスレッド・ユーザ・プロセスがバッファを再送しようとする場合、依然としてメモリはオリジナルの伝送によって使用中であるため、第2の伝送試行はブロックされることになる。

【0008】

他の可能な解決策は、専用バッファ内に伝送データのコピーを維持することによって、このメモリ・ブロック問題を避けようとするものであるかもしれない。その後、第1の伝送が長くかかり過ぎていとみなされた場合、このコピーを使用して第2の伝送を作成することができる。専用コピーには仮想メモリ・マネージャによってブロックされることなくアクセスすることができるが、この方式では、第1の伝送試行が実行される前に各伝送用のデータがコピーされなければならないため、問題に遭遇することのない大多数を含むあらゆる伝送のコストが増加することになる。

10

【発明の開示】

【発明が解決しようとする課題】

【0009】

こうした追加の処理コストは、ほとんどの現在のネットワークでは受け入れられないことがわかるであろう。

【課題を解決するための手段】

【0010】

したがって本発明は、第1の態様において、受信機に伝送されるデータ項目を格納するためのバッファと、前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構であって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構と、第1の伝送用にデータにアクセスするため、および前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するための、バッファの第1のアクセス機構と、再送用に前記データにアクセスするため、および前記タイマ機構内の前記間隔に設定された少なくとも第2のタイムアウト・クロックを開始するための、前記タイマ機構による前記信号発信に応答する前記バッファの第2のアクセス機構であって、前記第2のアクセス機構は前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよう試行するものである、第2のアクセス機構と、を有する、複数のパスを有するネットワークを介してデータを伝送するための装置を提供する。

20

30

【0011】

好ましくは、この装置は、前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへの参照のカウントを維持するための参照カウンタであって、前記参照カウンタは、前記第1のアクセス機構および前記第2のアクセス機構によるそれぞれの参照時に増分され、各データ伝送の完了時に減分されるものであり、前記参照カウンタは前記カウンタがゼロに達すると信号を発信するものである、参照カウンタと、前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへのアクセスを読み取ることが可能なように適合され、前記カウンタがゼロに達した際の前記参照カウンタの信号発信に
40

【0012】

好ましくは、この装置は、前記データの再送用に前記最適な間隔を決定するための分析機構と、前記データの再送用の前記最適な間隔に前記タイマ機構を調整するための調整機構と、をさらに有する。

【0013】

好ましくは、前記分析機構は、前記データの再送用に前記最適な間隔を決定するためにネットワーク監視データを使用するように動作可能である。

【0014】

50

好ましくは、前記ネットワークはストレージ・ネットワークを有する。

【0015】

好ましくは、前記ネットワークはインターネットを有する。

【0016】

第2の態様では、本発明は、受信機に伝送されるデータ項目を格納するためのバッファと、前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構であって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構と、第1の伝送用にデータにアクセスするため、および前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するための、バッファの第1のアクセス機構と、再送用に前記データにアクセスするため、および前記タイマ機構内の前記間隔に設定された少なくとも第2のタイムアウト・クロックを開始するための、前記タイマ機構による前記信号発信に応答する前記バッファの第2のアクセス機構であって、前記第2のアクセス機構は前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよう試行するものである、第2のアクセス機構と、を有する、複数のパスを有するネットワークを介してデータを伝送するための装置を含む、ストレージ・コントローラを提供する。

10

【0017】

第3の態様では、本発明は、受信機に伝送されるデータ項目を格納するためのバッファと、前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構であって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構と、第1の伝送用にデータにアクセスするため、および前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するための、バッファの第1のアクセス機構と、再送用に前記データにアクセスするため、および前記タイマ機構内の前記間隔に設定された少なくとも第2のタイムアウト・クロックを開始するための、前記タイマ機構による前記信号発信に応答する前記バッファの第2のアクセス機構であって、前記第2のアクセス機構は前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよう試行するものである、第2のアクセス機構と、を有する、複数のパスを有するネットワークを介してデータを伝送するための装置を含む、ネットワーク・アプリケーションを提供する。

20

30

【0018】

第2の態様のストレージ・コントローラおよび第3の態様のネットワーク・アプリケーションの好ましい機能は、第1の態様の装置のそれぞれの好ましい機能に対応する。

【0019】

第4の態様では、本発明は、受信機に伝送されるデータ項目を格納するためのバッファを提供するステップと、前記受信機宛てに前記データ項目を再送するための所定の最適な間隔を計るためのタイマ機構を提供するステップであって、前記間隔はネットワークからのいずれのエラー信号の受信に必要な間隔よりも短く、さらに前記タイマ機構は前記所定の最適な間隔の終わりに信号を発信するものである、タイマ機構を提供するステップと、第1の伝送用に第1のアクセス機構によってバッファ内のデータにアクセスするステップと、前記タイマ機構内の前記間隔に設定された第1のタイムアウト・クロックを開始するステップと、前記タイマ機構による前記信号発信に反応して第2のアクセス機構によってバッファ内のデータにアクセスするステップと、前記第2のアクセス機構によって前記データを再送するステップと、前記タイマ機構内の前記間隔に設定された第2のタイムアウト・クロックを開始するステップと、前記第1のアクセス機構によって使用されるパス要素を使用せずに前記データの伝送用のパスを選択するよう前記第2のアクセス機構によって試行するステップと、を有する、複数のパスを有するネットワークを介してデータを伝送するための方法を提供する。

40

【0020】

50

好ましくは、本方法は、前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへの参照のカウントを参照カウンタによって維持するステップであって、前記参照カウンタは、前記第1のアクセス機構および前記第2のアクセス機構によるそれぞれの参照時に増分され、各データ伝送の完了時に減分されるものであり、前記参照カウンタは前記カウンタがゼロに達すると信号を発信するものである、維持するステップと、メモリ・マネージャによって前記第1のアクセス機構および前記第2のアクセス機構による前記バッファへの読み取りアクセスを許可するステップと、前記カウンタがゼロに達した際の前記参照カウンタの信号発信にตอบสนองして前記バッファをフリー・バッファ・プールに戻すステップと、をさらに有する。

【0021】

10

好ましくは、本方法は、前記データの再送用に前記最適な間隔を決定するステップと、前記データの再送用の前記最適な間隔に前記タイマ機構を調整するステップと、をさらに有する。

【0022】

好ましくは、前記決定するステップは、前記データの再送用に前記最適な間隔を決定するためにネットワーク監視データを使用する。

【0023】

好ましくは、前記ネットワークはストレージ・ネットワークを有する。

【0024】

好ましくは、前記ネットワークはインターネットを有する。

20

【0025】

第5の態様では、本発明は、コンピュータ・システムにロードされ実行された場合に、第4の態様の方法のステップを実行するためのコンピュータ・プログラム・コードを有するコンピュータ・プログラムを提供する。

【0026】

したがって本システムは、好ましくは、同じアプリケーション・バッファからの複数のデータ伝送を同時に未処理とすることが可能なように構成される。

【0027】

ハードウェア制御ブロックなどのいくつかのリソースは、ハングされた伝送が持続している間、常時消費される。したがって好ましくは、こうしたリソースが独立した伝送パスに分離されるようにシステムが構成される。その結果、1つのパス上でハングされた伝送が、第2のパスへとデータを伝送する機能を妨げることはない。

30

【0028】

伝送がアクティブである（および認知されていない）間、タイマは応答の受信においていずれの異常な遅延も早期検出するように動作可能である。タイマがしきい値に達した場合は、代替パスの試行が有利である可能性があることを示す。したがって適切なことには、第1の伝送を待機またはブロックして完了する必要なしに、伝送は代替パスへと即時に再発行される。

【0029】

当業者であれば、起こり得る伝送障害に十分な即時の応答に失敗することと、他方で、システムが頻繁に応答しすぎる場合および再送のオーバーヘッドが増加しすぎた場合に、「偽肯定 (false positive)」に基づいた追加の再送によってシステムに負担をかけることとの間の均衡をとることによって、タイマ用に設定される間隔が特定のネットワークおよびそれに接続されたデバイスに好適な性能を最適化するように選択されるべきであることが明らかであろう。本発明の好ましい実施形態では、使用される間隔はこれを最適な値に設定するように調整することができる。最も好ましい実施形態では、タイマ間隔は、間隔がネットワーク性能に基づいた最適な値に調整されることを保証するようにネットワークの監視に基づいて設定される。他のまたは代替の実施形態では、最適な間隔を決定する際にサービス品質要件を考慮に入れることができる。

40

【0030】

50

さらに本発明の好ましい実施形態は、従来技術の仮想メモリ・システムが複数の伝送の進行を追跡することまたは伝送バッファへのアプリケーション・アクセスを監視することが不可能であるため、アプリケーション・バッファからの複数の同時伝送をサポートしていない、という事実によって発生する問題を軽減する。従来技術の従来システムでは、バッファが特定のデータ項目を含むように割り当てられ、ハードウェアI/Oデバイスに与えられたその実アドレスを有するようになると、ハードウェアI/Oデバイスは伝送持続期間に与えられたメモリ参照に依拠できなければならないため、たとえばバッファを含むページをディスクにスワップすることができないオペレーティング・システムを含んでいたとしても、バッファは実際には「フリーズ」され、ハードウェアI/Oデバイス以外のいずれのエンティティもアクセスすることはできない。本発明の好ましい実施形態は、ネットワーク・パスの途絶にもかかわらずシステム性能を向上させるために、タイムリーな再送を許可するように、ハードウェアI/Oデバイスおよびオペレーティング・システムの範囲外の複数のアクセスを制御するための機構を提供することによって、この制限を軽減する。

10

【0031】

本発明の好ましい実施形態は、データを余分にコピーする必要がないという他の利点を有する。代替パスが作動している場合、伝送パスのうちの1つの途絶または他の遅延にもかかわらず、システム内でのデータのタイムリーな可用性を維持することができる。

【0032】

次に、本発明の好ましい実施形態について添付の図面を参照しながら単なる例として説明する。

20

【発明を実施するための最良の形態】**【0033】**

本発明の現在で最も好ましい実施形態は、分散ストレージ・コントローラ内でインプリメントされるが、当業者であれば、クライアントおよびサーバ・コンピュータ・システムのネットワークを含むがこれらに限定されることのない他のネットワーク・システムでも、本発明が等しく有利に実施可能であることが明らかであろう。こうしたシステムは有線または無線とすることが可能であり、ローカル処理機能を有するデバイス、ならびに単純なI/Oデバイスなどのこうした機能のないデバイスを有することが可能である。

【0034】

分散ストレージ・コントローラの好ましい実施形態では、メモリ管理コンポーネントは32kBサイズまでのデータ部片用にI/Oバッファ記述を維持する。I/Oバッファは、SCSI書き込みオペレーションの一部として受信したホスト・システムからのデータなどの、何らかの外部ソースから受信したデータを含む。I/OバッファはPCIアドレスの分散・集合リストに変換することが可能であり、これらをファイバ・チャネル・アダプタへのデータ転送命令の一部として使用することができる。

30

【0035】

メモリ・マネージャは、分散・集合リストを構築するため、または伝送がアクティブな場合に、クライアントをブロックすることはない。バッファは、いずれかのI/Oハードウェア・デバイスによって使用されており、そのI/Oハードウェア・デバイスの伝送が確実な成功または確実な失敗のいずれかでまだ完了していない限り、「ピン固定状態(pinned)」を維持(すなわち、バッファに使用される仮想メモリが実のままであるように)しなければならない。

40

【0036】

図1は、アクセス機構要求に回答してメモリ・マネージャ(114)がアクセス可能なバッファ・メモリ(104)を有する、本発明の好ましい実施形態に従った装置(102)を示す図である。メモリ・マネージャ(114)はカウンタ(108)と通信している。カウンタ(108)は、アクセス機構(106、110)によるバッファ・メモリ(104)の参照回数のカウントを維持するものであり、カウントは各参照時に増分され、各アクセス機構のデータ伝送の完了時に減分される。さらにカウンタ(108)は、カウ

50

トがゼロに達すると信号を発するように適合される。アクセス機構(106、110)はそれぞれタイマ・クロック・インスタンス(116、120)に関連付けられる。タイマ・クロック・インスタンス(116、120)は、ネットワークから実エラー状態を戻すために必要となる「ラウンド・トリップ」時間よりも短い所定の最適な間隔を計るよう適合された、タイマ機構(122)内に提供される。さらにタイマ機構(122)は、各タイマ・クロック・インスタンス(116、120)によって設定された所定の最適な間隔の終わりに信号を発するようにも適合される。

【0037】

メモリ・マネージャ(114)は、書き込み動作中にバッファ・メモリをロックするように、アクセス機構(106、110)によるバッファ・メモリ(104)への読み取りアクセスを許可するように、およびカウンタがゼロに達したことをカウンタ(108)が信号発信する場合にバッファ・メモリをフリー・バッファ・プールに戻すように、適合される。

10

【0038】

複数の同時処理を追跡するために、メモリ・マネージャは同時アクセス機構の参照カウントを維持する。メモリ・バッファの同時アクセス機構はデータを読み取ることはできるが、データを書き込む(変更する)ことはできない。したがって、データを変更しようとするプロセスがなければ、データのソースを同時に読み取る複数の処理が同時にアクティブとなることができる。各プロセスは完了するとバッファへのその参照を「解放」し、参照カウントは減分される。最後のプロセスが完了すると参照カウントはゼロに達し、バッファはフリー・バッファのプールに入れられる。

20

【0039】

このようにして、伝送にデータのコンテンツの変更が含まれないことから、このメモリ・マネージャは複数のプロセスがバッファからデータを伝送できるようにする。

【0040】

メモリ・マネージャは、たとえオリジナルの要求または多くの以前の要求が依然としてアクティブであっても、第2または他のプロセスがバッファにアクセスするためにそれ専用のPCI分散・集合リストを構築するデータの同じバッファの伝送を開始できるようにする。

【0041】

本発明の好ましい実施形態では、各代替パスは、第1または他の以前のパスよりも成功する見込があるように選択される。これは、たとえば、別のインターフェース・アダプタまたは異なる物理ネットワークなどの完全に別のハードウェア要素のセットから代替パスを選択することによって保証される。さらにパスは、使用可能な帯域幅、コントローラまたはアプライアンスがアイドル状態であった期間中にping応答から推定された応答時間、およびその他などの、追加の基準を使用して選択することもできる。

30

【0042】

他の好ましい実施形態では、間隔チューナ(126)は、ネットワーク監視データ・アナライザ(128)によるデータの分析に基づいて、クロック(116、120)によって使用される間隔を調整するために使用される。

40

【0043】

好ましい実施形態のメモリ・マネージャは、伝送が失敗しそうな場合の迅速な再送の必要性と過度な再送オーバーヘッドを避ける必要性との均衡を取ることによって、性能を最適化するように適合されたタイマ機構と組み合わせられて、バッファからデータをタイムリーに再送することができる。同時に、好ましい実施形態のカウンタは、バッファ・メモリを割り振ること、「ピン固定状態」にすること、および他のアクセス機構による読み取りアクセスを妨げることなく解放することができる。

【0044】

好ましい実施形態の機構は、タイムアウト伝送が遅延し、永遠には失われない場合、伝送の受信機側でプロトコルをインプリメントし、データの重複受信に対処できるようにす

50

る必要がある。こうした方式は当分野でよく理解されており、TCP/IPなどのプロトコルですでにインプリメントされているため、ここではこれ以上詳しく説明しない。

【0045】

次に図2を見ると、本発明の好ましい実施形態に従った方法を示す単純な通信流れの例が示されている。

【0046】

通信流れには、タイマ、メモリ・マネージャ、アクセス機構1および2（複数のこうしたアクセス機構が存在可能であることを示すために他のアクセス機構3、4...nも示されている）、およびネットワークが含まれる。この流れは、1組の時間T1、T2、などで表して定義される。この通信流れは、第1のアクセス機構がタイムアウトし、第2のアクセス機構がデータを再送するために同じバッファへのアクセスを許可される、例示的なケースを示している。どちらの流れも（データの受信に成功するかまたは確実にエラーが戻されて）完了する。

10

【0047】

T1では、第1のアクセス機構がバッファにアクセスし、カウンタがそのアクセスを記録するために増分される。T2では、所定の最適なタイムアウト間隔をカウント・ダウンするためにクロックが開始される。バッファはT3でピン固定状態となり、T4でアクセス機構1がネットワークを介してデータを伝送する。T5では、アクセス機構1のタイムアウト間隔が満了し、アクセス機構2がバッファにアクセスするためにトリガされる。これはT6で発生し、アクセスを記録するためにカウンタが増分される。T7では、アクセス機構2のタイムアウトをカウント・ダウンするためにクロックが開始される。T8ではバッファがピン固定状態となり、T9ではアクセス機構2がデータを伝送する。T10ではアクセス機構1の伝送がネットワークによって認知され、アクセス機構1はバッファへのそのアクセスを解放する。アクセス機構1のアクセスの終わりを記録するためにカウンタが減分される。T12では、アクセス機構2の伝送がネットワークによって認知される。T13では、アクセス機構1はバッファへのアクセスを解放し、カウンタが減分される。T14ではカウンタがゼロになり、バッファはフリー・バッファ・プールへ戻される。当業者であれば、これが好ましい実施形態内で発生可能な基本的な流れのセットを説明するためだけに作成された非常に簡略化された例であることは明らかであろう。

20

【0048】

図3の論理流れは、バッファ・メモリ(104)と、書き込み動作中にバッファ・メモリ(104)をロックするように、バッファ・メモリ(104)への読み取りアクセスを許可するように、およびバッファ・メモリ(104)をフリー・バッファ・メモリのプールに戻すように適合されたメモリ・マネージャ(114)と、前記データの再送のために所定の最適な間隔を計るためのタイマ機構と、が提供され、前記間隔は前記ネットワークからのいずれのエラー信号を受け取るために必要な、および前記所定の最適な間隔の終わりに信号を発信するために必要な間隔よりも短い、好ましい実施形態に従ったデータの伝送のための方法を示すものである。

30

【0049】

ステップ(202)で、現在の好ましい実施形態に従ったランタイム論理プロセスに入る。ステップ(204)で、第1のアクセス機構が伝送用のバッファ・メモリ内のデータにアクセスする。ステップ(206)で、第1のアクセス機構によるバッファ・メモリ(104)へのアクセスを記録するためにカウンタが増分され、ステップ(208)で、データの再送のためにタイムアウト・クロックが所定の最適な間隔に設定され、その間隔は、ネットワークからのいずれのエラー信号の受信に必要なラウンド・トリップ時間よりも短い。

40

【0050】

ステップ(212)で、システムは、伝送が完了（伝送の正常な受信またはネットワークからの確実な失敗信号のいずれか）に達した旨のネットワークからの通知またはタイムアウト間隔の満了のいずれかを待つ。ステップ(214)で、タイムアウトのテストが実

50

行され、タイムアウト間隔の満了以前に完了を受信していない場合、論理はステップ(216)を開始し、クロックが停止される。ステップ(218)で、第2または他のアクセス機構が、ネットワークを介した再送のために同じバッファ・メモリ内の同じデータにアクセスする。ステップ(220)で、第2または他のアクセス機構がデータの伝送のために代替パスを選択し、この選択ステップは、可能であれば異なるパス要素の完全なセットを選択するように、およびそれが不可能な場合はできる限り多くの代替パス要素を選択するように、適合される。第1以外のアクセス機構は、それぞれ、以前のいずれかのアクセス機構が使用したパス要素は使用しないようにする。

【0051】

代替パスは、たとえば最低伝送コスト・パスまたは最高性能パス、またはおそらく最高のサービス品質保証を提供するパスを選択するように論理要素を最適化することによって、選択することもできる。こうしたサービス品質保証はネットワーク通信の分野では良く知られており、ここで説明する必要はない。現在の好ましい実施形態では、代替パスはネットワーク性能監視統計の分析に基づいて選択される。

10

【0052】

論理流れはステップ(206)に進み、ここで、第2または他のアクセス機構によるバッファ・メモリへのアクセスのステップを記録するためにカウントが増分され、以前と同様にクロックの設定などへと続行される。

【0053】

ステップ(214)でのいずれかの反復において、それぞれの伝送のタイムアウトが満了する前に完了(伝送の正常な受信またはネットワークからの確実な失敗信号のいずれか)があったことを示す、テストへの肯定応答があった場合、クロックはステップ(222)で停止され、カウントはステップ(224)で減分される。ステップ(226)で、カウントがゼロに達したか否かを判別するために他のテストが実行される。ゼロに達していない場合、論理プロセスのこの部分は次の完了によってトリガされるようにステップ(228)に戻る。いずれかの反復において、カウントがゼロに達したことが判別された場合、メモリ・マネージャ(114)はバッファ・メモリ(104)を解放する。

20

【0054】

次に図4を見ると、例示的な実施形態の他の好ましい改良が示されている。図4の論理プロセスはステップ(302)で始まり、ステップ(304)でネットワーク監視データを受け取る。こうしたデータは、当分野でよく知られた多数のネットワーク監視デバイス、システム、またはコンピュータ・プログラムのうちのいずれかによって提供可能である。このデータから、ラウンド・トリップ応答時間のステップ(306)で予測可能な評価が実行され、ステップ(308)で設定された任意のサービス品質パラメータが検査される。これらの入力に基づいて、ステップ(310)で最適な再送タイムアウト間隔が設定される。ステップ(312)で、プロセスのこの部分が終了する。当業者であれば明らかのように、このプロセスは反復可能であり、ネットワークからおよびサービス品質パラメータの任意の設定者から受信した最も新しいデータに従って、最適な間隔にリセットするような間隔で実行することができる。

30

【図面の簡単な説明】

40

【0055】

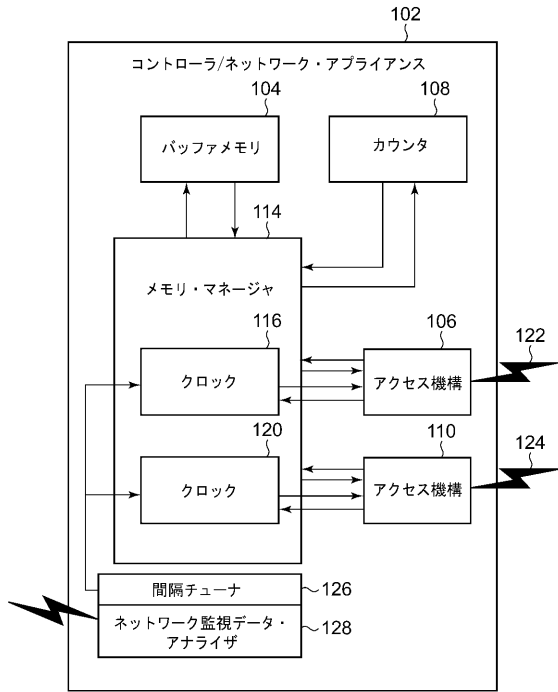
【図1】本発明の好ましい実施形態に従った装置を表すブロック図である。

【図2】本発明の好ましい実施形態に従った例示的通信流れを表す図である。

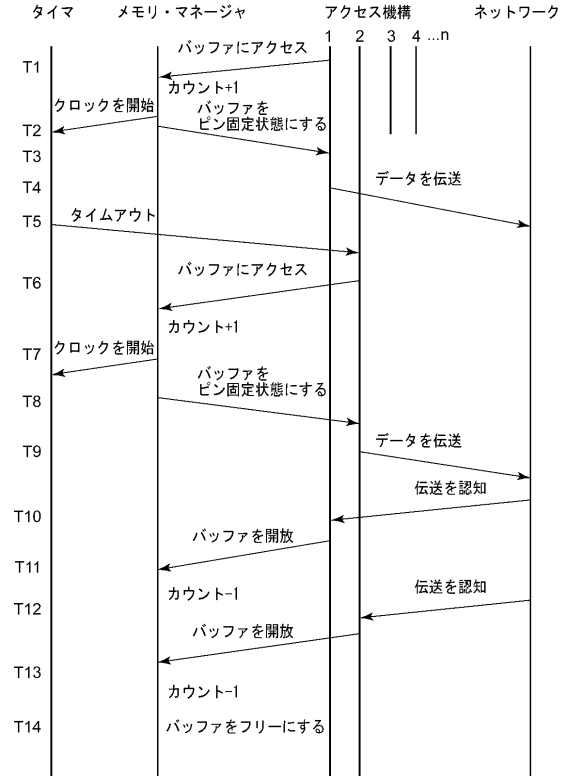
【図3】本発明の好ましい実施形態に従った論理流れを示す図である。

【図4】本発明の一実施形態の他の好ましい改良の論理流れを示す図である。

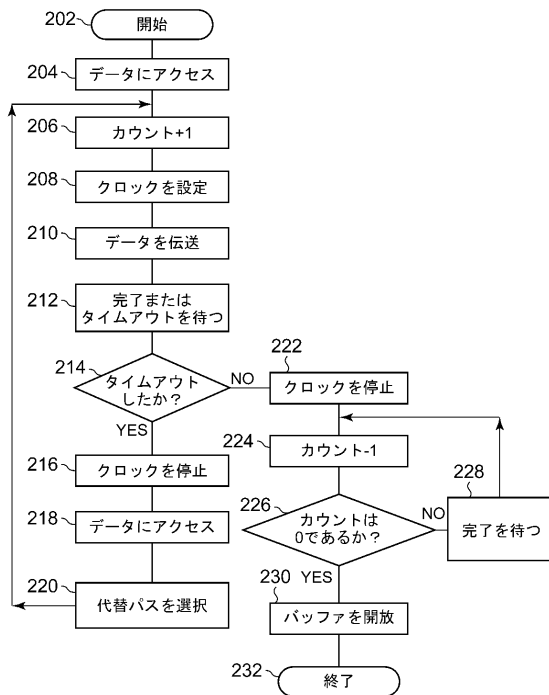
【 図 1 】



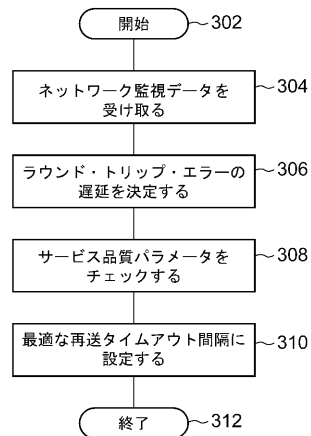
【 図 2 】



【 図 3 】



【 図 4 】



フロントページの続き

- (72)発明者 フエンテ、カルロス、フランシスコ
イギリス国ピーオー1 2ティーワイ ハンプシャー州ポーツマス ホワイト・ハート・ロード4
3
- (72)発明者 ジョーンズ、ロバート、マイケル
イギリス国ピーアール2 9ピーエフ ケント州ブロムリー ブロムリー・コモン102
- (72)発明者 パッシンガム、ウィリアム、ジョン
イギリス国エスオー30 2エヌエックス ハンプシャー州ボトリー プレコザ・ロード8
- (72)発明者 スケールズ、ウィリアム、ジェイムズ
イギリス国ピーオー16 8エイディー ハンプシャー州フェアラム ポーチェスター ポーチェ
スター・ロード5

審査官 安藤 一道

- (56)参考文献 国際公開第00/42746(WO, A1)
米国特許第5859959(US, A)
欧州特許出願公開第1089504(EP, A1)
国際公開第01/77830(WO, A1)

(58)調査した分野(Int.Cl., DB名)

H04L 29/08
H04L 12/56