

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4012545号
(P4012545)

(45) 発行日 平成19年11月21日(2007.11.21)

(24) 登録日 平成19年9月14日(2007.9.14)

(51) Int. Cl. F I
H O 4 L 12/56 (2006.01) H O 4 L 12/56 G

請求項の数 13 (全 52 頁)

(21) 出願番号	特願2004-533593 (P2004-533593)	(73) 特許権者	390009531
(86) (22) 出願日	平成15年8月5日(2003.8.5)		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2005-538588 (P2005-538588A)		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17年12月15日(2005.12.15)		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/GB2003/003418		
(87) 国際公開番号	W02004/023305	(74) 代理人	100086243
(87) 国際公開日	平成16年3月18日(2004.3.18)		弁理士 坂口 博
審査請求日	平成17年3月9日(2005.3.9)	(74) 代理人	100091568
(31) 優先権主張番号	10/235,691		弁理士 市位 嘉宏
(32) 優先日	平成14年9月5日(2002.9.5)	(74) 代理人	100108501
(33) 優先権主張国	米国 (US)		弁理士 上野 剛史

最終頁に続く

(54) 【発明の名称】 リモート・ダイレクト・メモリ・アクセス対応ネットワーク・インタフェース・コントローラの
スイッチオーバーとスイッチバックのサポート

(57) 【特許請求の範囲】

【請求項1】

分散データ処理システム中の各々のエンドノードの備える、リモート・ダイレクト・メモリ・アクセス対応ネットワーク・インタフェース・コントローラ(RNIC)である、主要RNICと代替RNIC間のスイッチオーバーを行う方法であって、

前記主要RNICと前記代替RNICが共通待ち行列ペアを共有するように、前記主要RNICと前記代替RNICに前記共通待ち行列ペアを作成するステップと

前記主要RNICからのチェックポイント・メッセージが代替RNICで所定期間内に受信されないことで、スイッチオーバー・イベントを検出するステップと、

前記スイッチオーバー・イベントの検出に回答して、前記共通待ち行列ペアの処理を前記主要RNICから前記代替RNICに切り替えるステップと 10
を備える方法。

【請求項2】

前記主要RNICにおいて、データ・セグメントを受信するステップと、

前記主要RNIC内の実効送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを更新するステップと、

実効送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを前記代替RNICに送信するステップと、

前記代替RNICにおいて、

前記主要 R N I C から、前記実効送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを受信するステップと、

チェックポイント・メッセージに含まれる命令コードから、前記チェックポイント・メッセージのタイプを、前記実効送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージと識別するステップと、

前記代替 R N I C において、保持している、実効送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、前記実効送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージの内容で更新するステップと、

前記代替 R N I C において、保持している、前記実効送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、コミット済み送信待ち行列スイッチオーバー・コンテキストのローカル・コピーにコピーするステップと、

コミット済み送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを前記主要 R N I C に送信するステップと、

前記主要 R N I C において、

前記代替 R N I C から、前記コミット済み送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを受信するステップと、

前記チェックポイント・メッセージに含まれる命令コードから、前記チェックポイント・メッセージのタイプを、前記コミット済み送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージと識別するステップと、

前記主要 R N I C において、保持している、コミット済み送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、前記コミット済み送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージの内容で更新するステップと、

送信待ち行列から前記データ・セグメントを送信するステップと

を含む請求項 1 記載の方法。

【請求項 3】

前記共通待ち行列ペアを作成するステップは、前記主要 R N I C と前記代替 R N I C がその待ち行列の範囲内で待ち行列を共有するように、前記主要 R N I C と前記代替 R N I C にある待ち行列の範囲を割り当てるステップを含む請求項 1 に記載の方法。

【請求項 4】

前記共通待ち行列ペアを作成するステップは、前記主要 R N I C と前記代替 R N I C に、ある範囲のメモリ変換保護テーブル・エントリを割り当てるステップを含む請求項 1 ないし 3 のいずれかに記載の方法。

【請求項 5】

前記主要 R N I C と前記代替 R N I C が共通完了待ち行列を共有するように、前記主要 R N I C と前記代替 R N I C に、前記共通完了待ち行列を作成するステップと、

スイッチオーバー・イベントが検出された場合は、前記共通完了待ち行列の操作を前記代替 R N I C に切り替えるステップと

をさらに備える請求項 1 ないし 4 のいずれかに記載の方法。

【請求項 6】

前記主要 R N I C から前記代替 R N I C に前記共通待ち行列ペアの扱いを切り替えるステップは、

前記主要 R N I C のアドレスを前記代替 R N I C のアドレス・テーブルに追加するステップと、

前記主要 R N I C が、前記主要 R N I C と前記代替 R N I C とに結合されたスイッチからアクセスできないようにするステップと、

前記主要 R N I C のアドレスを、前記スイッチ内で前記代替 R N I C のアドレスとして認識できるようにするステップと

を含む請求項 1 ないし 5 のいずれかに記載の方法。

【請求項 7】

前記主要 R N I C から前記代替 R N I C に前記共通待ち行列ペアの扱いを切り替えるス

10

20

30

40

50

テップは、

前記代替 R N I C が前記主要 R N I C として認識されるように、共通完了待ち行列コンテキストについての前記主要 R N I C と前記代替 R N I C の識別コンテキスト情報を変更するステップと、

前記代替 R N I C が前記主要 R N I C として認識されるように、共通待ち行列ペアコンテキストについての前記主要 R N I C と前記代替 R N I C の識別コンテキスト情報を変更するステップとを含む請求項 1 ないし 6 のいずれかに記載の方法。

【請求項 8】

前記主要 R N I C と前記代替 R N I C は、前記共通待ち行列のペアと共通完了待ち行列とのスイッチオーバー状態を識別する状態情報を、前記主要 R N I C および前記代替 R N I C 内に保持するステップを含む請求項 7 に記載の方法。

10

【請求項 9】

前記状態情報を保持するステップは、前記主要 R N I C と、前記代替 R N I C と、のスイッチオーバー状態、前記主要 R N I C の識別子、前記主要 R N I C のポート識別子、前記代替 R N I C の識別子、および前記代替 R N I C のポート識別子の少なくとも 1 つを含む、前記主要 R N I C および前記代替 R N I C 識別コンテキスト・データ構造を保持するステップを含む請求項 7 に記載の方法。

【請求項 10】

前記状態情報を保持するステップは、前記主要 R N I C と前記代替 R N I C の各々の前記共通完了待ち行列についてのエントリを有する共通完了待ち行列コンテキスト・テーブルを保持するステップを含む請求項 7 ないし 9 に記載の方法。

20

【請求項 11】

コンピュータ・システムにロードされ、実行されると、請求項 1 ないし 10 のいずれかに記載の方法のステップを前記コンピュータに行わせるコンピュータ・プログラム・コード要素を備えるコンピュータ・プログラム。

【請求項 12】

分散データ処理システム中の各々のエンドノードの備える、前記主要 R N I C と前記代替 R N I C の間のスイッチオーバーを行う装置であって、

前記主要 R N I C と前記代替 R N I C が共通待ち行列のペアを共有するように、前記主要 R N I C と前記代替 R N I C に前記共通待ち行列のペアを作成する手段と、

30

主要 R N I C からのチェックポイント・メッセージが代替 R N I C で所定期間内に受信されないことで、スイッチオーバー・イベントを検出する手段と、

前記スイッチオーバー・イベントの検出にตอบสนองして、前記共通待ち行列のペアの処理を前記主要 R N I C から前記代替 R N I C に切り替える手段と
を備える装置。

【請求項 13】

前記主要 R N I C において、
データ・セグメントを受信する手段と、

前記主要 R N I C 内の実効送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを更新する手段と、

40

実効送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを前記代替 R N I C に送信する手段と

前記代替 R N I C において、
前記主要 R N I C から、前記実効送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを受信する手段と、
チェックポイント・メッセージに含まれる命令コードから、前記チェックポイント・メッセージのタイプを、実効送信待ち行列スイッチオーバー送信コンテキスト更新送信チェックポイント・メッセージと識別する手段と、

前記代替 R N I C において、保持している、実効送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、前記実効送信待ち行列スイッチオーバー・コンテキスト

50

更新送信チェックポイント・メッセージの内容で更新する手段と、

前記代替RNICにおいて、保持している、前記実効送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、コミット済み送信待ち行列スイッチオーバー・コンテキストのローカル・コピーにコピーする手段と、

コミット済み送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージを前記主要RNICに送信する手段と、

前記主要RNICにおいて、

前記代替RNICから、前記コミット済み送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージを受信する手段と、

前記チェックポイント・メッセージに含まれる命令コードから、前記チェックポイント・メッセージのタイプを、前記コミット済み送信待ち行列スイッチオーバー・コンテキスト更新送信チェックポイント・メッセージと識別する手段と、

前記主要RNICにおいて、保持している、コミット済み送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、前記コミット済み送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージの内容で更新する手段と、

送信待ち行列から前記データ・セグメントを送信する手段と
を含む請求項12記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般には、ホスト・コンピュータと入出力（I/O）デバイス間の通信プロトコルに関する。より詳細には、本発明は、リモート・ダイレクト・メモリ・アクセス（RDMA）対応のネットワーク・インタフェース・コントローラ（NIC）が、主要なRDMA対応NIC（RNIC）と代替RNICで構成される冗長な構成をサポートできる方法を提供する。

【背景技術】

【0002】

インターネット・プロトコル（IP）ネットワークでは、ソフトウェアが、入出力デバイス、汎用コンピュータ（ホスト）、および特殊目的コンピュータとの間の通信に使用することができるメッセージ・パッシング機構を提供する。このメッセージ・パッシング機構は、トランスポート・プロトコル、上位プロトコル、およびアプリケーション・プログラミング・インタフェースで構成される。現在IPネットワークで使用される主要な標準的トランスポート・プロトコルは、伝送制御プロトコル（TCP）とユーザ・データグラム・プロトコル（UDP）である。TCPは、信頼性のあるサービスを提供し、UDPは、信頼性のないサービスを提供する。将来は、信頼性のあるサービスを提供するためにストリーム制御伝送プロトコル（SCTP）も使用されるものと予想される。デバイスまたはコンピュータで実行されるプロセスは、Socket、iSCSI、およびダイレクト・アクセス・ファイル・システム（DAFS）などの上位プロトコルを通じてIPネットワークにアクセスする。不都合なことに、TCP/IPソフトウェアは、相当量のプロセッサ・リソースとメモリ・リソースを消費する。この問題は、次の文献で広く扱われている（J. ケイ（J. Kay）、J. パスクアレ（J. Pasquale）「TCP/IPにおけるオーバーヘッドのプロファイリングと低減（Profiling and reducing processing overheads in TCP/IP）」1996年12月、IEEE/ACM Networking 紀要、Volume 4、No. 6、817～828頁、および、D. D. クラーク（D.D. Clark）、V. ヤコブソン（V. Jacobson）、J. ロムキー（J. Romkey）、H. サルウェン（H. Salwen）「TCP処理のオーバーヘッドの分析（An analysis of TCP processing overhead）」IEEE 通信マガジン（IEEE Communications Magazine）1989年6月、Volume 27、Issue 6 23～29頁、を参照）。今後も、アプリケーションによるネットワークの使用の増加、ネットワーク・セキュリティ・プロトコルの使用、基礎ファブリックの帯域幅がマイクロプロセッサとメモリの帯域幅を上回る速度で増大しつつある

10

20

30

40

50

こと等のいくつかの理由から、ネットワーク・スタックは、引き続き過度なリソースを消費していくと予想される。この問題に対処するために、業界は、ネットワーク・スタックの処理を、RDMA対応NIC(RNIC)に移している。

【0003】

業界でとられているオフロードの手法には2つある。第1の手法は、プロトコルをそれ以上追加せずに、既存のTCP/IPのネットワーク・スタックを使用するものである。この手法では、TCP/IPをハードウェアにオフロードすることができるが、不都合な点として、サイド・コピーを受け取る必要性はなくなる。上記の論文で言及されるように、コピーは、CPUの利用率に対する最大の要因の1つである。コピーの必要性をなくすために、業界は、TCPプロトコルおよびSCTPプロトコルに、フレーム化、直接データ配置(DDP)、およびリモート・ダイレクト・メモリ・アクセス(RDMA)を追加することからなる第2の手法を推進している。これら2つの手法をサポートするのに必要なRDMA対応NIC(RNIC)は同様であり、その大きな相違点は、第2の手法ではハードウェアが追加的なプロトコルをサポートしなければならないことである。

10

【0004】

RNICは、ノード間で通信するためにソケット、iSCSI、およびDAFSによって使用できるメッセージ・パッシング機構を提供する。ホスト・コンピュータあるいはデバイスで実行されるプロセスは、RNICの送受信作業待ち行列に送信/受信メッセージを入れることにより、IPネットワークにアクセスする。これらのプロセスは、「コンシューマ」とも呼ばれる。

20

【0005】

送信/受信の作業待ち行列(WQ)は、待ち行列のペア(QP)としてコンシューマに割り当てられる。メッセージは、従来のTCP、RDMA TCP、UDP、またはSCTPなどいくつかの異なるトランスポート・タイプで送信することができる。コンシューマは、RNICの送受信作業完了(WC)待ち行列を通じて、完了待ち行列(CQ)から、それらメッセージの結果を取得する。送信元のRNICは、発信メッセージを分割し、それらを宛先に送信する役割を担う。宛先RNICは、着信メッセージを再度組み立て、そのメッセージを、宛先のコンシューマによって指定されたメモリ空間に入れる役割を担う。これらのコンシューマは、RNIC verbを使用して、RNICによってサポートされる機能にアクセスする。verbを解釈し、RNICに直接アクセスするソフトウェアは、RNICインタフェース(RI)と称される。

30

【0006】

現在は、ホストCPU内のソフトウェアが、トランスポート層(例えばTCP)とネットワーク層(例えばIP)の処理の大半を行っている。現在、NICは通例、リンク層(例えばイーサネット(R))の処理と、恐らくは少量のトランスポート層またはネットワーク層のオフロード(例えばチェックサムのオフロード)を行っている。現在は、ホスト・ソフトウェアが、TCP/IP接続に関連するすべての状態情報をホスト・ローカル・メモリに保持している。すべての状態情報をホスト・ローカル・メモリに保持することにより、ホスト・ソフトウェアが、主要なNICと代替NICの間のスイッチオーバーおよびスイッチバックをサポートすることが可能になる。すなわち、主要NICに障害が発生した場合、ホスト・ソフトウェアは、すべての接続を代替NICに移し、通信処理を続ける。

40

【非特許文献1】J. ケイ(J. Kay)、J. パスクアレ(J. Pasquale)「TCP/IPにおけるオーバーヘッドのプロファイリングと低減(Profiling and reducing processing overheads in TCP/IP)」1996年12月、IEEE/ACM Networking 紀要、Vol 4、No. 6、817~828頁

【非特許文献2】D. D. クラーク(D. D. Clark)、V. ヤコブソン(V. Jacobson)、J. ロムキー(J. Romkey)、H. サルウェン(H. Salwen)「TCP処理のオーバーヘッドの分析(An analysis of TCP processing overhead)」IEEE通信マガジン(IEEE Communications Magazine)1989年6月、Volume 27、Issue. 6 23~

50

29頁

【発明の開示】

【発明が解決しようとする課題】

【0007】

RDMA対応NICは、他の汎用コンピュータまたはI/Oデバイスあるいはその両方と通信するためのより高性能のインタフェースを提供する。RNICは、トランスポート層(TCP等)とネットワーク層(例えばIP)をRNICにオフロードする。これらの層をRNICに移すと、ホスト・ソフトウェアは、もはや現在の機構を使用してスイッチオーバーとスイッチバックをサポートすることができなくなる。したがって、RNICが信頼性のあるトランスポート層(例えばTCP)の接続のスイッチオーバーとスイッチバックをサポートできるようにし、意図された、または意図されないRNICの停止が生じてても通信が継続できるようにする単純な機構が必要とされる。

10

【課題を解決するための手段】

【0008】

したがって、本発明は、第1の態様で、データ処理システム中のリモート・ダイレクト・メモリ・アクセス対応の主要ネットワーク・インタフェース・コントローラ(RNIC)と代替のRNIC間でスイッチオーバーを行う方法を提供し、この方法は、主要RNICと代替RNICが共通の待ち行列ペアを共有するように、主要RNICと代替RNICに共通の待ち行列ペアを作成するステップと、スイッチオーバー・イベントを検出するステップと、スイッチオーバー・イベントの検出にตอบสนองして、待ち行列ペアの扱いを主要RNICから代替RNICに切り替えるステップとを備える。

20

【0009】

好ましくは、共通の待ち行列ペアを作成するステップは、主要RNICと代替RNICがその待ち行列の範囲内で待ち行列を共有するように、主要RNICと代替RNICにある範囲の待ち行列を割り当てるステップを含む。

【0010】

好ましくは、待ち行列の範囲は、ある範囲の待ち行列ペアとある範囲の完了待ち行列を含む。

【0011】

好ましくは、共通の待ち行列ペアを作成するステップは、主要RNICと代替RNICに、ある範囲のメモリ変換保護テーブル・エントリを割り当てるステップを含む。

30

【0012】

この方法は、好ましくは、主要RNICと代替RNICが共通の完了待ち行列を共有するように、主要RNICと代替RNICに対して共通の完了待ち行列を作成するステップと、スイッチオーバー・イベントが検出された場合は、共通の完了待ち行列の操作を代替RNICに切り替えるステップとを備える。

【0013】

好ましくは、主要RNICと代替RNICは、ファブリックおよび共通システム・メモリの1つを介して、相互にチェックポイント・メッセージを送信する。

【0014】

好ましくは、ファブリックは、ローカル・エリア・ネットワーク、ワイド・エリア・ネットワーク、メモリ・マップI/O拡張ネットワーク、システム・エリア・ネットワークの1つである。

40

【0015】

好ましくは、チェックポイント・メッセージは、命令コード・フィールド、長さフィールド、および有効性検証フィールドを含む。

【0016】

好ましくは、命令コードは、チェックポイント・メッセージのタイプを、aliveメッセージ、実効送信待ち行列スイッチオーバー送信コンテキスト更新、コミット済み送信待ち行列スイッチオーバー送信コンテキスト更新、実効送信待ち行列スイッチオーバー肯

50

定応答コンテキスト更新、コミット済み送信待ち行列スイッチオーバー肯定応答コンテキスト更新、実効受信待ち行列スイッチオーバー受信コンテキスト更新、コミット済み受信待ち行列スイッチオーバー受信コンテキスト更新、実効受信待ち行列スイッチオーバー肯定応答コンテキスト更新、コミット済み受信待ち行列スイッチオーバー肯定応答コンテキスト更新、実効完了待ち行列スイッチオーバー・コンテキスト更新、および、コミット済み完了待ち行列スイッチオーバー・コンテキスト更新、の1つと識別する。

【0017】

好ましくは、スイッチオーバー・イベントは、主要RNICからのチェックポイント・メッセージが代替RNICで所定期間内に受信されないことである。

【0018】

好ましくは、主要RNICから代替RNICへの待ち行列ペアの扱いの切り替えは、主要RNICのアドレスを代替RNICのアドレス・テーブルに追加するステップと、主要RNICが、主要RNICと代替RNICとに結合されたスイッチからアクセスできないようにするステップと、主要RNICのアドレスを、スイッチ内で代替RNICのアドレスとして認識可能にするステップとを含む。

【0019】

好ましくは、主要RNICから代替RNICへの待ち行列ペアの扱いの切り替えは、代替RNICが主要RNICとして認識されるように、完了待ち行列コンテキストについての主要RNICと代替RNICの識別コンテキスト情報を変更するステップと、代替RNICが主要RNICとして認識されるように、待ち行列ペアコンテキストについての主要RNICと代替RNICの識別コンテキスト情報を変更するステップとを含む。

【0020】

この方法は、好ましくは、さらに、待ち行列のスイッチオーバー状態を識別する状態情報を主要RNICと代替RNIC内に保持するステップを備える。

【0021】

好ましくは、状態情報を保持するステップは、RNICのスイッチオーバー状態、主要RNICの識別子、主要RNICのポート識別子、代替RNICの識別子、および代替RNICのポート識別子の少なくとも1つを含む、主要RNICおよび代替RNIC識別コンテキスト・データ構造を保持するステップを含む。

【0022】

好ましくは、状態情報を保持するステップは、主要RNICと代替RNICの各共通の完了待ち行列についてのエントリを有する完了待ち行列コンテキスト・テーブルを保持するステップを含む。

【0023】

本発明は、好ましくは、主要RNICと代替RNICを有するデータ処理システム中でデータ・セグメントを送信する方法を提供し、この方法は、データ・セグメントを受信するステップと、主要RNIC内で実効送信待ち行列のスイッチオーバー・コンテキスト情報を更新するステップと、実効送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージを代替RNICに送信するステップと、代替RNICから、コミット済み送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージを受信するステップと、コミット済み送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを更新するステップと、送信待ち行列からデータ・セグメントを送信するステップとを備える。

【0024】

好ましくは、代替RNICは、ローカルの実効送信待ち行列スイッチオーバー・コンテキストを、実効送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージの内容で更新し、実効送信待ち行列スイッチオーバー・コンテキストを、コミット済み送信待ち行列スイッチオーバー・コンテキストにコピーする。

【0025】

10

20

30

40

50

第2の態様で、本発明は、コンピュータ・システムにロードされ、実行されると、第1の態様による方法のステップを前記コンピュータに行わせるコンピュータ・プログラム・コード要素を備えるコンピュータ・プログラムを提供する。第2の態様の好ましい特徴は、第1の態様の好ましい特徴の方法ステップに対応するステップを行うプログラム・コード要素を備える。

【0026】

第3の態様で、本発明は、主要RNICと代替RNICを有するデータ処理システム中でデータ・セグメントを送信する装置を提供し、この装置は、データ・セグメントを受信する手段と、主要RNIC内の実効送信待ち行列スイッチオーバー・コンテキスト情報を更新する手段と、実効送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージを代替RNICに送信する手段と、代替RNICから、コミット済み送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージを受信する手段と、コミット済み送信待ち行列スイッチオーバー・コンテキストのローカル・コピーを、コミット済み送信待ち行列スイッチオーバー・コンテキスト更新チェックポイント・メッセージの内容で更新する手段と、送信待ち行列からデータ・セグメントを送信する手段と、を備える。

10

【0027】

したがって、本発明の好ましい実施形態は、RNICのスイッチオーバーとスイッチバックをサポートする方法、コンピュータ・プログラム製品、および分散データ処理システムを提供し、この分散データ処理システムは、エンドノード、スイッチ、ルータ、およびそれらの構成要素を相互に接続するリンクを備え、エンドノードは、送信および受信の待ち行列ペアを使用してメッセージを送受信する。エンドノードは、メッセージをセグメントに分割し、リンクを通じてそのセグメントを送信することが好ましい。スイッチとルータは、エンドノードを相互に接続し、セグメントを適切なエンドノードにルーティングすることができる。エンドノードは、宛先でセグメントをメッセージに組み立て直すことができる。

20

【0028】

本発明の好ましい実施形態は、RNIC(RDMA対応RNIC)のスイッチオーバーとスイッチバックをサポートする機構を提供し、主要RNICで、意図される停止、または意図されない停止が生じた際に、提供されるこの好ましい機構を使用して、すべての継続中の接続が代替RNICに切り替えられ、代替RNICが、通信処理を継続する。また、本発明の好ましい実施形態で提供される機構を使用して、接続を主要RNICに再度切り替えることもできる。

30

【0029】

次いで、添付図面を参照して、本発明の好ましい実施形態を単なる例として説明する。

【発明を実施するための最良の形態】

【0030】

本発明の好ましい実施形態は、エンドノード、スイッチ、ルータ、およびそれらの構成要素を相互に接続するリンクを有する分散コンピューティング・システムを提供する。エンドノードは、インターネット・プロトコル・スイート・オフロード・エンジン、または従来のホスト・ソフトウェアに基づくインターネット・プロトコル・スイートである。各エンドノードは、送信および受信の待ち行列ペアを使用してメッセージを送受信する。エンドノードは、メッセージをフレームに分割し、リンクを通じてそのフレームを送信する。スイッチとルータは、エンドノードを相互に接続し、フレームを適切なエンドノードにルーティングする。エンドノードは、宛先でフレームを再度組み立ててメッセージにする。

40

【0031】

次いで図面、具体的には図1を参照すると、本発明の好ましい実施形態による分散コンピュータ・システムの図が示される。図1に表す分散コンピュータ・システムは、IPネットワーク100などのインターネット・プロトコル・ネットワーク(IPネット)の形をとり

50

、このシステムは単に説明の目的で提供され、以下に記載される本発明の実施形態は、多数の他のタイプおよび構成のコンピュータ・システムで実施することができる。例えば、本発明の好ましい実施形態を実施するコンピュータ・システムは、1つのプロセッサと数個の入出力（I/O）アダプタを備える小型サーバから、数百または数千個のプロセッサと数千個のI/Oアダプタを備える超並列スーパーコンピュータ・システムに及ぶことができる。さらに、本発明の好ましい実施形態は、インターネットまたはイントラネットで接続された遠隔コンピュータ・システムからなるインフラストラクチャで実施することができる。

【0032】

IPネット100は、分散コンピュータ・システム内のノードを相互に接続する広帯域幅、低待ち時間のネットワークである。ノードは、ネットワークの1つまたは複数のリンクに接続され、ネットワーク内でメッセージの発信元または宛先あるいはその両方を形成する任意のコンポーネントである。図の例では、IPネット100は、ホスト・プロセッサ・ノード102、ホスト・プロセッサ・ノード104、およびRAID（redundant array independent disk）サブシステム・ノード106の形態のノードを含む。図1に示すノードは、説明のみを目的とするものであり、IPネット100は、いくつの数のどのようなタイプの独立したプロセッサ・ノード、ストレージ・ノード、および特殊目的処理ノードにも接続することができる。どのノードもエンドノードとして機能することができ、エンドノードとは、ここでは、IPネット100中でメッセージを発信するか、最終的にメッセージを消費するデバイスと定義される。

【0033】

本発明の一実施形態では、分散コンピュータ・システム中のエラー処理機構が存在し、このエラー処理機構により、IPネット100などの分散コンピューティング・システム中のエンドノード間のTCPまたはSCTP通信が可能になる。

【0034】

本発明で使用されるメッセージとは、アプリケーションによって定義されるデータ交換の単位であり、これは協力プロセス間の通信の基本単位である。フレームは、インターネット・プロトコル・スイート・ヘッダまたはトレーラ（trailer）あるいはその両方によってカプセル化された1単位のデータである。ヘッダは、一般に、IPネット100を通じてフレームを誘導するための制御情報およびルーティング情報を提供する。トレーラは、一般に、内容が破損した状態でフレームが配信されないようにする制御データと巡回冗長検査（CRC）データを含んでいる。

【0035】

分散コンピュータ・システム内で、IPネット100は、ストレージ、プロセス間通信（IPC）、ファイル・アクセス、ソケットなど、各種形態のトラフィックをサポートする通信と管理のインフラストラクチャを含む。図1に示すIPネット100は、交換通信ファブリック116を含み、通信ファブリック116は、セキュリティが確保された、遠隔から管理される環境で、多くのデバイスが、広帯域幅、低待ち時間で同時にデータを転送することを可能にする。エンドノードは、複数のポートを通じて通信することができ、また、IPネット・ファブリックを通る複数の経路を利用することができる。図1に示す複数のポートと、IPネットを通る経路は、耐障害性と、高帯域幅のデータ転送のために用いることができる。管理と保守のための交換通信ファブリック116のコンポーネントへのアクセスは、コンソール110を通じて行うことができる。

【0036】

図1のIPネット100は、スイッチ112、スイッチ114、およびルータ117を含む。スイッチは、複数のリンクをともに接続し、レイヤ2の宛先アドレス・フィールドを使用して、あるリンクから別のリンクへのフレームのルーティングを可能にするデバイスである。イーサネット（R）がリンクとして使用される場合、宛先フィールドは、メディア・アクセス制御（MAC）アドレスと称される。ルータは、レイヤ3の宛先アドレス・フィールドに基づいてフレームをルーティングするデバイスである。インターネット・

10

20

30

40

50

プロトコル (I P) がレイヤ 3 プロトコルとして使用される場合は、宛先アドレス・フィールドは I P アドレスになる。

【 0 0 3 7 】

一実施形態では、リンクは、エンドノード、スイッチ、ルータなど任意の 2 つのネット・ファブリック要素間の全二重チャネルである。適切なリンクの例には、これらに限定しないが、銅線ケーブル、光ケーブル、およびバックプレーンとプリント回路基盤のプリント回路銅線トレースが含まれる。

【 0 0 3 8 】

信頼性のあるサービス・タイプ (T C P および S C T P) の場合、ホスト・プロセッサ・エンドノードや I / O アダプタ・エンドノードなどのエンドノードは、要求フレームを生成し、肯定応答フレームを返す。スイッチとルータは、送信元から宛先にフレームを転送する。

10

【 0 0 3 9 】

図 1 に示す I P ネット 1 0 0 では、ホスト・プロセッサ・ノード 1 0 2、ホスト・プロセッサ・ノード 1 0 4、および R A I D サブシステム・ノード 1 0 6 は、I P ネット 1 0 0 とのインタフェースをとる少なくとも 1 つの I P S O E を含む。一実施形態では、各 I P S O E は、I P ネット 1 0 0 で送信されるソース・フレームまたはシンク・フレームに対して十分な詳細で I P S O E を実装するエンドポイントである。ホスト・プロセッサ・ノード 1 0 2 は、ホスト I P S O E 1 1 8 および I P S O E 1 2 0 の形態の I P S O E を含む。ホスト・プロセッサ・ノード 1 0 4 は、I P S O E 1 2 2 と I P S O E 1 2 4 を含む。ホスト・プロセッサ・ノード 1 0 2 は、バス・システム 1 3 4 によって相互に接続された中央演算処理装置 1 2 6 ~ 1 3 0 とメモリ 1 3 2 も含む。ホスト・プロセッサ・ノード 1 0 4 も同様に、バス・システム 1 4 4 によって相互に接続された中央演算処理装置 1 3 6 ~ 1 4 0 とメモリ 1 4 2 を含む。

20

【 0 0 4 0 】

I P S O E 1 1 8 は、スイッチ 1 1 2 との接続を提供し、一方、I P S O E 1 2 4 はスイッチ 1 1 4 との接続を提供し、I P S O E 1 2 0 および 1 2 2 は、スイッチ 1 1 2 および 1 1 4 との接続を提供する。

【 0 0 4 1 】

一実施形態では、I P スイート・オフロード・エンジンは、ハードウェア、またはハードウェアとオフロード・マイクロプロセッサの組み合わせとして実施される。この実施では、I P スイートの処理が I P S O E にオフロードされる。この実施では、通信プロトコルに伴う従来のオーバーヘッドなしに、交換ネットワークを通じて複数の同時の通信も許可される。一実施形態では、図 1 の I P S O E と I P ネット 1 0 0 が、オペレーティング・システムのカーネル・プロセスを関与させずに、分散コンピュータ・システムのコンシューマにプロセッサ・コピーのないデータ転送を提供し、ハードウェアを使用して信頼性と耐障害性のある通信を提供する。

30

【 0 0 4 2 】

図 1 に示すように、ルータ 1 1 7 は、他のホストまたは他のルータとのワイド・エリア・ネットワーク (W A N) またはローカル・エリア・ネットワーク (L A N) 接続、あるいはその両方の接続に結合される。

40

【 0 0 4 3 】

この例では、図 1 の R A I D サブシステム・ノード 1 0 6 は、プロセッサ 1 6 8、メモリ 1 7 0、I P スイート・オフロード・エンジン (I P S O E) 1 7 2、および冗長な、またはストライピングされた、あるいはその両方の複数の記憶ディスク・ユニット 1 7 4 を含む。

【 0 0 4 4 】

I P ネット 1 0 0 は、ストレージ、プロセッサ間通信、ファイル・アクセス、およびソケットのためのデータ通信に対応する。I P ネット 1 0 0 は、広帯域幅で拡張可能性があり、極めて低待ち時間の通信をサポートする。ユーザ・クライアントは、オペレーティン

50

グ・システムのカーネル・プロセスを回避し、効率的なメッセージ・パッシング・プロトコルを可能にするIPSOEなどのネット通信コンポーネントに直接アクセスすることができる。IPネット100は、現在のコンピューティング・モデルに適合されており、新しい形態のストレージ、クラスタ、および一般的なネットワーキング通信のための基礎単位である。さらに、図1のIPネット100では、ストレージ・ノード間で通信する、または、ストレージ・ノードが分散コンピュータ・システム中のプロセッサ・ノードのいずれかまたはすべてと通信することができる。IPネットワーク100にストレージが付加されるので、ストレージ・ノードは、IPネットワーク100のどのホスト・プロセッサ・ノードとも、実質的に同じ通信能力を有する。

【0045】

一実施形態では、図1に示すIPネット100は、チャンネル・セマンティクスとメモリ・セマンティクスをサポートする。チャンネル・セマンティクスは、時に、送受信動作またはプッシュ通信動作とも称される。チャンネル・セマンティクスは、送信元デバイスがデータをプッシュし、宛先デバイスがそのデータの最終的な宛先を判定する従来のI/Oチャンネルで用いられる通信タイプである。チャンネル・セマンティクスでは、送信元プロセスから送信されるフレームが宛先プロセスの通信ポートを指定するが、宛先プロセスのメモリ空間のどこにそのフレームを書き込むかは指定しない。したがって、チャンネル・セマンティクスでは、宛先プロセスが、送信されたデータを入れる場所をあらかじめ割り当てる。

【0046】

メモリ・セマンティクスでは、送信元のプロセスが、リモート・ノードの宛先プロセスの仮想アドレス空間を直接読み出すか、または仮想アドレス空間に書き込む。リモートの宛先プロセスは、データのためのバッファの位置を伝えるだけでよく、どのデータの転送にも関与する必要がない。したがって、メモリ・セマンティクスでは、送信元プロセスは、宛先プロセスの宛先バッファのメモリ・アドレスを含んだデータ・フレームを送信する。メモリ・セマンティクスでは、宛先プロセスは、送信元プロセスが自身のメモリにアクセスする許可を事前に与える。

【0047】

一般に、チャンネル・セマンティクスとメモリ・セマンティクスはどちらも、記憶、クラスタ、および一般的なネットワーキング通信に必要とされる。典型的な記憶動作では、チャンネル・セマンティクスとメモリ・セマンティクスの組み合わせを用いる。図1に示す分散コンピュータ・システムの例示的な記憶動作では、ホスト・プロセッサ・ノード102などのホスト・プロセッサ・ノードが、チャンネル・セマンティクスを使用してRAIDサブシステムIPSOE172にディスク書き込みコマンドを送信することにより記憶動作を開始する。RAIDサブシステムは、コマンドを調べ、メモリ・セマンティクスを使用してホスト・プロセッサ・ノードのメモリ空間から直接データ・バッファを読み出す。データ・バッファが読み出されると、RAIDサブシステムは、チャンネル・セマンティクスを使用して、I/O完了メッセージをホスト・プロセッサ・ノードにプッシュする。

【0048】

一例示の実施形態では、図1に示す分散コンピュータ・システムは、仮想アドレスと仮想メモリの保護機構を用いてすべてのメモリへの正確で適正なアクセスを保証する動作を行う。このような分散コンピュータ・システムで実行されるアプリケーションは、どの動作にも物理的なアドレス指定を使用する必要がない。

【0049】

次いで図2に、本発明の好ましい実施形態によるホスト・プロセッサ・ノードの機能ブロック図を示す。ホスト・プロセッサ・ノード200は、図1のホスト・プロセッサ・ノード102などのホスト・プロセッサ・ノードの一例である。

【0050】

この例では、図2に示すホスト・プロセッサ・ノード200は、ホスト・プロセッサ・ノード200で実行されるプロセスであるコンシューマの集合202~208を含む。ホスト・プロセッサ・ノード200は、IPスイート・オフロード・エンジン(IPSOE

10

20

30

40

50

) 210とIPSOE212も含む。IPSOE210は、ポート214および216を含み、一方、IPSOE212はポート218および220を含む。各ポートは、リンクに接続する。これらのポートは、図1のIPネット100などの1つのIPネットのサブネット、または複数のIPネットのサブネットに接続することができる。

【0051】

コンシューマ202~208は、verbインタフェース222とメッセージおよびデータ・サービス224を介してIPネットにメッセージを転送する。verbインタフェースとは、基本的に、IPスイート・オフロード・エンジンの機能を抽象的に記述するものである。オペレーティング・システムは、verb機能の一部またはすべてを、自身のプログラミング・インタフェースを通じて公開することができる。基本的に、このインタフェースは、ホストの振る舞いを定義する。また、ホスト・プロセッサ・ノード200は、メッセージおよびデータ・サービス224を含み、これは、verb層より高レベルのインタフェースであり、IPSOE210およびIPSOE212を通じて受信されるメッセージとデータを処理するために使用される。メッセージおよびデータ・サービス224は、メッセージおよび他のデータを処理するためにコンシューマ202~208とのインタフェースを提供する。

10

【0052】

次いで図3に、本発明の好ましい実施形態によるIPスイート・オフロード・エンジンの図を示す。図3に示すIPスイート・オフロード・エンジン300Aは、IPSOEポート312A~316Aにメッセージを転送するために使用される、待ち行列のペア(QP)の集合302A~310Aを含む。IPSOEポート312A~316Aへのデータのバッファリングは、例えばIPバージョン6仕様のトラフィック・クラス・フィールドなどのネットワーク層のサービス品質フィールド(QOSF)318A~334Aを使用して運ばれる。ネットワーク層のサービス品質フィールドはそれぞれ、独自のフロー制御を有する。IETF(インターネット・エンジニアリング・タスク・フォース)標準のネットワーク・プロトコルを使用して、ネットワークに接続されたすべてのIPスイート・オフロード・エンジンのポートのリンク・アドレスとネットワーク・アドレスを設定する。そのようなプロトコルの2つは、アドレス解決プロトコル(ARP)とダイナミック・ホスト構成プロトコルである。メモリ変換および保護(MTP)338Aは、仮想アドレスを物理アドレスに変換し、アクセス権を確認する機構である。ダイレクト・メモリ・アクセス(DMA)340Aは、待ち行列のペア302A~310Aとの関係でメモリ350Aを使用して、ダイレクト・メモリ・アクセス動作を提供する。

20

30

【0053】

図3に示すIPSOE 300AなどのIPスイート・オフロード・エンジンは、1つで数千個の待ち行列ペアをサポートすることができる。各待ち行列のペアは、送信作業待ち行列(SWQ)と受信作業待ち行列(RWQ)で構成される。送信作業待ち行列は、チャンネル・セマンティクス・メッセージとメモリ・セマンティクス・メッセージの送信に使用される。受信作業待ち行列は、チャンネル・セマンティクス・メッセージを受信する。コンシューマは、オペレーティング・システムに固有のプログラミング・インタフェース(ここでは「verb」と称する)を呼び出して、作業要求(WR)を作業待ち行列に入れる。

40

【0054】

図4に、本発明の好ましい実施形態によるスイッチ300Bなどのスイッチを示す。スイッチ300Bは、QOSF306Bとして識別されるIPバージョン4のサービス・タイプ・フィールド等のリンク層またはネットワーク層のサービス品質フィールドを通じて複数のポート304と通信するパケット・リレー302Bを含む。一般に、スイッチ300Bなどのスイッチは、同じスイッチ上のあるポートから他のポートにフレームをルーティングすることができる。

【0055】

同様に、図5に、本発明の好ましい実施形態によるルータ300Cを示す。ルータ30

50

0 Cは、Q O S F 3 0 6 Cと識別されるI Pバージョン4のサービス・タイプ・フィールドなどのネットワーク層のサービス品質フィールドを通じて複数のポート3 0 4 Cと通信する、G R Hパケット・リレー3 0 2 Cなどのフレーム・リレーを含む。スイッチ3 0 0 Bと同様に、ルータ3 0 0 Cは、一般に、同じルータ上のあるポートから他のポートにフレームをルーティングすることができる。

【 0 0 5 6 】

次いで図6に、本発明の好ましい実施形態による作業要求の処理を説明する図を示す。図6では、コンシューマ4 0 6からの要求およびコンシューマ4 0 6に対する要求を処理するために、受信作業待ち行列4 0 0、送信作業待ち行列4 0 2、および完了待ち行列4 0 4がある。このコンシューマ4 0 6からの要求は、最終的にはハードウェア4 0 8に送られる。この例では、コンシューマ4 0 6は、作業要求4 1 0および4 1 2を生成し、作業の完了4 1 4を受け取る。図6に示すように、作業待ち行列に入れられた作業要求を作業待ち行列要素(W Q E)と呼ぶ。

10

【 0 0 5 7 】

送信作業待ち行列4 0 2は、I Pネット・ファブリックで送信されるデータを記述する作業待ち行列要素(W Q E)4 2 2 ~ 4 2 8を含んでいる。受信作業待ち行列4 0 0は、I Pネット・ファブリックから入ってくるチャンネル・セマンティクス・データをどこに置くかを記述する作業待ち行列要素(W Q E)4 1 6 ~ 4 2 0を含む。作業待ち行列要素は、I P S O Eのハードウェア4 0 8によって処理される。

【 0 0 5 8 】

v e r bは、完了待ち行列4 0 4から完了した作業を取り出す機構も提供する。図6に示すように、完了待ち行列4 0 4は、完了待ち行列要素(C Q E)4 3 0 ~ 4 3 6を含む。完了待ち行列要素は、これまでに完了した作業待ち行列要素についての情報を含む。完了待ち行列4 0 4を使用して、複数の待ち行列ペアについての単一の完了通知点を作成する。完了待ち行列要素は、完了待ち行列にあるデータ構造である。この要素は、完了した作業待ち行列要素を記述する。完了待ち行列要素は、待ち行列と完了した具体的な作業待ち行列要素を判定するのに十分な情報を含む。完了待ち行列のコンテキストは、個々の完了待ち行列へのポインタ、待ち行列の長さ、個々の完了待ち行列を管理するのに必要なその他の情報を含む情報のブロックである。

20

【 0 0 5 9 】

図6に示す送信作業待ち行列4 0 2にサポートされる作業要求の例は以下である。送信作業要求は、ローカル・データ・セグメントを、リモート・ノードの受信作業待ち行列要素によって参照されるデータ・セグメントにプッシュするチャンネル・セマンティクス動作である。例えば、作業待ち行列要素4 2 8は、データ・セグメント4 4 8 3、データ・セグメント5 4 4 0、およびデータ・セグメント6 4 4 2への参照を含む。送信作業待ち行列のデータ・セグメントはそれぞれ、実質的に連続したメモリ領域の一部を含む。ローカル・データ・セグメントを参照するために使用される仮想アドレスは、そのローカル待ち行列のペアを作成したプロセスのアドレス・コンテキスト内にある。

30

【 0 0 6 0 】

リモート・ダイレクト・メモリ・アクセス(R D M A)の読み出し作業要求は、リモート・ノードの実質的に連続したメモリ空間を読み出すメモリ・セマンティクス動作を提供する。メモリ空間は、メモリ領域の一部であっても、メモリ・ウィンドウの一部であってもよい。メモリ領域は、仮想アドレスと長さによって定義された、実質的に連続したメモリ・アドレスのあらかじめ登録されたセットを参照する。メモリ・ウィンドウは、あらかじめ登録された領域に結び付けられた、実質的に連続したメモリ・アドレスのセットを参照する。

40

【 0 0 6 1 】

R D M A読み出しの作業要求は、リモートのエンドノードの実質的に連続したメモリ空間から読み出し、実質的に連続したローカル・メモリ空間にそのデータを書き込む。送信作業要求と同様に、ローカル・データ・セグメントを参照するためにR D M Aの読み出し

50

作業待ち行列要素によって使用される仮想アドレスは、そのローカル待ち行列のペアを作成したプロセスのアドレス・コンテキストにある。リモートの仮想アドレスは、そのRDMA読み出し作業待ち行列要素が対象とするリモートの待ち行列ペアを所有するプロセスのアドレス・コンテキスト内にある。

【0062】

RDMA書き込みの作業待ち行列要素は、リモート・ノードの実質的に連続したメモリ空間に書き込むメモリ・セマンティクス動作を提供する。例えば、受信作業待ち行列400中の作業待ち行列要素416は、データ・セグメント1 444、データ・セグメント2 446、およびデータ・セグメント3 448を参照する。RDMA書き込み作業待ち行列要素は、ローカルの実質的に連続したメモリ空間とそのローカル・メモリ空間が書き込まれるリモートのメモリ空間の仮想アドレスとのスキャット・リストを含んでいる。RDMA FetchOpの作業待ち行列要素は、リモートのワードにアトミックな動作を行うメモリ・セマンティクス動作を提供する。RDMA FetchOpの作業待ち行列要素は、RDMA Read、Modify、およびRDMA Write動作を組み合わせた動作である。RDMA FetchOpの作業待ち行列要素は、比較し、等しければ交換する(Compare and Swap if equal)など数個の読み出し - 変更 - 書き込みの動作をサポートすることができる。RDMA FetchOpは、現在のRDMA over IPの標準化の取り組みでは含まれていないが、一部の実施では付加価値のある機能として使用することが可能であるためここに記載する。

【0063】

バインド(非バインド)のリモート・アクセス・キー(R_Key)の作業待ち行列要素は、メモリ・ウィンドウをメモリ領域に関連付ける(関連付けを解除する)ことによりメモリ・ウィンドウを変更(破壊)するコマンドを、IPスイート・オフロード・エンジンのハードウェアに提供する。R_Keyは、各RDMAアクセスの一部であり、リモート・プロセスがバッファへのアクセスを許可したことを確認するために使用される。

【0064】

一実施形態では、図6に示す受信作業待ち行列400は、1タイプの作業待ち行列要素だけをサポートし、これを受信作業待ち行列要素と呼ぶ。受信作業待ち行列要素は、入ってくる送信メッセージが書き込まれるローカルのメモリ空間を記述するチャネル・セマンティクス動作を提供する。受信作業待ち行列要素は、いくつかの実質的に連続したメモリ空間を記述するスキャット・リストを含む。送られてきた送信メッセージは、そのメモリ空間に書き込まれる。仮想アドレスは、そのローカル待ち行列のペアを作成したプロセスのアドレス・コンテキスト内にある。

【0065】

プロセッサ間通信の場合は、ユーザモードのソフトウェア・プロセスが、メモリ中のバッファがある場所から、直接待ち行列ペアを通じてデータを転送する。一実施形態では、待ち行列のペアを通じたこの転送は、オペレーティング・システムを回避し、ホスト命令サイクルをほとんど消費しない。待ち行列のペアは、オペレーティング・システム・カーネルを関与させずに、プロセッサ・コピーのないデータ転送を可能にする。プロセッサ・コピーのないデータ転送は、広帯域幅で低待ち時間の通信を効率的にサポートする。

【0066】

待ち行列のペアが作成される際、その待ち行列のペアは、選択されたタイプのトランスポート・サービスを提供するように設定される。一実施形態では、本発明の好ましい実施形態を実施する分散コンピュータ・システムは、TCP、SCTP、UDPの3タイプのトランスポート・サービスをサポートする。

【0067】

TCPおよびSCTPは、ローカルの待ち行列ペアを、ただ1つのリモート待ち行列のペアに関連付ける。TCPおよびSCTPでは、プロセスが、TCPおよびSCTPがIPネット・ファブリックを介して通信する各プロセスの待ち行列ペアを作成することが必要とされる。したがって、N個のホスト・プロセッサ・ノードがそれぞれP個のプロセス

10

20

30

40

50

を含み、各ノードのP個のプロセスすべてがすべての他のノードのすべてのプロセスとの通信を要求する場合、各ホスト・プロセッサ・ノードは、 $P^2 \times (N - 1)$ 個の待ち行列ペアを必要とする。さらに、プロセスは、待ち行列のペアを、同じIPSOEの別の待ち行列ペアに関連付けることができる。

【0068】

TCPまたはSCTPを用いて分散プロセス間で通信する分散コンピュータ・システムの一部を図7に概略的に示す。図7の分散コンピュータ・システム500は、ホスト・プロセッサ・ノード1、ホスト・プロセッサ・ノード2、およびホスト・プロセッサ・ノード3を含む。ホスト・プロセッサ・ノード1は、プロセスA 510を含む。ホスト・プロセッサ・ノード3は、プロセッサC 520とプロセッサD 530を含む。ホスト・プロセッサ・ノード2は、プロセッサE 540を含む。

10

【0069】

ホスト・プロセッサ・ノード1は、待ち行列ペア4、6、および7を含み、待ち行列ペアはそれぞれ、送信作業待ち行列と受信作業待ち行列を有する。ホスト・プロセッサ・ノード2は、待ち行列ペア9を有し、ホスト・プロセッサ・ノード3は、待ち行列ペア2および5を有する。分散コンピュータ・システム500のTCPまたはSCTPは、ローカルの待ち行列ペアをただ1つのリモート待ち行列ペアに関連付ける。したがって、待ち行列ペア4を使用して待ち行列ペア2と通信し、待ち行列ペア7を使用して待ち行列ペア5と通信し、待ち行列ペア6を使用して待ち行列ペア9と通信する。

【0070】

20

TCPまたはSCTPである送信待ち行列に置かれたWQEは、関連付けられた待ち行列ペアの受信WQEによって参照される受信メモリ空間にデータを書き込ませる。RDMA動作は、関連付けられた待ち行列ペアのアドレス空間に作用する。

【0071】

本発明の一実施形態では、ハードウェアが連続番号を保持し、すべてのフレーム転送に肯定応答を返すので、TCPまたはSCTPに信頼性を持たせられる。ハードウェアとIPネットのドライバ・ソフトウェアの組み合わせが、失敗した通信を再度試みる。待ち行列ペアのプロセス・クライアントは、ビット誤り、受信のアンダーラン、およびネットワークの輻輳がある場合でも、信頼性のある通信を得る。IPネット・ファブリックに代替の経路が存在する場合は、ファブリック・スイッチ、リンク、またはIPスイート・オフロード・エンジン・ポートの障害があっても、信頼性のある通信を維持することができる。

30

【0072】

また、肯定応答を用いて、IPネット・ファブリックを介して確実にデータを配信することができる。肯定応答は、プロセス・レベルの肯定応答、すなわち受信側プロセスがデータを消費したことを証明する肯定応答であっても、そうでなくともよい。あるいは、肯定応答は、データがその宛先に到達したことを知らせるだけの肯定応答であってもよい。

【0073】

ユーザ・データグラム・プロトコルは、コネクションレスである。UDPは、新しいスイッチ、ルータ、およびエンドノードを発見し、所与の分散コンピュータ・システムに組み込むために管理アプリケーションによって用いられる。UDPは、TCPまたはSCTPが提供する信頼性の保証は提供しない。そのため、UDPは、各エンドノードにTCPおよびSCTPに比べて少ない情報を維持して動作する。

40

【0074】

次いで図8に、本発明の好ましい実施形態によるデータ・フレームの図を示す。データ・フレームは、IPネット・ファブリックを通じてルーティングされる情報の単位である。データ・フレームは、エンドノード間の構造物であり、したがってエンドノードによって作成され、消費される。IPSOEを宛先とするフレームの場合、データ・フレームは、IPネット・ファブリックのスイッチとルータには生成も消費もされない。代わりに、IPSOEを宛先とするデータ・フレームについては、スイッチとルータは、単に、要求

50

フレームまたは肯定応答フレームを最終的な宛先のより近くに移動し、その過程でリンク・ヘッダ・フィールドを変更する。ルータは、フレームがサブネットの境界を越える際にフレームのネットワーク・ヘッダを変更してよい。1つのサブネットを横断する間、単一のフレームは、単一のサービス・レベルにとどまる。

【0075】

メッセージ・データ600は、データ・セグメント1 602、データ・セグメント2 604、データ・セグメント3 606を含み、これらは、図6に示すデータ・セグメントと同様である。この例では、これらのデータ・セグメントはフレーム608を形成し、フレーム608は、データ・フレーム612中のフレーム・ペイロード610に入れられる。また、データ・フレーム612は、エラーの検査に使用される巡回冗長検査(CRC)614を含む。また、ルーティング・ヘッダ616とトランスポート・ヘッダ618が、データ・フレーム612中にある。ルーティング・ヘッダ616は、データ・フレーム612の送信元ポートと宛先ポートを識別するために使用される。この例のトランスポート・ヘッダ618は、データ・フレーム612の連続番号と送信元ポート番号および宛先ポート番号を指定する。連続番号は、通信が確立される時に初期化され、フレーム・ヘッダ、DDP/RDMAヘッダ、データ・ペイロード、およびCRCの1バイトごとに1ずつ増分する。この例のフレーム・ヘッダ620は、フレームに関連付けられた宛先待ち行列ペアの番号、および、DDP/RDMA(Direct Data PlacementまたはRemote Direct Memory Accessあるいはその両方)ヘッダとデータ・ペイロードとCRCとを合計した長さを指定する。DDP/RDMAヘッダ622は、データ・ペイロードのメッセージ識別子と配置情報を指定する。メッセージ識別子は、1つのメッセージの一部分であるすべてのフレームに一定である。メッセージ識別子の例には、例えば、send、write R DMA、およびread R DMAが含まれる。

【0076】

図9に、要求と肯定応答の処理の例を説明するために、分散コンピュータ・システム700などの分散コンピュータ・システムの一部を示す。図9の分散コンピュータ・システム700は、プロセスA 716を実行するホスト・プロセッサ・ノード702と、プロセスB 718を実行するホスト・プロセッサ・ノード704を含む。ホスト・プロセッサ・ノード702は、IPSOE706を含む。ホスト・プロセッサ・ノード704は、IPSOE708を含む。図9の分散コンピュータ・システムは、スイッチ712と714を含むIPネット・ファブリック710を含む。IPネット・ファブリックは、IPSOE706をスイッチ712に結合するリンク、スイッチ712をスイッチ714に結合するリンク、およびIPSOE708をスイッチ714に結合するリンクを含む。

【0077】

この例示的では、ホスト・プロセッサ・ノード702は、クライアント・プロセスAを含む。ホスト・プロセッサ・ノード704は、クライアント・プロセスBを含む。クライアント・プロセスAは、送信待ち行列724と受信待ち行列726からなる待ち行列ペア23 720を通じてホストIPSOE706と対話する。クライアント・プロセスBは、送信待ち行列728と受信待ち行列730からなる待ち行列ペア24 722を通じてホストIPSOE708と対話する。待ち行列のペア23および24は、送信作業待ち行列と受信作業待ち行列を含むデータ構造である。

【0078】

プロセスAは、待ち行列ペア23の送信待ち行列に作業待ち行列要素を入れることによりメッセージ要求を開始する。そのような作業待ち行列要素を図6に示す。クライアント・プロセスAのメッセージ要求は、送信作業待ち行列要素に含まれるギャザー・リストによって参照される。ギャザー・リストの各データ・セグメントは、それぞれ図8で個々にメッセージ部分1、2、および3を保持するデータ・セグメント1、2、3として表すようなメッセージの一部を含む、実質的に連続したローカル・メモリ領域の一部をポイントする。

【0079】

10

20

30

40

50

ホストIPSOE706のハードウェアは、作業待ち行列要素を読み出し、実質的に連続したバッファに記憶されたメッセージを、図8に示すデータ・フレームのようなデータ・フレームに区分する。データ・フレームはIPネット・ファブリックを通じてルーティングされ、信頼性のある転送サービスのために、最終的な宛先エンドノードによって肯定応答が返される。肯定応答が成功しない場合、データ・フレームは、送信元のエンドノードから再送信される。データ・フレームは、送信元エンドノードによって生成され、宛先エンドノードによって消費される。

【0080】

図10を参照すると、本発明の好ましい実施形態による分散ネットワーク・システムで使用されるネットワーク・アドレッシングを説明する図が示される。ホスト名は、ホスト・プロセッサ・ノードやI/Oアダプタ・ノードなどのホスト・ノードの論理的な識別を提供する。ホスト名は、メッセージが、そのホスト名によって指定されるエンドノードに存在するプロセスを宛先とするように、メッセージのエンドポイントを識別する。したがって、1つのノードにつき1つのホスト名があるが、1つのノードは複数のIPSOEを有することができる。

10

【0081】

1つのリンク層アドレス(例えばイーサネット(R)のメディア・アクセス層アドレス)804が、エンドノード・コンポーネント802の各ポート806に割り当てられる。コンポーネントは、IPSOE、スイッチ、またはルータである。すべてのIPSOEコンポーネントおよびルータ・コンポーネントは、MACアドレスを持たなければならない。スイッチ上のメディア・アクセス・ポイントにもMACアドレスが割り当てられる。

20

【0082】

エンドノード・コンポーネント802の各ポート806に、1つのネットワーク・アドレス(例えばIPアドレス)812が割り当てられる。コンポーネントは、IPSOE、スイッチ、またはルータである。すべてのIPSOEおよびルータ・コンポーネントは、ネットワーク・アドレスを持たなければならない。スイッチのメディア・アクセス・ポイントにもMACアドレスが割り当てられる。

【0083】

スイッチ810の各ポートは、関連付けられたリンク層アドレスを持たない。しかし、スイッチ810は、関連付けられたリンク層アドレス816とネットワーク層アドレス808を有するメディア・アクセス・ポート814を有することができる。

30

【0084】

本発明の好ましい実施形態による分散コンピュータ・システムの一部を図11に示す。分散コンピュータ・システム900は、サブネット902とサブネット904を含む。サブネット902は、ホスト・プロセッサ・ノード906、908、および910などのエンドノードを含む。サブネット904は、ホスト・プロセッサ・ノード912および914などのエンドノードを含む。サブネット902は、スイッチ916および918を含む。サブネット904は、スイッチ920および922を含む。

【0085】

ルータは、サブネットを作成し、接続する。例えば、サブネット902は、ルータ924および926によりサブネット904に接続される。例示的な一実施形態では、サブネットは、最高で216個のエンドノード、スイッチ、およびルータを有する。

40

【0086】

サブネットは、1つの単位として管理されるエンドノードとカスケード・スイッチの群と定義される。通例、サブネットは、単一の地理的範囲または機能範囲を占める。例えば、1つの部屋にある1つのコンピュータ・システムをサブネットと定義することができる。一実施形態では、サブネット内のスイッチは、非常に高速のメッセージのワームホール・ルーティングまたはカットスルー・ルーティングを行うことができる。

【0087】

サブネット内のスイッチは、入ってくるメッセージ・フレームをスイッチが迅速かつ効

50

率的にルーティングできるようにそのサブネット内で一意である宛先リンク層アドレス（例えばMACアドレス）を調べる。一実施形態では、スイッチは、比較的単純な回路であり、通例は単一の集積回路として実施される。サブネットは、カスケード・スイッチから形成された数百または数千のエンドノードを有することができる。

【0088】

図11に示すように、はるかに大規模なシステムへの拡張のために、サブネットは、ルータ924および926などのルータと接続される。ルータは、宛先のネットワーク層アドレス（例えばIPアドレス）を解釈し、フレームをルーティングする。

【0089】

スイッチの例示的な実施形態を図4に概略的に示す。スイッチまたはルータの各I/P経路は、ポートを有する。一般に、スイッチは、同じスイッチ上のあるポートから他のポートにフレームをルーティングすることができる。サブネット902やサブネット904などのサブネット内では、送信元ポートから宛先ポートまでの経路は、宛先ホストIPSOEポートのリンク層アドレス（例えばMACアドレス）によって決定される。サブネット間では、経路は、宛先のIPSOEポートのネットワーク層アドレス（IPアドレス）と、ルータポートのリンク層アドレス（例えばMACアドレス）によって決定され、その経路を使用して宛先のサブネットに到達する。

【0090】

一実施形態では、要求フレームとその要求フレームに対応する肯定応答（ACK）フレームによって使用される経路は、対称である必要はない。明白なルーティングを用いる一実施形態では、スイッチは、リンク層アドレス（例えばMACアドレス）に基づいて出力ポートを選択する。一実施形態では、スイッチは、そのスイッチ中のすべての入力ポートに対して1つのルーティング決定基準のセットを用いる。一例示的な実施形態では、ルーティング決定基準は、1つのルーティング・テーブルに含まれる。代替実施形態では、スイッチは、入力ポートごとに別個の基準のセットを用いる。

【0091】

本発明の好ましい実施形態の分散コンピュータ・システムにおけるデータ・トランザクションは、通例、数個のハードウェア・ステップとソフトウェア・ステップからなる。クライアント・プロセスのデータ・トランスポート・サービスは、ユーザモードまたはカーネルモードのプロセスである。クライアント・プロセスは、図3、5、および8に示す待ち行列ペアのような1つまたは複数の待ち行列ペアを通じて、IPスイート・オフロード・エンジン・ハードウェアにアクセスする。クライアント・プロセスは、オペレーティング・システムに固有のプログラミング・インタフェースを呼び出し、このインタフェースをここでは「verb」と呼ぶ。verbを実施するソフトウェア・コードは、所与の待ち行列ペアの作業待ち行列に作業待ち行列要素を入れる。

【0092】

作業待ち行列要素の挿入には多くの可能な方法があり、また可能な作業待ち行列要素のフォーマットも多く存在し、それにより様々なコスト/パフォーマンスの設計点が可能になるが、相互運用性には影響を与えない。ただし、ユーザ・プロセスは、明確に定義された方式でverbと通信しなければならない。また、IPネット・ファブリックを通じて送信されるデータのフォーマットとプロトコルは、デバイスが異種のベンダ環境で相互運用できるように十分に詳細に規定されなければならない。

【0093】

一実施形態では、IPSOEハードウェアは、作業待ち行列要素の挿入を検出し、その作業待ち行列要素にアクセスする。この実施形態では、IPSOEハードウェアは、作業待ち行列要素の仮想アドレスを変換および検証し、データにアクセスする。

【0094】

発信されるメッセージは、1つまたは複数のデータ・フレームに分割される。一実施形態では、IPSOEハードウェアは、DDP/RDMAヘッダ、フレーム・ヘッダ、およびCRC、トランスポート・ヘッダ、およびネットワーク・ヘッダを各フレームに追加す

10

20

30

40

50

る。トランスポート・ヘッダは、連続番号と他のトランスポート情報を含む。ネットワーク・ヘッダは、宛先IPアドレスや他のネットワーク・ルーティング情報などのルーティング情報を含む。リンク・ヘッダは、宛先リンク層アドレス（例えばMACアドレス）や他のローカル・ルーティング情報を含む。

【0095】

TCPまたはSCTPを用いる場合には、要求データ・フレームがその宛先エンドノードに到達すると、肯定応答データ・フレームが宛先エンドノードによって使用されて、その要求データ・フレームの送信者に、要求データ・フレームが宛先で有効性が検証され、受け付けられたことを知らせる。肯定応答データ・フレームは、有効で、受け付けられた1つまたは複数の要求データ・フレームを肯定応答する。要求者は、肯定応答を受け取る前に、複数の未処理の要求データ・フレームを有することができる。一実施形態では、複数の未処理のメッセージ、すなわち要求データ・フレームの数は、待ち行列のペアが作成される時に決定される。

10

【0096】

本発明の好ましい実施形態を実施する階層アーキテクチャ1000の一実施形態を図12の形で概略的に図示する。図12の階層アーキテクチャ図は、各種のデータ通信経路の層と、層の間で渡されるデータと制御情報の編成を示す。

【0097】

IPSOEエンドノード・プロトコル層（例えばエンドノード1011により用いられる）は、コンシューマ1003によって定義される上層プロトコル1002、トランスポート層1004、ネットワーク層1006、リンク層1008、および物理層1010を含む。スイッチ層（例えばスイッチ1013によって用いられる）は、リンク層1008と物理層1010を含む。ルータ層（例えばルータ1015によって用いられる）は、ネットワーク層1006、リンク層1008、および物理層1010を含む。

20

【0098】

階層アーキテクチャ1000は、一般に、コンシューマ1003と1005の間でデータを転送するコンシューマ動作1012を実現するために、従来の通信スタックの概要に従う。例えばエンドノード1011のプロトコル層に関して、上層プロトコル1002は、verbを使用してトランスポート層1004でメッセージを作成する。トランスポート層1004は、メッセージ1014をネットワーク層1006に渡す。ネットワーク層1006は、ネットワーク・サブネット間でフレームをルーティングする1016。リンク層1008は、ネットワーク・サブネット内でフレームをルーティングする1018。物理層1010は、ビットまたはビットのグループを他のデバイスの物理層に送信する。これらの各層は、上または下の層が各自の機能をどのように行うかを認識しない。

30

【0099】

コンシューマ1003と1005は、エンドノード間の通信のために他の層を用いるアプリケーションまたはプロセスを表す。トランスポート層1004は、終端間のメッセージ移動を提供する。一実施形態では、トランスポート層は、上記のように、従来のTCP、RDMA over TCP、SCTP、およびUDPの4タイプのトランスポート・サービスを提供する。ネットワーク層1006は、1つまたは複数のサブネットを通じて宛先エンドノードへのフレームのルーティングを行う。リンク層1008は、複数のリンクを介して、フローが制御され1020、エラーが検査され、優先順位がつけられたフレームの配信を行う。

40

【0100】

物理層1010は、技術に依存したビット送信を行う。ビットまたはビットのグループは、リンク1022、1024、および1026を介して物理層間で渡される。リンクは、プリント回路銅線トレース、銅線ケーブル、光ケーブル、または他の適切なリンクを用いて実施することができる。図13に、本発明の例示的实施形態による、ホスト・ソフトウェアとのRNICインタフェースを説明する図を示す。verbコンシューマ1156は、verb1152と、verbドライバおよびライブラリ1148を通じて、主要R

50

N I C 1 1 0 0 と代替 R N I C 1 1 0 4 にアクセスする。v e r b コンシューマ 1 1 5 6 は、それぞれ R N I C 管理照会 v e r b 1 1 2 4 と 1 1 3 6 を呼び出すことにより、主要 R N I C 1 1 0 0 と代替 R N I C 1 1 0 4 がスイッチオーバー/スイッチバックをサポートしていることを判断する。R N I C 管理照会 v e r b は、R N I C の機能を返し、この機能は、この例示的实施形態では、R N I C がスイッチオーバー/スイッチバック (S / S) をサポートしていることを示すフィールドを含む。v e r b コンシューマ 1 1 5 6 は次いで、R N I C 管理変更 v e r b 1 1 2 4 および 1 1 3 6 を使用して、ある範囲の待ち行列ペア (Q P)、完了待ち行列 (C Q)、およびメモリ変換および保護テーブル (T P T) エントリを、S / S と非 S / S サポートに割り当てる。R N I C 変更 v e r b の完了が成功すると、主要 R N I C 1 1 0 0 と代替 R N I C 1 1 0 4 は、共通の Q P、C Q、およびメモリ T P T の範囲を共有する。

10

【 0 1 0 1 】

v e r b コンシューマ 1 1 5 6 は、C Q 作成 v e r b を使用して C Q を作成し、この作成では、その C Q が、C Q 1 1 7 6 および 1 1 8 0 などのように S / S をサポートするか、または Q P 1 1 6 8 および 1 1 8 8 のようにサポートしないかを選択する。v e r b コンシューマが S / S サポートを選択した場合、C Q 作成 v e r b は、主要 R N I C、主要 R N I C のポート、代替の R N I C、および代替 R N I C ポートを識別する追加的な修飾子を含む。R N I C 識別子は、E U I - 6 4 ビット識別子とすることができる。ポート識別子は、ポート番号とすることができる。別の代替法は、R N I C の M A C アドレスを使用して R N I C とポートの両方を識別するものである。v e r b コンシューマが S / S サポートを選択した場合、R N I C とポートの識別子は、C Q コンテキスト 1 1 1 0 および 1 1 1 8 に置かれる。v e r b コンシューマが S / S サポートを選択しない場合は、R N I C とポートの識別子は、C Q コンテキスト 1 1 0 6 および 1 1 2 2 に置かれない。

20

【 0 1 0 2 】

v e r b コンシューマ 1 1 5 6 は、メモリ領域登録 v e r b (例えばメモリ領域登録、共有メモリ領域登録、物理メモリ領域登録) の 1 つを使用してメモリ領域を登録し、この際に、そのメモリ領域が、メモリ領域 1 1 2 8 および 1 1 3 2 のように S / S をサポートするか、またはメモリ領域 1 1 9 4 および 1 1 9 6 のようにサポートしないかを選択する。v e r b コンシューマが S / S サポートを選択した場合、メモリ領域登録 v e r b は、主要 R N I C、主要 R N I C ポート、代替 R N I C、および代替 R N I C ポートを識別する追加的な修飾子を含む。R N I C 識別子は、例えば E U I - 6 4 ビット識別子である。ポート識別子にはポート番号を使用することができる。別の代替法は、R N I C の M A C アドレスを使用して N I C とポートの両方を識別するものである。v e r b コンシューマが S / S サポートを選択した場合、R N I C とポートの識別子は、そのメモリ領域のメモリ変換および保護テーブル (T P T) のエントリ 1 1 2 8 および 1 1 3 2 に置かれることができる。v e r b コンシューマが S / S サポートを選択しない場合、R N I C とポートの識別子は、メモリ領域のメモリ T P T エントリ 1 1 9 4 および 1 1 9 6 に置かれない。

30

【 0 1 0 3 】

v e r b コンシューマ 1 1 5 6 は、Q P 作成 v e r b を使用して Q P を作成し、この際に、その Q P が Q P 1 1 7 2 および 1 1 8 4 のように S / S をサポートするか、または Q P 1 1 6 4 および 1 1 9 2 のようにサポートしないかを選択する。v e r b コンシューマが S / S サポートを選択した場合、Q P 作成 v e r b は、主要 R N I C、主要 R N I C ポート、代替 R N I C、および代替 R N I C ポートを識別する追加的な修飾子を含む。R N I C 識別子は、例えば E U I - 6 4 ビット識別子である。ポート識別子にはポート番号を用いることができる。別の代替法は、R N I C の M A C アドレスを使用して R N I C とポートの両方を識別するものである。v e r b コンシューマが S / S サポートを選択した場合、R N I C とポートの識別子は、Q P コンテキスト 1 1 1 2 および 1 1 1 6 に置かれる。v e r b コンシューマが S / S サポートを選択しない場合、R N I C とポートの識別子は、Q P コンテキスト 1 1 0 8 および 1 1 2 0 に置かれない。

40

50

【0104】

S / S Q P が例えば停止のために生じたスイッチオーバーの後に代替 R N I C なしで動作している時、v e r b コンシューマ 1 1 5 6 は、Q P 変更 v e r b を使用して Q P 1 1 7 2 および 1 1 8 4 などの S / S Q P を停止させ、代替の R N I C に同じ Q P 番号を再度割り当てる。Q P 変更 v e r b は、主要 R N I C、主要 R N I C ポート、代替 R N I C、および代替 R N I C ポートを識別する修飾子を含む。R N I C 識別子は、例えば E U I - 6 4 ビット識別子である。ポート識別子にはポート番号を用いることができる。別の代替法は、R N I C の M A C アドレスを用いて N I C とポートの両方を識別するものである。R N I C とポートの識別子は、Q P コンテキスト 1 1 1 2 および 1 1 1 6 に置かれる。

10

【0105】

S / S R N I C が例えば停止のために生じたスイッチオーバーの後に代替 R N I C なしで動作している場合、v e r b コンシューマ 1 1 5 6 は、再同期登録 v e r b を使用して、1 1 2 8 などのすべてのメモリ T P T エントリを基本的に代替 R N I C に登録する。

【0106】

代替 R N I C なしに動作している所与の S / S C Q に関連付けられたすべての S / S Q P が停止されると、v e r b コンシューマ 1 1 5 6 は、C Q 変更 v e r b を使用して、C Q 1 1 7 6 および 1 1 8 0 などの S / S C Q を停止させ、代替の R N I C に同じ C Q 番号を再度割り当てる。C Q 変更 v e r b は、主要 R N I C、主要 R N I C ポート、代替 R N I C、および代替 R N I C ポートを識別する修飾子を含む。R N I C 識別子は例えば E U I - 6 4 ビット識別子である。ポート識別子にはポート番号を使用することができる。別の代替法は、R N I C の M A C アドレスを使用して R N I C とポートの両方を識別するものである。R N I C とポートの識別子は、C Q コンテキスト 1 1 1 0 および 1 1 1 8 に置かれる。

20

【0107】

S / S C Q が停止されると、C Q 変更 v e r b を使用して、C Q をアクティブな状態（使用可能状態など）にし、C Q に関連付けられた各 S / S Q P について、C Q 変更 v e r b を使用して Q P をアクティブな状態（送信可能状態など）にする。

【0108】

図 1 4 に、本発明の好ましい実施形態による、2 つの R N I C によって維持されるスイッチオーバーとスイッチバックに関連する例示的な接続状態を説明する図を提供する。好ましい実施形態では、各主要 R N I C、代替 R N I C、または単一の R N I C は、図 1 4 に示す状態情報を保持する。状態情報は、例えば、R N I C 内、R N I C からアクセスできるホストの記憶場所、それら 2 つの組み合わせなどに保持することができる。図 1 4 は、状態情報が R N I C に保持される例示的事例を示す。

30

【0109】

Q P 1 1 7 2 などの各 R N I C Q P は、Q P C T 1 2 0 0 などの Q P コンテキスト・テーブル（Q P C T）内に、Q P C E 0 1 2 0 4、Q P C E 1 1 2 0 8、および Q P C E N 1 2 1 2 などの Q P コンテキスト・エントリを有する。Q P C T 1 2 0 0 は、Q P C T アドレス 1 2 4 8 などの Q P C T のアドレスと、Q P C T 長 1 2 5 2 など Q P C T の長さを含む、Q P C T レジスタ 1 2 4 4 などの Q P コンテキスト・テーブル・レジスタを通じてアクセスされる。Q P コンテキスト・テーブル 1 2 0 0 の各エントリは、固定されたサイズであるが、Q P C T エントリのサイズが可変の Q P C T も使用することができる。Q P C T の各エントリは、そのエントリについて参照される Q P 番号に関連付けられた関連するコンテキスト情報を保持する。例えば、Q P C E N 1 2 1 2 は、Q P 番号 N に関連付けられた Q P 状態情報を含んでいる。Q P 番号 N の状態には、従来の Q P コンテキスト状態 1 2 2 0 などの従来の Q P コンテキスト状態、作業 S Q スイッチオーバー / スイッチバック・コンテキスト 1 2 2 4 などの実効 S Q スイッチオーバー / スイッチバック・コンテキスト状態、コミット済み S Q スイッチオーバー / スイッチバック・コンテキスト 1 2 2 8 などのコミット済み S Q スイッチオーバー / スイッチバック・コンテキスト

40

50

状態、作業 R Q スイッチオーバー / スイッチバック・コンテキスト 1 2 3 2 などの作業 R Q スイッチオーバー / スイッチバック・コンテキスト状態、コミット済み R Q スイッチオーバー / スイッチバック・コンテキスト 1 2 3 6 などのコミット済み R Q スイッチオーバー / スイッチバック・コンテキスト状態、および、主要 R N I C と代替 R N I C の識別コンテキスト 1 2 4 0 などの主要 R N I C および代替 R N I C の識別コンテキスト状態が含まれる。

【 0 1 1 0 】

従来の Q P コンテキスト状態 1 2 2 0 は、Q P の状態、Q P に関連付けられたサービスのタイプ、Q P に関連付けられた T C P ソースポートと宛先番号、Q P に関連付けられた I P の送信元アドレスおよび宛先アドレスなどの、周知の Q P コンテキスト状態情報を含む。

10

【 0 1 1 1 】

実効 S Q スイッチオーバー / スイッチバック・コンテキスト 1 2 2 4 は、次の情報を含み、この情報は Q P が作成される際にゼロに設定される。

【 0 1 1 2 】

【表 1】

現在の実効送信 W Q E 番号
現在の送信 W Q E 中への実効バイト・オフセット
次の送信 T C P セグメントの実効連続番号
次の送信 T C P セグメントの値の実効ウィンドウ
現在の実効受信 A c k W Q E 番号
現在の受信 A c k W Q E 中への実効バイト・オフセット
次の R c v A c k T C P セグメントの実効連続番号
次の R c v A c k T C P セグメントの実効ウィンドウ値

20

【 0 1 1 3 】

すべての実効 S Q スイッチオーバー / スイッチバック・フィールドは、代替 R N I C によりチェックポイント A C K を通じてチェックポイントされていない情報を含む。主要 R N I C と代替 R N I C は両方とも、これらのフィールドをすべて含む。「現在の実効送信 W Q E 番号」は、現在送信待ち行列で処理されている W Q E を識別する。「現在の送信 W Q E への実効バイト・オフセット」は、現在の実効送信 W Q E について処理される次のバイトを識別する。「次の送信 T C P セグメントの実効連続番号」は、現在の送信 W Q E 中への実効バイト・オフセットに使用される T C P 連続番号を識別する。「次の送信 T C P セグメントの実効ウィンドウ値」は、現在の送信 W Q E への実効バイト・オフセットから送信が開始される次の T C P セグメントに利用できる T C P ウィンドウのサイズを識別する。「現在の実効受信 A C K W Q E 番号」は、リモート・ノードによって肯定応答が返された、送信待ち行列で現在処理中の W Q E を識別する。「現在の受信 A C K W Q E 中への実効バイト・オフセット」は、現在の実効受信 A C K W Q E についてリモート・ノードから肯定応答される次のバイトを識別する。「次の受信 A C K T C P セグメントの実効連続番号」は、リモート・ノードから受け取られることが予想され、現在の受信 A C K W Q E 中への実効バイト・オフセットに関連付けられた次の T C P 連続番号を識別する。「次の受信 A C K T C P セグメントの実効ウィンドウ値」は、リモート・ノードから送信され、現在の受信 A C K W Q E 中への実効バイト・オフセットに関連付けられた最後の T C P ウィンドウ・サイズの結果生じた送信ウィンドウの変化を識別する。

30

40

【 0 1 1 4 】

コミット済み S Q スイッチオーバー / スイッチバック・コンテキスト 1 2 2 8 は、次の

50

情報を含み、この情報は、QPが作成される際にゼロに設定される。

【0115】

【表2】

コミット済み現在の送信WQE番号
現在の送信WQE中へのコミット済みバイト・オフセット
次の送信TCPセグメントのコミット済み連続番号
次の送信TCPセグメントの値のコミット済みウィンドウ
コミット済み現在の受信AckWQE番号
現在の受信AckWQE中へのコミット済みバイト・オフセット
次のRcvAckTCPセグメントのコミット済み連続番号
次のRcvAckTCPセグメントのコミット済みウィンドウ値

10

【0116】

すべてのコミット済みSQスイッチオーバー/スイッチバック・フィールドは、代替RNICによってチェックポイントACKを通じてチェックポイントされた情報を含む。主要RNICと代替RNICはいずれも、これらのフィールドをすべて含む。「コミット済み現在の送信WQE番号」は、チェックポイントされ、現在送信待ち行列で処理中のWQEを識別する。「現在の送信WQEへのコミット済みバイト・オフセット」は、コミット済み現在の送信WQEについてチェックポイントされる次のバイトを識別する。「次の送信TCPセグメントのコミット済み連続番号」は、現在の送信WQE中へのコミット済みバイト・オフセットに使用されるTCP連続番号を識別する。「次の送信TCPセグメントのコミット済みウィンドウ値」は、現在の送信WQE中へのコミット済みバイト・オフセットから送信が開始される次のTCPセグメントに利用できるTCPウィンドウのサイズを識別する。「コミット済み現在の受信ACK WQE番号」は、リモート・ノードによって肯定応答が返され、送信待ち行列でチェックポイントされた最も新しいWQEを識別する。「現在の受信ACK WQEへのコミット済みバイト・オフセット」は、コミット済み現在の受信ACK WQEについてリモート・ノードによって肯定応答が返される次のバイトを識別する。「次の受信ACK TCPセグメントのコミット済み連続番号」は、リモート・ノードから受け取られることが予想される、現在の受信ACK WQE中へのコミット済みバイト・オフセットに関連付けられた次のTCP連続番号を識別する。「次の受信ACK TCPセグメントのコミット済みウィンドウ値」は、リモート・ノードから送り返された最後のTCPウィンドウのサイズの結果生じ、現在の受信ACK WQE中へのコミット済みバイト・オフセットに関連付けられた、送信ウィンドウの変化を識別する。

20

30

【0117】

実効RQスイッチオーバー/スイッチバック・コンテキスト1232は、次の情報を含み、この情報はQPが作成される時にゼロに設定される。

40

【0118】

【表 3】

現在の実効受信WQE番号
現在の受信WQE中への実効バイト・オフセット
次の受信TCPセグメントの実効連続番号
次の受信TCPセグメントの実効ウィンドウ
現在の実効送信AckWQE番号
現在の送信AckWQE中への実効バイト・オフセット
次の送信AckTCPセグメントの実効連続番号
次の送信AckTCPセグメントの実効ウィンドウ値

10

【0119】

すべての実効RQスイッチオーバー/スイッチバック・フィールドは、代替RNICによってチェックポイントACKを通じてチェックポイントされていない情報を含む。主要RNICと代替RNICはともに、これらのフィールドをすべて含む。「現在の実効受信WQE番号」は、現在受信待ち行列で処理中のWQEを識別する。「現在の受信WQE中への実効バイト・オフセット」は、現在の実効受信WQEについて処理される次のバイトを識別する。「次の受信TCPセグメントの実効連続番号」は、現在の受信WQE中への実効バイト・オフセットに予想されるTCP連続番号を識別する。「次の受信TCPセグメントの実効ウィンドウ値」は、現在の受信WQE中への実効バイト・オフセットから受信が開始される次のTCPセグメントに利用可能なTCPウィンドウ・サイズを識別する。「現在の実効送信ACK WQE番号」は、肯定応答がリモート・ノードに送信された、現在受信待ち行列で処理されているWQEを識別する。「現在の送信ACK WQEへの実効バイト・オフセット」は、現在の実効ACK WQEについてリモート・ノードに対して肯定応答が返される次のバイトを識別する。「次の送信ACK TCPセグメントの実効連続番号」は、リモート・ノードに送信されることになる、現在の送信ACK WQE中への実効バイト・オフセットに関連付けられた次のTCP連続番号を識別する。「次の送信ACK TCPセグメントの実効ウィンドウ値」は、リモート・ノードに送信されることになる、現在の送信ACK WQE中への実効バイト・オフセットに関連付けられた次のTCPウィンドウのサイズを識別する。

20

30

【0120】

コミット済みRQスイッチオーバー/スイッチバック・コンテキスト1236は、次の情報を含み、この情報は、QPが作成される際にゼロに設定される。

【0121】

【表4】

コミット済み現在の受信WQE番号
現在の受信WQE中へのコミット済みバイト・オフセット
次の受信TCPセグメントのコミット済み連続番号
次の受信TCPセグメントのコミット済みウィンドウ値
コミット済み現在の送信AckWQE番号
現在の送信AckWQE中へのコミット済みバイト・オフセット
次の送信AckTCPセグメントのコミット済み連続番号
次の送信AckTCPセグメントのコミット済みウィンドウ値

10

【0122】

すべてのコミット済みRQスイッチオーバー/スイッチバック・フィールドは、チェックポイントACKを通じて代替RNICによってチェックポイントされていない情報を含む。主要RNICと代替RNICはともに、これらのフィールドをすべて含む。「コミット済み現在の受信WQE番号」は、チェックポイントされ、現在受信待ち行列で処理中の最も新しいWQEを識別する。「現在の受信WQEへのコミット済みバイト・オフセット」は、コミット済み現在の受信WQEについて処理される次のバイトを識別する。「次の受信TCPセグメントのコミット済み連続番号」は、現在の受信WQE中へのコミット済みバイト・オフセットに予想されるTCP連続番号を識別する。「次の受信TCPセグメントのコミット済みウィンドウ値」は、現在の受信WQE中へのコミット済みバイト・オフセットから受信が開始される次のTCPセグメントに利用できるTCPウィンドウのサイズを識別する。「コミット済み現在の送信ACK WQE番号」は、肯定応答がリモート・ノードに送信され、受信待ち行列でチェックポイントされた最も新しいWQEを識別する。「現在の送信ACK WQE中へのコミット済みバイト・オフセット」は、コミット済み現在の送信ACK WQEについて、リモート・ノードに対して肯定応答すべき次のバイトを識別する。「次の送信ACK TCPセグメントのコミット済み連続番号」は、リモート・ノードに送信されることになる、現在の送信ACK WQE中へのコミット済みバイト・オフセットに関連付けられた次のTCP連続番号を識別する。「次の送信ACK TCPセグメントのコミット済みウィンドウ値」は、リモート・ノードに送信すべき、現在の送信ACK WQE中へのコミット済みバイト・オフセットに関連付けられた次のTCPウィンドウのサイズを識別する。

20

30

【0123】

主要RNICおよび代替RNIC識別コンテキスト1240は、次の情報を含む。

【0124】

【表5】

RNICのスイッチオーバー/スイッチバック状態
主要RNICの識別子
主要RNICのポート識別子
代替RNICの識別子
代替RNICのポート識別子

40

【0125】

QPが作成される際に、RNICスイッチオーバー/スイッチバック(S/S)状態が、QP作成verbの入力修飾子として渡される。RNIC S/S状態がゼロの場合、

50

QPは、S/Sが無効にされ、主要フィールドだけが有効になる。RNICのS/S状態が非ゼロの場合は、QPはS/Sが使用可能にされ、主要フィールドと代替フィールドの両方が有効になる。

【0126】

主要RNIC識別子は、主要RNICを一意に識別するために使用される。例えば、RNIC識別子は、RNICのEIU-64値とすることができる。主要RNICポートの識別子を使用して、主要RNICでQPが関連付けられたポートを一意に識別する。例えば、RNICポート識別子は、ポート番号、またはMACアドレスとIPアドレスの組み合わせとすることができる。代替RNIC識別子は、代替RNICを一意に識別するために使用される。代替RNICポート識別子を使用して、代替RNICでQPが関連付けら

10

【0127】

QP1176などの各RNIC CQは、CQCT1252などのCQコンテキストテーブル(CQCT)に、CQCE0 1256、CQCE1 1260、およびCQCE N 1264などのCWコンテキスト・エントリを有する。CQCTは、CQCTアドレス1288などのCQCTのアドレスと、CQCT長1292などのCQCTの長さを含むCQCTレジスタ1284などのCQコンテキスト・テーブル・レジスタを通じてアクセスされる。QPコンテキスト・テーブルの各エントリは、固定サイズであるが、CQCTエントリのサイズが可変のCQCTも使用することができる。CQCTの各エントリは、そのエントリについて参照されるCQ番号に関連付けられた関連するコンテキスト情報

20

【0128】

従来のCQコンテキスト状態1268は、CQの状態とCQ中の総エントリ数など、周

30

【0129】

実効CQスイッチオーバー/スイッチバック・コンテキスト1272は、次の情報を含み、この情報はQPが作成される時にゼロに設定される。

【0130】

【表6】

現在の実効完了CQE番号

40

【0131】

実効CQスイッチオーバー/スイッチバック・フィールドは、代替RNICによりチェックポイントACKを通じてチェックポイントされていない情報を含む。主要RNICと代替RNICとともに、これらのフィールドをすべて含む。「現在の実効完了CQE番号」は、現在完了待ち行列で処理されているCQEを識別する。

【0132】

「コミット済みCQスイッチオーバー/スイッチバック・コンテキスト1272」は、次の情報を含み、この情報はQPが作成される際にゼロに設定される。

【0133】

【表7】

コミット済み現在の完了CQE番号

【0134】

コミット済みCQスイッチオーバー/スイッチバック・フィールドは、代替RNICによりチェックポイントACKを通じてチェックポイントされた情報を含む。主要RNICと代替RNICはともに、これらのフィールドをすべて含む。「現在の実効完了CQE番号」は、すでにチェックポイントされ、現在完了待ち行列で処理中のCQEを識別する。

10

【0135】

主要RNICおよび代替RNIC識別コンテキスト1280は、次の情報を含む。

【0136】

【表8】

RNICのスイッチオーバー/スイッチバック状態
主要RNICの識別子
主要RNICのポート識別子
代替RNICの識別子
代替RNICのポート識別子

20

【0137】

QPが作成される際に、RNICスイッチオーバー/スイッチバック(S/S)状態がQP作成verbの入力修飾子として渡される。RNIC S/S状態がゼロの場合は、CQは、S/Sが使用不可にされ、主要フィールドのみが有効になる。この場合、S/Sが使用不可にされたRNICを有するQPだけをそのCQに関連付けることができる。

【0138】

RNICのS/S状態が非ゼロの場合は、CQはS/Sが使用可能にされ、主要フィールドと代替フィールドの両方が有効になる。この場合、一致する主要RNICと代替RNICを有するQPのみをCQに関連付けることができる。

30

【0139】

主要RNIC識別子は、主要RNICを一意に識別するために使用される。例えば、RNIC識別子は、RNICのEIU-64値とすることができる。主要RNICポート識別子は、主要RNICでCQが関連付けられたポートを一意に識別するために使用される。例えば、RNICポート識別子は、ポート番号、またはMACアドレスとIPアドレスの組み合わせとすることができる。代替RNIC識別子は、代替RNICを一意に識別するために使用される。代替RNIC識別子は、代替RNICでCQが関連付けられたポートを一意に識別するために使用される。

40

【0140】

図15に、本発明の好ましい実施形態による、2つのRNIC間の接続の例示的な初期化プロセスを説明するフローチャートを提供する。図15に示すように、この動作は、RNICがスイッチオーバー/スイッチバックをサポートするかどうかを判定するためのRNICの照会から開始する(ステップ1302)。このステップでは、主要RNICに選択されることになるRNICと、代替RNICに選択されることになるRNICの両方に照会する。S/Sをサポートし、ホストがS/S RNICとして設定したいと望む各RNICについて、ホストは、RNIC変更verbを発行する(ステップ1304)。RNIC変更verbは、S/Sと非S/Sに対して設定される、QP番号の範囲、CQ番号の範囲、およびメモリ変換および保護テーブル・エントリの範囲を選択する。

50

【0141】

作成された各CQに対して、ホストは、CQ作成verbを発行する(ステップ1306)。CQ作成verbは、そのCQに関連付けられた、主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する。メモリ領域登録verbの1つを通じて登録された各メモリ領域に対して、ホストは、そのメモリ領域に関連付けられた主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する(ステップ1308)。

【0142】

作成された各QPに対して、ホストはQP作成verbを発行する(ステップ1310)。QP作成verbは、そのQPに関連付けられた、RNIC状態(主要または代替) 10、主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する。

【0143】

メモリ・ウィンドウ割り振りverbを通じて割り振られた各メモリ・ウィンドウについて(ステップ1312)、ホストは、そのメモリ・ウィンドウに関連付けられた主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する(ステップ1312)。この結果、2つのRNIC、すなわち主要RNICと代替RNICの間の接続が初期化される。

【0144】

図16は、本発明の好ましい実施形態による2つのRNIC間の接続の例示的な再同期プロセスのフローチャートである。この再同期プロセスを使用して、S/Sをサポートするように構成されているが、代替RNICなしで動作している主要RNICで、S/S QPとして設定されているQPに代替RNICを割り当てる。 20

【0145】

図16に示すように、この動作は、S/Sと非S/Sに対して設定される、QP番号の範囲、CQ番号の範囲、およびメモリ変換および保護テーブルのエントリの範囲を決定するために主要RNICに照会することから開始する(ステップ1320)。代替RNICにも照会して、スイッチオーバー/スイッチバックをサポートするかどうかを判定する(これもステップ1320の一部)。サポートする場合、動作はステップ1322に進み、サポートしない場合は動作が終了される。 30

【0146】

ステップ1322で、ホストは、代替RNICに対してRNIC変更verbを発行し、S/Sと非S/Sに対して設定されるQP番号の範囲、CQ番号の範囲、およびメモリ変換および保護テーブル・エントリの範囲を選択する(ステップ1322)。代替RNICで作成される各CQに、ホストは、CQ作成verbを発行する(ステップ1324)。CQ作成verbは、そのCQに関連付けられた主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する。

【0147】

主要RNICと代替RNICとに共有されることになる各メモリ領域について、ホストは、共有メモリ登録verbを使用して、代替RNICのすべての共有メモリ領域を登録する(ステップ1326)。メモリ領域登録verbの1つを通じてメモリ領域が登録されると、ホストは、そのメモリ領域に関連付けられる主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する。 40

【0148】

代替RNICで作成される各QPに対して、ホストは、QP作成verbを発行する(ステップ1328)。QP作成verbは、そのQPに関連付けられるRNICの状態(「代替」)、主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する。メモリ・ウィンドウ割り振りverbを通じて割り振られる各メモリ・ウィンドウについて、ホストは、そのメモリ・ウィンドウに関連付けられる主要RNIC、主要RNICポート、代替RNIC、および代替RNICポートを指定する(ステップ 50

1330)。

【0149】

主要RNICで再同期すべき各QPについて、ホストは、QP変更verbを発行して、QPを停止させ、そのQPに、RNIC状態(「主要」)、主要RNICおよびポート、代替RNICおよびポートを割り当てる(ステップ1332)。主要RNICで再同期される各CQについて、そのCQに割り当てられたすべてのQPが停止されると、ホストは、CQ変更verbを発行してそのCQを停止させ、そのCQに、RNICの状態(「主要」)、主要RNICおよびポート、代替RNICおよびポートを割り当てる(ステップ1334)。

【0150】

主要RNICで再同期される各CQに対して、ホストは、CQ変更verbを発行してCQを再度アクティブ化する(ステップ1336)。主要RNICで再同期されることが必要な各QPに対して、ホストは、QP変更verbを発行してそのQPを再度アクティブ化する(ステップ1338)。

【0151】

図17に、本発明の好ましい実施形態による、主要RNICと代替RNIC間で使用される例示的なチェックポイント・メッセージを示す概略図を提供する。図13から、主要RNIC1100と代替RNIC1104は、チェックポイント・メッセージ1106を使用してRNICの状態とQPの状態を伝える。チェックポイント・メッセージは、主要RNICと代替RNICを相互に接続するファブリックを通じて直接通信することができる。チェックポイント・メッセージは、例えば、主要RNICと代替RNICの両方からアクセス可能な共有システム・メモリ領域を通じて間接的に通信することもできる。

【0152】

図13に、主要RNICと代替RNICを相互に接続するファブリックを直接通じてチェックポイント・メッセージが送信される例示的事例を示す。このファブリックは、ローカル・エリア・ネットワークでよく、メッセージは、1つまたは複数のTCP接続を通じて送信することができる。ファブリックは、PCI、PCI-X、PCI-Express等のメモリ・マップのI/O拡張ネットワークであってもよい。最後に、ファブリックは、InfiniBandなどのシステム・エリア・ネットワークとすることもできる。図13には、ファブリックがローカル・エリア・ネットワークである例示的事例を示す。

【0153】

図17に、チェックポイント・メッセージ1340など、すべてのチェックポイント・メッセージに共通のフィールドを示す。すべてのチェックポイント・メッセージの最初のフィールドは、OpCode1344などの命令コードであり、これは、そのメッセージに含まれるチェックポイント情報のタイプを記述する。メッセージの次のフィールドは、長さ1348などの長さであり、これはメッセージ長をバイト単位で記述する。メッセージの最後のフィールドは、CRC32 1352などの周知のiSCSIの32ビットの巡回冗長検査(CRC32)であり、メッセージの有効性の検証に使用される。各メッセージは、長さ1348とCRC32 1352の間で搬送されるフィールドを定義する。OpCodeフィールドと長さフィールドのサイズは、図より小さくとも大きくともよい。使用されるCRCの多項式は、図と異なってよい。

【0154】

次の表は、チェックポイント・プロセスで使用されるメッセージ・タイプを定義し、その後の項では、チェックポイント・メッセージの内容を含めて、チェックポイントのプロセスを説明する。各OpCodeと長さの値は、本発明の主旨から逸脱せずに、下記の表に示す値と異なってよい。

【0155】

10

20

30

40

【表 9】

OpCode	長さ	メッセージ
x0000	x0000	I am alive
x0001	x0004	予備
x0002	x0004	実効SQ S/S送信コンテキスト更新
x0003	x0004	コミット済みSQ S/S送信コンテキスト更新
x0004	x0004	実効SQ S/SAckコンテキスト更新
x0005	x0004	コミット済みSQ S/SAckコンテキスト更新
x0006	x0004	実効RQ S/S受信コンテキスト更新
x0007	x0004	コミット済みRQ S/S受信コンテキスト更新
x0008	x0004	実効RQ S/SAckコンテキスト更新
x0009	x0004	コミット済みRQ S/SAckコンテキスト更新
x000A	x0004	実効CQ S/Sコンテキスト更新
x000B	x0004	コミット済みCQ S/Sコンテキスト更新

10

【 0 1 5 6 】

20

図 1 8 に、本発明の好ましい実施形態による、2つの R N I C によって使用される例示的な送信チェックポイント・メッセージの流れとプロセスを説明する概略図を提供する。R N I C は、Q P W Q E 処理アルゴリズムを使用して、送信される T C P セグメントに関連付けられた Q P コンテキストを調べる。以下は、Q P コンテキスト 1 4 1 2 の S Q 1 4 0 8 から送信された送信 T C P セグメントをチェックポイントするために使用されるチェックポイント・メッセージの流れである。

【 0 1 5 7 】

S Q 1 4 0 8 が、アウトバウンドのスケジューリングを行える状態の T C P セグメントを有すると仮定する。主要 R N I C 1 4 0 4 は、保持している実効 S Q S / S コンテキスト 1 2 2 4 のローカル・コピーの次の 4 フィールドを更新する。

30

【 0 1 5 8 】

【表 1 0】

現在の実効送信WQE番号
現在の送信WQE中への実効バイト・オフセット
次の送信TCPセグメントの実効連続番号
次の送信TCPセグメントの値についての実効ウィンドウ

40

【 0 1 5 9 】

上記の 4 つのフィールドがローカル Q P コンテキストで更新されると、主要 R N I C 1 4 0 4 は、実効 S Q S / S コンテキスト更新送信メッセージ 1 4 1 6 を代替 R N I C 1 4 8 0 に送信する。O p C o d e 1 4 5 4 は、x 0 0 0 2、すなわち実効 S Q S / S 更新送信メッセージ 1 4 1 6 に設定される。長さフィールドは、x 0 0 0 4、すなわち実効 S Q S / S コンテキスト更新送信メッセージの長さに設定される。長さに続く 5 つのフィールドは以下のように設定される。

【 0 1 6 0 】

【表 1 1】

QP番号
現在の実効送信WQE番号
現在の送信WQE中への実効バイト・オフセット
次の送信TCPセグメントの実効連続番号
次の送信TCPセグメントの値の実効ウィンドウ

10

【0 1 6 1】

代替RNIC1480は、実効SQ S / Sコンテキスト更新送信メッセージ1416を受信し、その有効性を検証する。実効SQ S / Sコンテキスト更新送信メッセージ1416が有効である、例えばCRC32 1458が有効である場合、代替RNIC1480は、保持している実効SQ S / Sコンテキストのローカル・コピーを、実効SQ S / Sコンテキスト送信メッセージ1416の内容で更新する。代替RNIC1480における実効SQ S / Sコンテキスト1224のローカル・コピーの更新では、実効SQ S / Sコンテキスト1224の次の4フィールドを更新する。

【0 1 6 2】

【表 1 2】

20

現在の実効送信WQE番号
現在の送信WQE中への実効バイト・オフセット
次の送信TCPセグメントの実効連続番号
次の送信TCPセグメントの値の実効ウィンドウ

【0 1 6 3】

代替RNIC1480は次いで、保持している実効SQ S / Sコンテキスト1224のローカル・コピーを、コミット済みSQ S / Sコンテキスト1228のローカル・コピーにコピーする。以下に示すコミット済みSQ S / Sコンテキスト1228の4つのフィールドが、代替RNIC1480で更新される。

30

【0 1 6 4】

【表 1 3】

コミット済み現在の送信WQE番号
現在の送信WQE中へのコミット済みバイト・オフセット
次の送信TCPセグメントのコミット済み連続番号
次の送信TCPセグメントの値についてのコミット済みウィンドウ

40

【0 1 6 5】

上記の4フィールドがローカルQPコンテキストで更新されると、代替RNIC1480は、コミット済みSQ S / Sコンテキスト更新送信メッセージ1424を主要RNIC1404に送信する。OpCode1454は、x0003、すなわちコミット済みSQ S / Sコンテキスト更新送信メッセージ1424に設定される。長さフィールドは、x0004、すなわちコミット済みSQ S / Sコンテキスト更新送信メッセージの長さに設定される。長さに続く5つのフィールドは以下である。

【0 1 6 6】

【表 1 4】

QP番号
コミット済み現在の送信WQE番号
現在の送信WQE中へのコミット済みバイト・オフセット
次の送信TCPセグメントのコミット済み連続番号
次の送信TCPセグメントの値のコミット済みウィンドウ

10

【0 1 6 7】

主要 R N I C 1 4 0 4 は、コミット済み S Q S / S コンテキスト更新送信メッセージ 1 4 2 4 を受信し、有効性を検証する。コミット済み S Q S / S コンテキスト送信メッセージ 1 4 2 4 が有効である場合、例えば C R C 3 2 が有効である場合、主要 R N I C 1 4 0 4 は、保持しているコミット済み S Q S / S コンテキストのローカル・コピーを、コミット済み S Q S / S コンテキスト・メッセージ 1 4 2 4 の内容で更新する。すなわち、保持しているコミット済み S Q S / S コンテキスト 1 2 2 8 の次の 4 フィールドを更新する。

【0 1 6 8】

【表 1 5】

20

コミット済み現在の送信WQE番号
現在の送信WQE中へのコミット済みバイト・オフセット
次の送信WQEセグメントのコミット済み連続番号
次の送信TCPセグメントの値のコミット済みウィンドウ

【0 1 6 9】

S Q 1 4 0 8 は、T C P セグメントをアウトバウンド・スケジューラに置き、スケジューラは、次回 S Q 1 4 0 8 から T C P セグメントを送信する際にその T C P セグメントを送信する 1 4 0 0。

30

【0 1 7 0】

コミット済み S Q S / S コンテキスト更新送信メッセージ 1 4 2 4 が無効であるか、チェックポイントの期限内に受信されない場合、主要 R N I C 1 4 0 4 は、最高でチェックポイント再試行回数によって定義される回数までそれを再送信する。代替 R N I C 1 4 8 0 が、重複した実効 S Q S / S コンテキスト更新送信メッセージ 1 4 1 6 を受信した場合、代替 R N I C 1 4 8 0 は、1 つ前のコミット済み S Q S / S コンテキスト更新送信メッセージ 1 4 2 4 を再送信する。この時点で、アウトバウンドの S e n d は、チェックポイントされている。

【0 1 7 1】

40

次いで、T C P A C K セグメントが主要 R N I C 1 4 0 4 によって受信される際に使用され、Q P コンテキスト 1 4 1 2 の S Q 1 4 0 8 に関連付けられたチェックポイント・メッセージの流れを説明する。まず、着信する T C P A C K セグメントが、中間にある速度一致バッファに受信される。T C P A C K セグメントが有効でない場合は、破棄される。T C P A C K セグメントが有効である場合、主要 R N I C 1 4 0 4 は、保持する実効 S Q S / S コンテキスト 1 2 2 4 のローカル・コピーの以下の 4 フィールドを更新する。

【0 1 7 2】

【表 16】

現在の実効受信AckWQE番号
現在の受信AckWQE中への実効バイト・オフセット
次のRcvAckTCPセグメントの実効連続番号
次のRcvAckTCPセグメントの実効ウィンドウ値

【0173】

上記の4フィールドがローカルのQPコンテキストで更新されると、主要RNIC1404は、実効SQ S/Sコンテキスト更新Ackメッセージ1432を代替RNIC1480に送信する。OpCode1454は、x0004、すなわち実効SQ S/Sコンテキスト更新Ackメッセージ1432に設定される。長さフィールドは、x0004、すなわち実効SQ S/Sコンテキスト更新Ackメッセージの長さに設定される。長さに続く5つのフィールドは以下のように設定される。

【0174】

【表 17】

QP番号	20
現在の実効受信AckWQE番号	
現在の受信AckWQE中への実効バイト・オフセット	
次のRcvAckTCPセグメントの実効連続番号	
次のRcvAckTCPセグメントの実効ウィンドウ値	

【0175】

代替RNIC1480は、実効SQ S/Sコンテキスト更新Ackメッセージ1432を受信し、有効性を検証する。実効SQ S/Sコンテキスト・メッセージ1416が有効である場合、例えばCRC32が有効である場合、代替RNIC1480は、保持する実効SQ S/Sコンテキスト1224のローカル・コピーを、実効SQ S/Sコンテキスト・メッセージ1432の内容で更新する。代替RNIC1480における実効SQ S/Sコンテキスト1224のローカル・コピーの更新では、実効SQ S/Sコンテキスト1224の以下の4フィールドが更新される。

【0176】

【表 18】

現在の実効受信AckWQE番号	40
現在の受信AckWQE中への実効バイト・オフセット	
次のRcvAckTCPセグメントの実効連続番号	
次のRcvAckTCPセグメントの実効ウィンドウ値	

【0177】

代替RNIC1480は次いで、自身の実効SQ S/Sコンテキスト1224のローカル・コピーを、コミット済みSQ S/Sコンテキスト1228のローカル・コピーにコピーする。コミット済みSQ S/Sコンテキスト1228の以下に示す4つのフィールドが、代替RNIC1480で更新される。

【0178】

【表 19】

コミット済み現在の受信AckWQE番号
現在の受信AckWQE中へのコミット済みバイト・オフセット
次のRcvAckTCPセグメントのコミット済み連続番号
次のRcvAckTCPセグメントのコミット済みウィンドウ値

【0179】

上記の4フィールドがローカルQPコンテキストで更新されると、代替RNIC1480は、コミット済みSQ S / Sコンテキスト更新Ackメッセージ1436を主要RNIC1404に送信する。OpCode1454は、x0005、すなわちコミット済みSQ S / Sコンテキスト更新Ackメッセージ1436に設定される。長さフィールドは、x0004、すなわちコミット済みSQ S / Sコンテキスト更新Ackメッセージの長さに設定される。長さに続く5つのフィールドは以下のように設定される。

【0180】

【表 20】

QP番号
コミット済み現在の受信AckWQE番号
現在の受信AckWQE中へのコミット済みバイト・オフセット
次のRcvAckTCPセグメントのコミット済み連続番号
次のRcvAckTCPセグメントのコミット済みウィンドウ値

【0181】

主要RNIC1404は、コミット済みSQ S / Sコンテキスト更新Ackメッセージ1436を受信し、有効性を検証する。コミット済みSQ S / SコンテキストAckメッセージ1436が有効である、例えばCRC32が有効である場合、主要RNIC1404は、保持するコミット済みSQ S / Sコンテキストのローカル・コピーを、コミット済みSQ S / SコンテキストAckメッセージ1436の内容で更新する。すなわち、保持する実効SQ S / Sコンテキストの次の4フィールドを更新する。

【0182】

【表 21】

コミット済み現在の受信AckWQE番号
現在の受信AckWQE中へのコミット済みバイト・オフセット
次のRcvAckTCPセグメントのコミット済み連続番号
次のRcvAckTCPセグメントのコミット済みウィンドウ値

【0183】

コミット済みSQ S / Sコンテキスト更新Ackメッセージ1436が無効であるか、チェックポイント期限内に受信されない場合、主要RNIC1404は、最高でチェックポイント再試行回数によって定義される回数までそれを再送信する。代替RNIC1480が、重複した実効SQ S / Sコンテキスト更新Ackメッセージ1432を受信した場合、代替RNIC1480は、1つ前のコミット済みSQ S / Sコンテキスト更新Ackメッセージ1436を再送信する。この時点で、着信するAckは、チェックポイントされている。

10

20

30

40

50

【0184】

図19に、本発明の好ましい実施形態による、2つのRNICによって使用される例示的な受信チェックポイント・メッセージの流れとプロセスを説明する概略図を提供する。RNICは、周知のTCP/IPルックアップ・アルゴリズムを使用することにより、着信TCPセグメントに関連付けられたQPコンテキストを検索する。着信TCPセグメントが受信されると、TCP/IPクインタプル(quintuple)検索を使用して、その着信TCPセグメントに関連付けられたQPを判定する。

【0185】

以下は、着信TCPセグメントがQPコンテキスト1512のRQ1508を対象とする場合に使用される例示的なチェックポイント・メッセージの流れである。着信TCPセグメントは、中間にある速度一致バッファに受信される。そのTCPセグメントが有効である場合は、最終的な宛先に置かれ、有効でない場合は破棄される。

10

【0186】

TCPセグメントが最終的な宛先に置かれると、主要RNIC1504は、保持している実効RQS/Sコンテキスト1232のローカル・コピーの以下の4フィールドを更新する。

【0187】

【表22】

現在の実効受信WQE番号
現在の受信WQE中への実効バイト・オフセット
次の受信TCPセグメントの実効連続番号
次の受信TCPセグメントの実効ウィンドウ値

20

【0188】

上記の4フィールドがローカルQPコンテキストで更新されると、主要RNIC1504は、実効RQS/Sコンテキスト更新受信メッセージ1516を代替RNIC1580に送信する。OpCode1554は、x0006、すなわち実効RQS/Sコンテキスト更新受信メッセージ1516に設定される。長さフィールドは、x0004、すな

30

【0189】

【表23】

QP番号
現在の実効受信WQE番号
現在の受信WQE中への実効バイト・オフセット
次の受信TCPセグメントの実効連続番号
次の受信TCPセグメントの実効ウィンドウ値

40

【0190】

代替RNIC1580は、実効RQS/Sコンテキスト更新受信メッセージ1516を受信し、有効性を検証する。実効RQS/Sコンテキスト受信メッセージ1516が有効である場合、例えばCRC321562が有効である場合、代替RNIC1580は、保持する実効RQS/Sコンテキスト1232のローカル・コピーを、実効RQS/Sコンテキスト受信メッセージ1516の内容で更新する。代替RNIC1580における実効SQS/Sコンテキスト1232のローカル・コピーの更新では、実効SQS/Sコンテキスト1232の以下の4フィールドが更新される。

50

【 0 1 9 1 】

【表 2 4】

現在の実効受信WQE番号
現在の受信WQE中への実効バイト・オフセット
次の受信TCPセグメントの実効連続番号
次の受信TCPセグメントの実効ウィンドウ値

【 0 1 9 2 】

10

代替RNIC1580は次いで、保持する実効RQ S/Sコンテキスト1232のローカル・コピーを、コミット済みRQ S/Sコンテキスト1236のローカル・コピーにコピーする。コミット済みRQ S/Sコンテキスト1236の以下に示す4つのフィールドが、代替RNIC1580で更新される。

【 0 1 9 3 】

【表 2 5】

コミット済み現在の受信WQE番号
現在の受信WQE中へのコミット済みバイト・オフセット
次の受信TCPセグメントのコミット済み連続番号
次の受信TCPセグメントのコミット済みウィンドウ値

20

【 0 1 9 4 】

上記の4フィールドがローカルQPコンテキストで更新されると、代替RNIC1580は、コミット済みRQ S/Sコンテキスト更新受信メッセージ1524を主要RNIC1504に送信する。OpCode1554は、x0007、すなわちコミット済みRQ S/Sコンテキスト更新受信メッセージ1524に設定される。長さフィールドは、x0004、すなわちコミット済みRQ S/Sコンテキスト更新受信メッセージの長さに設定される。長さに続く5つのフィールドは以下のように設定される。

30

【 0 1 9 5 】

【表 2 6】

QP番号
コミット済み現在の受信WQE番号
現在の受信WQE中へのコミット済みバイト・オフセット
次の受信TCPセグメントのコミット済み連続番号
次の受信TCPセグメントのコミット済みウィンドウ値

40

【 0 1 9 6 】

主要RNIC1504は、コミット済みRQ S/Sコンテキスト更新受信メッセージ1524を受信し、有効性を検証する。コミット済みRQ S/Sコンテキスト受信メッセージ1524が有効である場合、例えばCRC32が有効である場合、主要RNIC1504は、保持するコミット済みRQ S/Sコンテキスト1236のローカル・コピーを、コミット済みRQ S/Sコンテキスト受信メッセージ1524の内容で更新する。すなわち、それが保有するコミット済みRQ S/Sコンテキスト1236の次の4フィールドを更新する。

【 0 1 9 7 】

50

【表 2 7】

コミット済み現在の受信WQE番号
現在の受信WQE中へのコミット済みバイト・オフセット
次の受信TCPセグメントのコミット済み連続番号
次の受信TCPセグメントのコミット済みウィンドウ値

【 0 1 9 8 】

コミット済みRQ S / Sコンテキスト更新受信メッセージ1524が無効であるか、
 チェックポイントの期限内に受信されない場合、主要RNIC1504は、最高でチェ
 ックポイント再試行回数によって定義される回数までそれを再送信する。代替RNIC15
 80が、重複した実効RQ S / Sコンテキスト更新受信メッセージ1516を受信した
 場合、代替RNIC1580は、1つ前のコミット済みRQ S / Sコンテキスト更新受
 信メッセージ1524を再送信する。この時点で、着信するSendは、チェックポイント
 されている。

10

【 0 1 9 9 】

以下は、主要RNIC1504からTCP ACKセグメントが送信され、QPコンテ
 クスト1512のRQ1508に関連付けられる場合に使用される例示的なチェックポ
 イント・メッセージの流れである。主要RNIC1504は、保持している実効RQ S /
 Sコンテキスト1232のローカル・コピーの以下の4フィールドを更新する。

20

【 0 2 0 0 】

【表 2 8】

現在の実効送信AckWQE番号
現在の送信AckWQE中への実効バイト・オフセット
次の送信AckTCPセグメントの実効連続番号
次のAckTCPセグメントの実効ウィンドウ値

30

【 0 2 0 1 】

上記の4フィールドがローカルQPコンテキストで更新されると、主要RNIC150
 4は、実効RQ S / Sコンテキスト更新Ackメッセージ1532を代替RNIC15
 80に送信する。OpCode1554は、x0008、すなわち実効RQ S / Sコン
 テキスト更新Ackメッセージ1532に設定される。長さフィールドは、x0004、
 すなわち実効RQ S / Sコンテキスト更新Ackメッセージの長さに設定される。長さ
 に続く5つのフィールドは以下のように設定される。

【 0 2 0 2 】

【表 2 9】

QP番号
現在の実効送信AckWQE番号
現在の送信AckWQE中への実効バイト・オフセット
次の送信AckTCPセグメントの実効連続番号
次のAckTCPセグメントの実効ウィンドウ値

40

【 0 2 0 3 】

代替RNIC1580は、実効RQ S / Sコンテキスト更新Ackメッセージ153
 2を受信し、有効性を検証する。実効RQ S / SコンテキストAckメッセージ153

50

2 が有効である場合、例えばCRC32が有効である場合、代替RNIC1580は、保持している実効RQ S / Sコンテキスト1232のローカル・コピーを、実効RQ S / SコンテキストAckメッセージ1532の内容で更新する。代替RNIC1580における実効RQ S / Sコンテキスト1232のローカル・コピーの更新では、実効RQ S / Sコンテキスト1232の以下の4フィールドが更新される。

【0204】

【表30】

現在の実効送信AckWQE番号	
現在の送信AckWQE中への実効バイト・オフセット	10
次の送信AckTCPセグメントの実効連続番号	
次のAckTCPセグメントの実効ウィンドウ値	

【0205】

代替RNIC1580は次いで、保持する実効RQ S / Sコンテキスト1232のローカル・コピーを、コミット済みRQ S / Sコンテキスト1236のローカル・コピーにコピーする。コミット済みRQ S / Sコンテキスト1236の以下に示す4フィールドが、代替RNIC1580で更新される。

【0206】

【表31】

コミット済み現在の送信AckWQE番号	
現在の送信AckWQE中へのコミット済みバイト・オフセット	
次の送信AckTCPセグメントのコミット済み連続番号	
次のAckTCPセグメントのコミット済みウィンドウ値	

【0207】

上記の4フィールドがローカルQPコンテキストで更新されると、代替RNIC1580は、コミット済みRQ S / Sコンテキスト更新Ackメッセージ1536を主要RNIC1504に送信する。OpCode1454は、x0009、すなわちコミット済みRQ S / Sコンテキスト更新Ackメッセージ1536に設定される。長さフィールドは、x0004、すなわちコミット済みRQ S / Sコンテキスト更新Ackメッセージの長さに設定される。長さに続く5つのフィールドは以下のように設定される。

【0208】

【表32】

QP番号	
コミット済み現在の送信AckWQE番号	
現在の送信AckWQE中へのコミット済みバイト・オフセット	
次の送信AckTCPセグメントのコミット済み連続番号	
次のAckTCPセグメントのコミット済みウィンドウ値	

【0209】

主要RNIC1504は、コミット済みRQ S / Sコンテキスト更新Ackメッセージ1536を受信し、有効性を検証する。コミット済みRQ S / SコンテキストAckメッセージ1536が有効である、例えばCRC32が有効である場合、主要RNIC1

10

20

30

40

50

504は、保持するコミット済みRQ S / Sコンテキスト1236のローカル・コピーを、コミット済みRQ S / SコンテキストAckメッセージ1536の内容で更新する。すなわち、保持するコミット済みRQ S / Sコンテキストの以下の4フィールドを更新する。

【0210】

【表33】

コミット済み現在の送信AckWQE番号
現在の送信AckWQE中へのコミット済みバイト・オフセット
次の送信AckTCPセグメントのコミット済み連続番号
次のAckTCPセグメントのコミット済みウィンドウ値

10

【0211】

コミット済みRQ S / Sコンテキスト更新Ackメッセージ1536が無効であるか、チェックポイントの期限内に受信されない場合、主要RNIC1504は、最高でチェックポイント再試行回数によって定義される回数までそれを再送信する。代替RNIC1580が、重複した実効RQ S / Sコンテキスト更新Ackメッセージ1532を受信した場合、代替RNIC1580は、1つ前のコミット済みRQ S / Sコンテキスト更新Ackメッセージ1536を再送信する。この時点で、着信するAckは、チェックポイントされている。

20

【0212】

図20に、本発明の好ましい実施形態による、2つのRNICに使用される例示的な完了チェックポイント作業の流れとプロセスを説明する概略図を提供する。RNICは、周知のQP WQE処理アルゴリズムを使用することにより、完了した作業待ち行列要素に関連付けられたCQコンテキストを検索し、周知のCQ CQE処理アルゴリズムを使用してCQEを作成する。

【0213】

以下は、RNIC1604によりホストにサーフェス (surface) されたCQ1608に付加された完了待ち行列要素をチェックポイントするために使用されるチェックポイント・メッセージの流れである。CQ1608が、CQ1608に挿入可能なCQEを有するとする。CQ1608は、そのCQEをCQ1608に挿入し、保持している実効CQ S / Sコンテキスト1272のローカル・コピーの以下のフィールドを更新する。

30

【0214】

【表34】

現在の実効完了CQE番号

【0215】

上記のフィールドがローカルのCQコンテキストで更新されると、主要RNIC1604は、代替RNIC1680に、実効CQ S / Sコンテキスト更新メッセージ1616を送信する。OpCode1654は、x000A、すなわち、実効CQコンテキスト更新メッセージ1616に設定される。長さフィールドは、x0004、すなわち実効CQ S / Sコンテキスト更新メッセージの長さに設定される。長さに続く2つのフィールドは以下のように設定される。

40

【0216】

【表 3 5】

CQ番号
現在の実効完了CQE番号

【 0 2 1 7 】

代替RNIC1680は、実効CQ S / Sコンテキスト更新メッセージ1616を受信し、有効性を検証する。実効CQ S / Sコンテキスト更新メッセージ1616が有効である場合、例えばCRC32 1662が有効である場合、代替RNIC1680は、保持する実効CQ S / Sコンテキスト1272のローカル・コピーを、実効CQ S / Sコンテキスト・メッセージ1516の内容で更新する。代替RNIC1680における実効CQ S / Sコンテキスト1272のローカル・コピーの更新では、実効CQ S / Sコンテキスト1272の以下のフィールドが更新される。

10

【 0 2 1 8 】

【表 3 6】

現在の実効完了CQE番号

20

【 0 2 1 9 】

代替RNIC1680は次いで、自身の実効CQ S / Sコンテキスト1624のローカル・コピーを、コミット済みCQ S / Sコンテキスト1276のローカル・コピーにコピーする。コミット済みCQ S / Sコンテキスト1276の以下に示すフィールドが、代替RNIC1680で更新される。

【 0 2 2 0 】

【表 3 7】

コミット済み現在の完了CQE番号

30

【 0 2 2 1 】

上記のフィールドがローカルQPコンテキストで更新されると、代替RNIC1680は、コミット済みCQ S / Sコンテキスト更新メッセージ1624を主要RNIC1604に送信する。OpCode1654は、x000B、すなわちコミット済みCQ S / Sコンテキスト更新メッセージ1624に設定される。長さフィールドは、x0004、すなわち、コミット済みCQ S / Sコンテキスト更新メッセージの長さに設定される。長さに続く2つのフィールドは以下のように設定される。

【 0 2 2 2 】

【表 3 8】

40

CQ番号
コミット済み現在の完了CQE番号

【 0 2 2 3 】

主要RNIC1604は、コミット済みCQ S / Sコンテキスト更新メッセージ1624を受信し、有効性を検証する。コミット済みCQ S / Sコンテキスト・メッセージ1624が有効である、例えばCRC32が有効である場合、主要RNIC1604は、保持するコミット済みCQ S / Sコンテキスト1276のローカル・コピーを、コミッ

50

ト済みCQ S / Sコンテキスト・メッセージ1624の内容で更新する。すなわち、保持するコミット済みCQ S / Sコンテキスト1276の次のフィールドを更新する。

【0224】

【表39】

コミット済み現在の完了CQE番号

【0225】

コミット済みCQ S / Sコンテキスト更新メッセージ1624が無効であるか、チェックポイントの期限内に受信されない場合、主要RNIC1604は、最高でチェックポイント再試行回数によって定義される回数までそれを再送信する。代替RNIC1680が、重複した実効CQ S / Sコンテキスト更新メッセージ1616を受信した場合、代替RNIC1680は、1つ前のコミット済みCQ S / Sコンテキスト更新メッセージ1624を再送信する。この時点で、CQEは、チェックポイントされている。

10

【0226】

図21に、本発明の好ましい実施形態による、2つのRNICで使用される、主要RNICから代替RNICへのスイッチオーバー・メッセージの流れとプロセスを説明する概略図を提供する。主要RNIC1700などの主要RNICと、代替RNIC1796などの代替RNICは、時間Nごとに、「I Am Alive」チェックポイント・メッセージ1704および1772などの「I Am Alive」チェックポイント・メッセージを発行し、Nは、RNIC照会verbを通じて問い合わせることができ、RNIC変更verbを通じて変更できるプログラム可能なRNIC属性である。

20

【0227】

以下は、主要RNIC1700から代替RNIC1796に切り替えるために使用される例示的なスイッチオーバー・メッセージの流れである。まず、代替RNIC1796が、「I Am Alive」のカウント・ダウン・タイマを設定する。代替RNIC1796が、「I Am Alive」カウント・ダウン・タイマが切れる前に「I Am Alive」メッセージ1704を受信した場合は、タイマがリセットされる。「I Am Alive」メッセージのOpCode1754は、x0000、すなわち「I Am Alive」メッセージ1704および1750に設定される。長さフィールドは、x0003、すなわち「I Am Alive」メッセージの長さに設定される。長さに続く3つのフィールドは、以下のように設定される。

30

【0228】

【表40】

RNIC識別子
RNICポート番号
RNICポート番号状態

40

【0229】

複数ポートを有するRNICの場合、「I Am Alive」メッセージは、下に示すように、長さフィールドの後に可変数のフィールドを搬送するように実施することができる（Nは、RNICによってサポートされる最大ポート数）。

【0230】

【表 4 1】

RNIC識別子
RNICポートの数
RNICポート番号1
RNICポート番号状態1
RNICポート番号2
RNICポート番号状態2
RNICポート番号N
RNICポート番号状態N

10

【 0 2 3 1】

「I Am Alive」カウント・ダウン・タイマが切れる前に代替RNIC 1796が「I Am Alive」メッセージ1704を受信しない場合、代替RNIC 1796は、スイッチオーバー要求非同期イベント1708をホストに送信することにより、スイッチオーバーを開始する。ホストは、スイッチオーバー要求非同期イベント1708を受信し、RNIC照会1712を通じて主要RNIC 1700にアクセスすることを試みる。主要RNIC 1700は、RNIC照会の結果1716を返す。

20

【 0 2 3 2】

RNICの照会1712が成功し、RNIC照会の結果1716が、主要RNIC 1700が完全に機能しており、エラーのない状態にあることを示す場合、ホストは、

A) スwitch管理プロトコルを通じて、代替RNIC 1796が、それが接続されたスイッチからアクセスできないようにすることを要求し、

B) RNIC変更1776 verbを発行することにより、代替RNIC 1796のLANアドレス(例えばイーサネット(R)のMACアドレス)を、主要RNIC 1700のLANアドレス・テーブルに追加し、

C) ネットワーク・アドレス変更プロトコル(例えば gratuitous ARP 応答)を通じて、追加的なLANアドレス(それまで代替RNIC 1796によって使用されていたアドレス)が、主要RNIC 1700のアドレスの1つとして認識されることを要求し、

30

D) 主要RNIC 1700のCQ変更1790を使用して、CQC 1712の状態を変更し、

E) 主要RNIC 1700が、CQC 1712の主要RNICおよび代替RNIC識別コンテキストのRNICの状態を代替状態から主要状態に変えることにより、CQC 1712を主要状態にし、

F) 主要RNIC 1700のQP変更1788を使用して、QPC 1716の状態を変更し、

G) 主要RNIC 1700が、QPC 1716の主要RNICおよび代替RNIC識別コンテキストのRNICの状態を、代替RTS状態から主要RTS状態に変えることにより、QPC 1716を主要RTS状態にする。

40

【 0 2 3 3】

ステップDおよびEは、主要RNIC 1700を主要RNICとして、代替RNIC 1796を代替RNICとして作成されたすべてのCQに行われる。CQ変更1790には最適化を行うことができ、この最適化は、個々のCQ変更 verb ではなくCQのリストを変更するCQセット変更 verb としてCQ変更 verb 1790が発行されるものである。

【 0 2 3 4】

ステップFおよびGは、主要RNIC 1700を主要RNICとして、代替RNIC 1

50

796を代替RNICとして作成されたすべてのQPに行われる。QP変更1788には最適化を行うことができ、この最適化は、複数の個々のQP変更verbを送信するのではなくQPのリストを変更するQPセット変更verbを発行するものである。

【0235】

RNIC照会1712が失敗した(例えば主要RNIC1700がRNIC照会結果1716を返さない)、またはRNIC照会1712は成功したもののRNIC照会の結果1716にエラーがある、あるいは主要RNIC1700がエラー状態にあることを示す場合、ホストは、

A) RNIC変更1720を発行することにより、主要RNIC1700のLANアドレス(例えばイーサネット(R)のMACアドレス)を、代替RNIC1796のLAN 10
アドレス・テーブルに追加し、

B) 周知のスイッチ管理プロトコルを通じて、主要RNIC1700が、それが接続されたスイッチからアクセスできないようにすることを要求し、

C) 周知のネットワーク・アドレス変更プロトコル(例えばgratuitousARP応答)を通じて、追加的なLANアドレス(すなわちそれまで主要RNIC1700によって使用されていたアドレス)が、代替RNIC1796のアドレスの1つとして認識されることを要求し、

D) 代替RNIC1796のCQ変更1782を使用して、CQC1724の状態を変更し、

E) 代替RNIC1796が、CQC1724の主要RNICおよび代替RNIC識別 20
コンテキストのRNICの状態を代替状態から主要状態に変えることにより、CQC1824を主要状態にし、

F) 代替RNIC1796のQP変更1780を使用して、QPC1720の状態を変更し、

G) 代替RNIC1796が、QPC1720の主要RNICおよび代替RNIC識別
コンテキストのRNICの状態を代替RTS状態から主要RTS状態に変えることにより、QPC1720を主要RTS状態にする。

【0236】

ステップDおよびEは、主要RNIC1700を主要RNICとして、代替RNIC1796を代替RNICとして作成されたすべてのCQに行われる。CQ変更1782には 30
最適化を行うことができ、この最適化は、複数の個々のCQ変更verbではなくCQのリストを変更するCQセット変更verbを発行するものである。

【0237】

ステップFおよびGは、主要RNIC1700を主要RNICとして、代替RNIC1796を代替RNICとして作成されたすべてのQPに行われる。QP変更1780には最適化を行うことができ、この最適化は、個々のQP変更verbを発行するのではなくQPのリストを変更するQPセット変更verbを発行するものである。

【0238】

次いで図22に、本発明の好ましい実施形態による、2つのRNICによって使用される例示的な代替RNIC使用不可メッセージの流れとプロセスを説明する概略図を提供する。図22では、初めに、QPC1816が主要QPCであり、QPC1820がその代替QPCである。図22では、初めに、CQC1812が主要CQCであり、CQC1824がその代替CQCである。 40

【0239】

以下は、完全に動作していない代替RNIC1896を使用不可にするために使用される、例示的な代替RNIC使用不可メッセージの流れである。初めに、主要RNIC1800が、「I Am Alive」のカウント・ダウン・タイマを設定する。主要RNIC1700が、「I Am Alive」カウント・ダウン・タイマが切れる前に「I Am Alive」メッセージ1804を受信した場合は、タイマがリセットされる。「I Am Alive」メッセージは、Op Code1854がx0000、すなわち「I Am Alive」メッセージ1804および18 50

50に設定される。長さフィールドは、x0003、すなわち「I Am Alive」メッセージの長さに設定される。長さに続く3つのフィールドは、以下のように設定される。

【0240】

【表42】

RNIC識別子
RNICポート番号
RNICポート番号状態

10

【0241】

複数ポートを有するRNICの場合、「I Am Alive」メッセージは、下に示すように、長さフィールドの後に可変数のフィールドを搬送するように実施することができる（Nは、RNICによってサポートされる最大ポート数）。

【0242】

【表43】

RNIC識別子
RNICポートの数
RNICポート番号1
RNICポート番号状態1
RNICポート番号2
RNICポート番号状態2
RNICポート番号N
RNICポート番号状態N

20

【0243】

「I Am Alive」カウント・ダウン・タイマが切れる前に代替RNIC1800が「I Am Alive」メッセージ1804を受信しない場合、代替RNIC1800は、代替RNIC使用不可非同期イベント1808をホストに送信することにより、代替RNICの使用不可を開始する。ホストは、代替RNIC使用不可非同期イベント1808を受信し、RNIC照会1812を通じて代替RNIC1896にアクセスすることを試みる。主要RNIC1896は、RNIC照会の結果1816を返す。

30

【0244】

RNIC照会1812が成功し、RNIC照会の結果1816が、代替RNIC1896が完全に機能しており、エラーのない状態にあることを示す場合、ホストは、

A) RNIC変更1820を発行することにより、主要RNIC1800のLANアドレス（例えばイーサネット(R)のMACアドレス）を、代替RNIC1896のLAN

40

アドレス・テーブルに追加し、
B) 周知のスイッチ管理プロトコルを通じて、主要RNIC1800が、それが接続されたスイッチからアクセスできないようにすることを要求し、

C) 周知のネットワーク・アドレス変更プロトコル（例えば gratuitous ARP 応答）を通じて、追加的なLANアドレス（すなわちそれまで主要RNIC1800によって使用されていたアドレス）が、代替RNIC1896のアドレスの1つとして認識されることを要求し、

D) 代替RNIC1896のCQ変更1882を使用して、CQC1824の状態を変更し、

E) 代替RNIC1896が、CQC1824の主要RNICおよび代替RNIC識別

50

コンテキストのRNICの状態を代替状態から主要状態に変えることにより、CQC1824を主要状態にし、

F)代替RNIC1896のQP変更1880を使用して、QPC1820の状態を変更し、

G)代替RNIC1896が、QPC1820の主要RNICおよび代替RNIC識別コンテキストのRNICの状態を代替RTS状態から主要RTS状態に変えることにより、QPC1820を主要RTS状態にする。ステップDおよびEは、主要RNIC1800を主要RNICとして、代替RNIC1896を代替RNICとして作成されたすべてのCQに行われる。CQ変更1882には最適化を行うことができ、この最適化は、複数の個々のCQ変更verbを発行するのではなくCQのリストを変更するCQセット変更verbを発行するものである。

10

【0245】

ステップFおよびGは、主要RNIC1800を主要RNICとして、代替RNIC1896を代替RNICとして作成されたすべてのQPに行われる。QP変更1880には最適化を行うことができ、この最適化は、複数の個々のQP変更verbを発行するのではなくQPのリストを変更するQPセット変更verbを発行するものである。

【0246】

RNIC照会1812が失敗した(例えば代替RNIC1896がRNIC照会結果1816を返さない)か、または、RNICの照会1812は成功したが、RNIC照会の結果1816にエラーがある、あるいは代替RNIC1896がエラー状態にあることを示す場合、ホストは、

20

A)周知のスイッチ管理プロトコルを通じて、代替RNIC1896が、それが接続されたスイッチからアクセスできないようにすることを要求し、

B)RNIC変更1876を発行することにより、代替RNIC1896のLANアドレス(例えばイーサネット(R)のMACアドレス)を、主要RNIC1800のLANアドレス・テーブルに追加し、

C)周知のネットワーク・アドレス変更プロトコル(例えばgratuitousARP応答)を通じて、追加的なLANアドレス(すなわちそれまで代替RNIC1896によって使用されていたアドレス)が、主要RNIC1800のアドレスの1つとして認識されることを要求し、

30

D)主要RNIC1800のCQ変更1890を使用して、CQC1812の状態を変更し、

E)主要RNIC1800が、CQC1812の主要RNICおよび代替RNIC識別コンテキストのRNICの状態を代替状態から主要状態に変えることにより、CQC1812を主要状態にし、

F)主要RNIC1800のQP変更1888を使用して、QPC1816の状態を変更し、

G)主要RNIC1800が、QPC1816の主要RNICおよび代替RNIC識別コンテキストのRNICの状態を代替RTS状態から主要RTS状態に変えることにより、QPC1816を主要RTS状態にする。

40

【0247】

ステップDおよびEは、主要RNIC1800を主要RNICとして、代替RNIC1896を代替RNICとして作成されたすべてのCQに行われる。CQ変更1890には最適化を行うことができ、この最適化は、複数の個々のCQ変更verbを発行するのではなくCQのリストを変更するCQセット変更verbを発行するものである。

【0248】

ステップFおよびGは、主要RNIC1800を主要RNICとして、代替RNIC1896を代替RNICとして作成されたすべてのQPに行われる。QP変更1888には最適化を行うことができ、この最適化は、個々のQP変更verbではなくQPのリストを変更するQPセット変更verbを発行するものである。

50

【0249】

このように、本発明の好ましい実施形態では、RNICのスイッチオーバーとスイッチバックのサポートが提供される。意図された、または意図されない停止が主要RNICで発生した際に、本発明の好ましい実施形態で提供されるこの機構を使用して、すべての継続中の接続が代替のRNICに切り替えられ、代替RNICが通信処理を継続する。

【0250】

重要な点として、本発明について完全に機能するデータ処理システムとの関連で説明したが、当業者は、本発明の好ましい実施形態のプロセスは、命令のコンピュータ可読媒体の形態および各種形態で配布されることが可能であり、本発明は、配布を実施するために実際に使用される信号担持媒体の特定の種類に関係なく、等しく適用可能であることを認識されるであろう。コンピュータ可読媒体の例には、フロッピー（R）・ディスク、ハード・ディスク・ドライブ、RAM、CD-ROM、DVD-ROMなどの記録型媒体と、例えば無線周波や光波伝送などの伝送形態を使用する、デジタルおよびアナログの通信リンク、有線または無線の通信リンクなどの伝送型媒体が含まれる。コンピュータ可読媒体は、特定のデータ処理システムで実際に使用するために復号される、符号化形式の形をとることもできる。

【0251】

本発明の説明は例証と説明の目的で提示され、完全なものでも、ここに開示される形態の本発明に限定されるものでもない。当業者には、多くの変更形態と変形形態が自明である。

【図面の簡単な説明】

【0252】

【図1】本発明の好ましい実施形態により表した分散コンピュータ・システムの図である。

【図2】本発明の好ましい実施形態によるホスト・プロセッサ・ノードの機能ブロック図である。

【図3】本発明の好ましい実施形態によるIPスイート・オフロード・エンジンの図である。

【図4】本発明の好ましい実施形態によるスイッチの図である。

【図5】本発明の好ましい実施形態によるルータの図である。

【図6】本発明の好ましい実施形態による作業要求の処理を説明する図である。

【図7】TCPまたはSCTPトランスポートが使用される、本発明の好ましい実施形態による分散コンピュータ・システムの一部の図である。

【図8】本発明の好ましい実施形態によるデータ・フレームの図である。

【図9】本発明の好ましい実施形態による、例示的な要求と肯定応答のトランザクションを説明する分散コンピュータ・システムの一部の図である。

【図10】本発明の好ましい実施形態による分散ネットワーク・システムで使用されるネットワーク・アドレッシング指定を説明する図である。

【図11】本発明の好ましい実施形態におけるサブネットを含む分散コンピュータ・システムの一部の図である。

【図12】本発明の好ましい実施形態で使用される階層化通信アーキテクチャの図である。

【図13】本発明の好ましい実施形態によるホスト・ソフトウェアとのRNICインタフェースを表す概略図である。

【図14】本発明の好ましい実施形態による、2つのRNICによって維持されるスイッチオーバーとスイッチバックに関連する例示的な接続状態を表す概略図である。

【図15】本発明の好ましい実施形態による、2つのRNIC間の接続の例示的な初期化プロセスを概説するフローチャートである。

【図16】本発明の好ましい実施形態による、2つのRNIC間の接続の例示的な再同期プロセスを概説するフローチャートである。

10

20

30

40

50

【図17】本発明の好ましい実施形態による、主要RNICと代替RNIC間で使用されるチェックポイント・メッセージを説明する例示的な概略図である。

【図18】本発明の好ましい実施形態による、2つのRNICによって使用される送信のチェックポイント・メッセージの流れとプロセスを説明する例示的な概略図である。

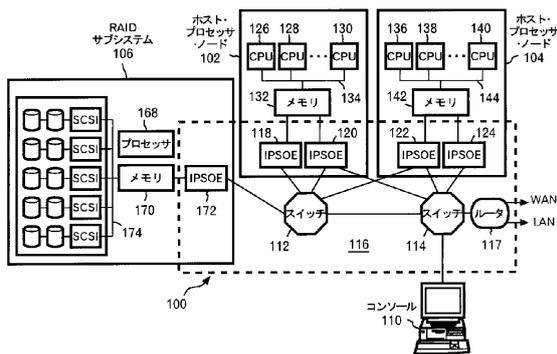
【図19】本発明の好ましい実施形態による、2つのRNICによって使用される受信のチェックポイント・メッセージの流れとプロセスを説明する例示的な概略図である。

【図20】本発明の好ましい実施形態による、2つのRNICによって使用される完了のチェックポイント・メッセージの流れとプロセスを説明する例示的な概略図である。

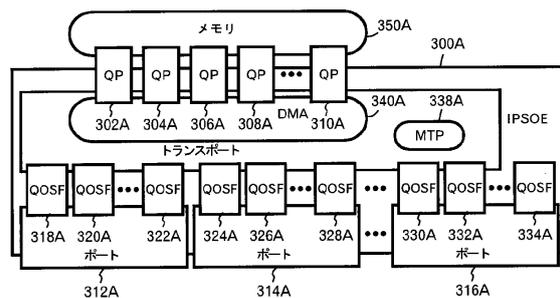
【図21】本発明の好ましい実施形態による、2つのRNICによって使用される主要RNICから代替RNICへのスイッチオーバー・メッセージの流れとプロセスを説明する例示的な概略図である。

【図22】本発明の好ましい実施形態による、代替QPアクティブ化メッセージの流れとプロセスを説明する例示的な概略図である。

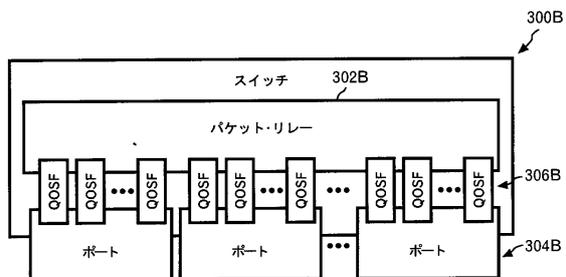
【図1】



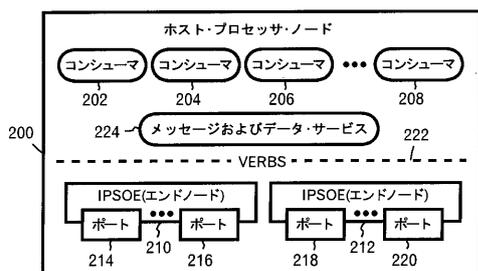
【図3】



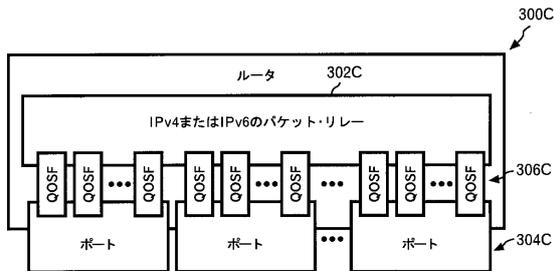
【図4】



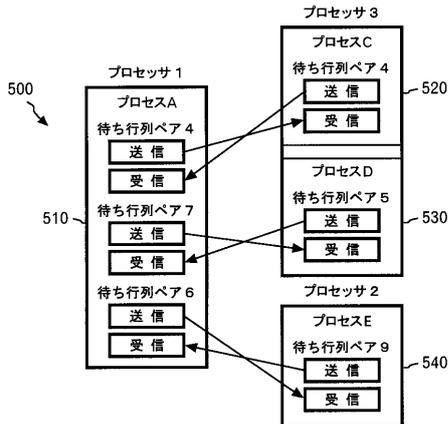
【図2】



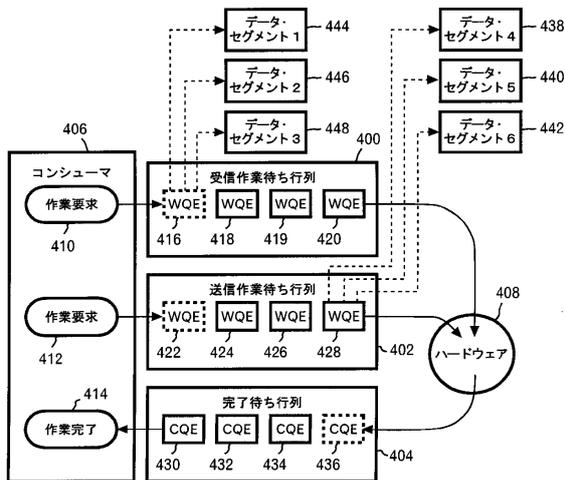
【図5】



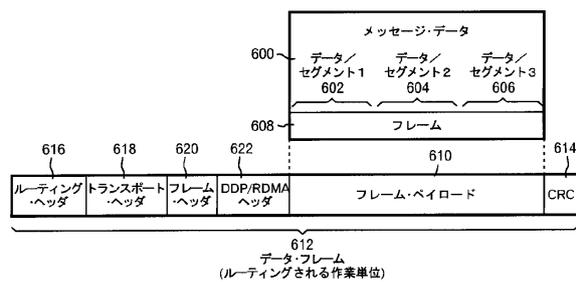
【図7】



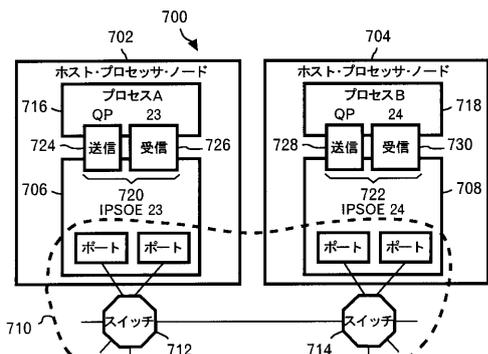
【図6】



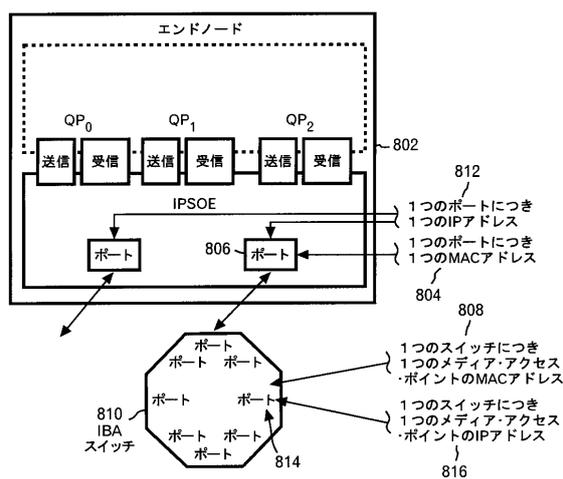
【図8】



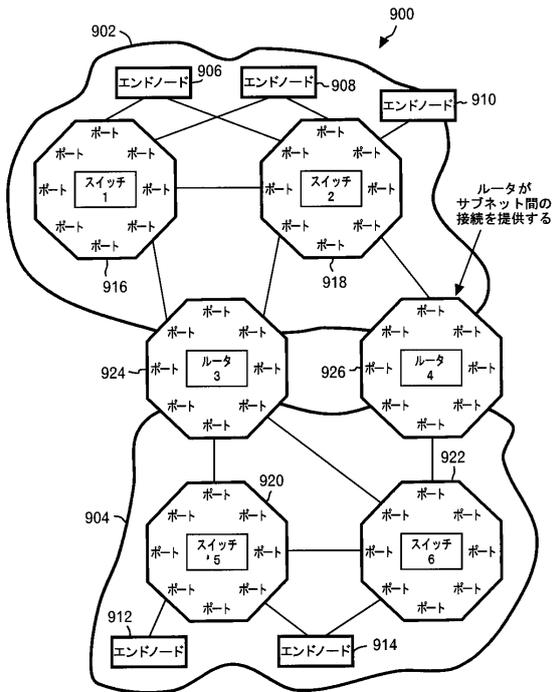
【図9】



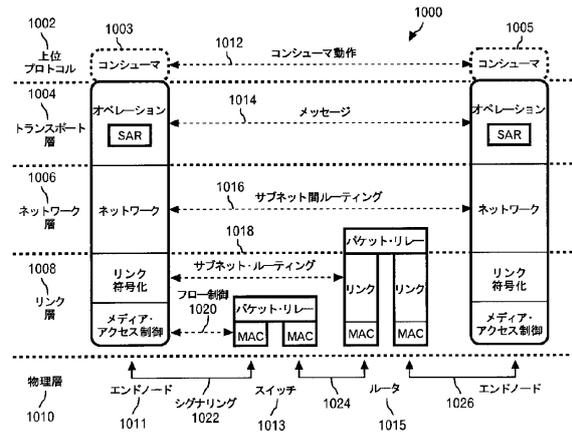
【図10】



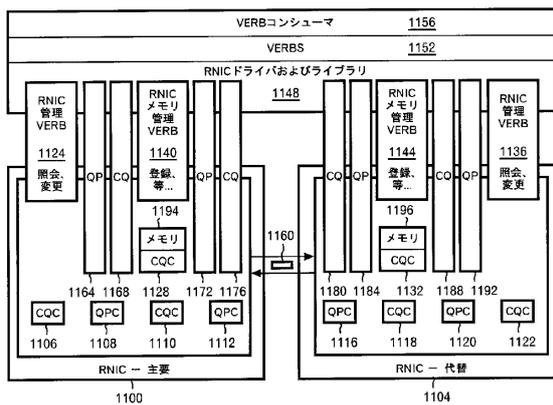
【 図 1 1 】



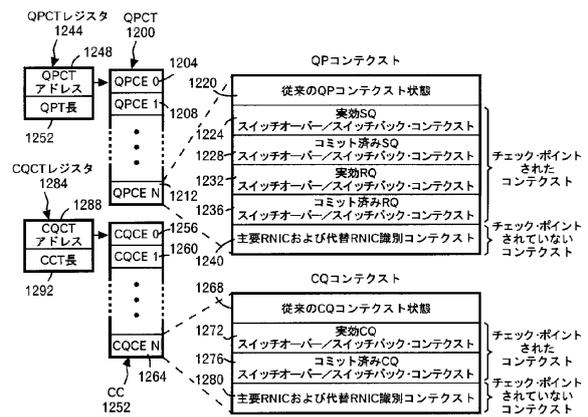
【 図 1 2 】



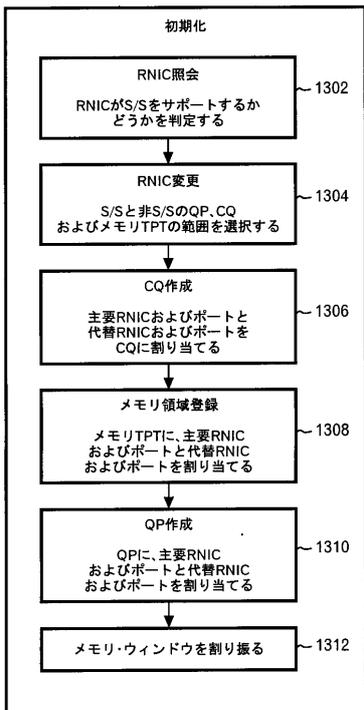
【 図 1 3 】



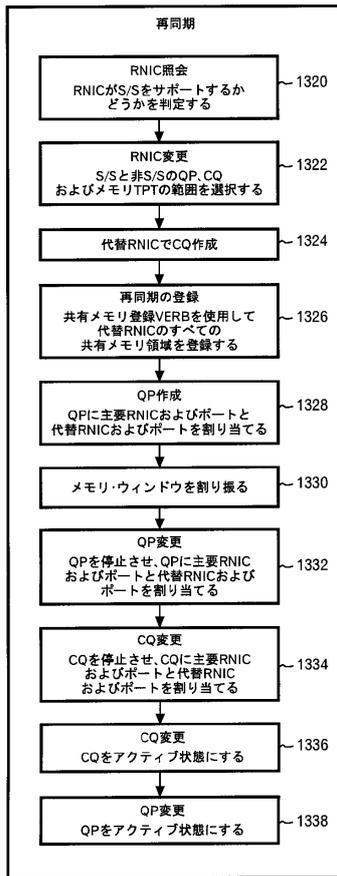
【 図 1 4 】



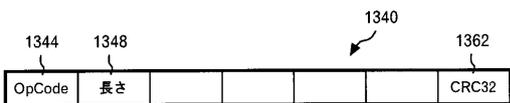
【 図 1 5 】



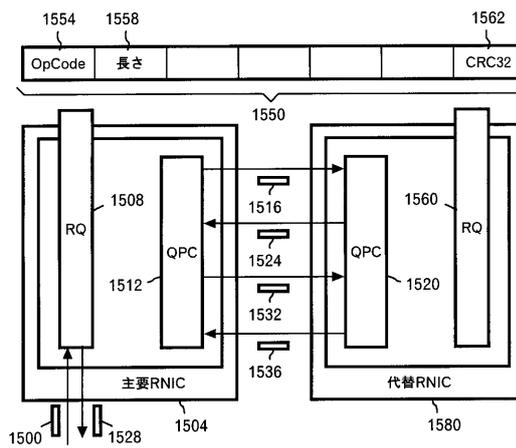
【 図 1 6 】



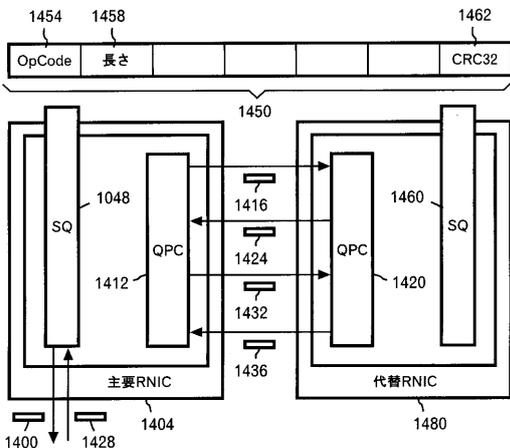
【 図 1 7 】



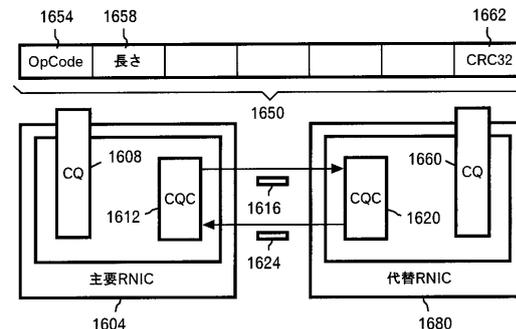
【 図 1 9 】



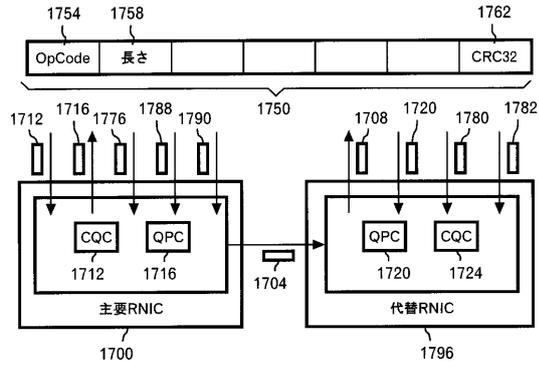
【 図 1 8 】



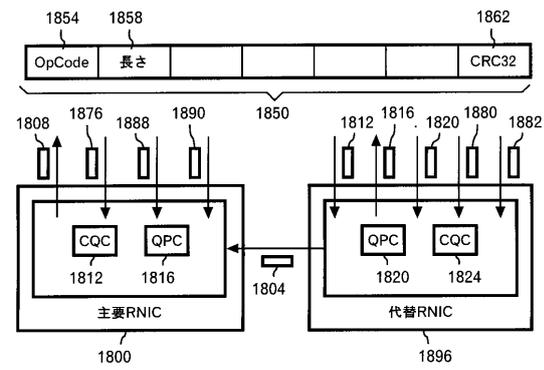
【 図 2 0 】



【 図 2 1 】



【 図 2 2 】



フロントページの続き

- (72)発明者 ボイド、ウイリアム、トッド
アメリカ合衆国12603 ニューヨーク州ポーキープシー カーメン・ドライブ 37
- (72)発明者 ジョゼフ、ダグラス
アメリカ合衆国06811 コネチカット州ダンバリー ブラグドン・アベニュー 6
- (72)発明者 コ、マイケル、アンソニー
アメリカ合衆国95120 カリフォルニア州サンノゼ クイーンズブリッジ・コート 1064
- (72)発明者 レシオ、レナト、ジョン
アメリカ合衆国78759 テキサス州オースティン ウィニペグ・コーブ 6707

審査官 吉田 隆之

- (56)参考文献 特開平8-221289(JP,A)
特開平11-68783(JP,A)
特開平11-275224(JP,A)
特開平5-167679(JP,A)
特開昭63-165950(JP,A)
特開2003-216592(JP,A)
特公平7-72882(JP,B2)
実公昭62-17879(JP,Y1)
信学技報 SSE91-114
信学誌, Vol.73 No.11, p1179-1184

- (58)調査した分野(Int.Cl., DB名)
H04L 12/00