

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4264866号
(P4264866)

(45) 発行日 平成21年5月20日(2009.5.20)

(24) 登録日 平成21年2月27日(2009.2.27)

(51) Int.Cl.

F I

G 0 6 F 13/00 (2006.01)

G 0 6 F 13/00 3 5 3 C

H 0 4 L 29/06 (2006.01)

H 0 4 L 13/00 3 0 5 A

請求項の数 28 (全 35 頁)

(21) 出願番号 特願2000-568012 (P2000-568012)
 (86) (22) 出願日 平成10年11月20日(1998.11.20)
 (65) 公表番号 特表2002-524005 (P2002-524005A)
 (43) 公表日 平成14年7月30日(2002.7.30)
 (86) 国際出願番号 PCT/US1998/024943
 (87) 国際公開番号 W02000/013091
 (87) 国際公開日 平成12年3月9日(2000.3.9)
 審査請求日 平成17年11月21日(2005.11.21)
 (31) 優先権主張番号 09/141, 713
 (32) 優先日 平成10年8月28日(1998.8.28)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 501082978
 アラクリテック・インコーポレイテッド
 アメリカ合衆国, カリフォルニア州・95
 012, サン・ノゼ, イースト・ギッシュ
 ・ロード・234
 (74) 代理人 100098062
 弁理士 梅田 明彦
 (72) 発明者 ブーチャー, ローレンス, ビー
 アメリカ合衆国, カリフォルニア州・95
 070, サラトガ, モントルボ・ドライブ
 ・20605
 (72) 発明者 ブライトマン, スチーブン, イー, ジェイ
 アメリカ合衆国, カリフォルニア州・95
 132, サン・ノゼ, アーレン・コート・
 3733

最終頁に続く

(54) 【発明の名称】 通信を高速化するインテリジェントネットワークインタフェース装置及びシステム

(57) 【特許請求の範囲】

【請求項1】

ハードウェアとプロセッサとプロトコル処理層のスタックとを有するホストコンピュータによりネットワーク通信を行うための方法であって、

ネットワークから前記ホストコンピュータによって、複数のヘッダと、前記ヘッダのプロトコル処理に従って前記ホストコンピュータの行先に配置するためのデータとを含むメッセージパケットを受信する過程と、

前記ハードウェアにより、前記ヘッダの要約を作成することを含めて、前記複数のヘッダをグループとして処理する過程と、

前記データを前記ホストコンピュータの行先に格納するために、前記要約に基づいて、前記プロトコル層によって前記パケットを処理するか、又は前記プロトコル層による処理を回避するかを選択する過程と、

前記データを前記ヘッダなしで前記要約に従って、前記プロトコル層により前記ヘッダを処理することなく、前記ホストコンピュータ上のアプリケーションにより提供される前記行先に転送する過程とを有することを特徴とする方法。

【請求項2】

前記パケットに対応する伝送制御プロトコル(TCP)コネクションのための通信制御ブロックを作成する過程と、前記データを前記行先に送出するために、前記要約と前記通信制御ブロックとを整合させる過程とをさらに含むことを特徴とする請求項1に記載の方法。

10

20

【請求項 3】

前記パケットに対応する伝送制御プロトコル（TCP）コネクションのための通信制御ブロックを作成する過程をさらに含み、前記データを前記行先に送出する過程は、前記通信制御ブロックによって前記データを誘導する過程を含むことを特徴とする請求項 1 に記載の方法。

【請求項 4】

前記通信制御ブロックを参照することにより、前記ホストコンピュータから前記ネットワークに第 2 のメッセージを送信する過程をさらに含むことを特徴とする請求項 3 に記載の方法。

【請求項 5】

通信プロセッサを備える装置によりネットワークに接続され、プロトコル処理スタックを有するホストコンピュータにより通信を処理するための方法であって、

前記ホストコンピュータによって前記ネットワークから前記プロトコル処理スタックに対応する一連のヘッダとデータとを含むメッセージフレームを受信する過程と、

前記フレームについて発信元及び行先のインターネットプロトコル（IP）アドレスと発信元及び行先の伝送制御プロトコル（TCP）ポートとを決定することを含めて、前記装置によって前記一連のヘッダをバイトのストリームとして解析する過程と、

前記装置により前記データを前記ヘッダなしで、前記 IP アドレス及び TCP ポートに基づいて、前記プロトコル処理スタックによって前記ヘッダを処理することなく、かつ前記ホストコンピュータによって前記データをコピーすることなく、前記ホストコンピュータ上で実行中のアプリケーションにより制御されるホストメモリ内の行先に送出する過程とを有することを特徴とする方法。

【請求項 6】

前記ホストメモリ内の行先に前記ヘッダなしで前記データを送出する前記過程が、前記ホストコンピュータによって、前記フレームを含むメッセージのための通信制御ブロックを作成する過程と、

前記装置に前記通信制御ブロックを格納する過程と、

前記 IP アドレス及び TCP ポートを前記通信制御ブロックと整合させる過程とをさらに含むことを特徴とする請求項 5 に記載の方法。

【請求項 7】

前記データを前記行先に送出する前に、前記 IP アドレス及び TCP ポートのハッシュを前記通信制御ブロックを識別するハッシュテーブルと比較する過程をさらに含むことを特徴とする請求項 6 に記載の方法。

【請求項 8】

前記装置を経由して、前記ホストコンピュータから前記ネットワークに送信データを伝送する過程をさらに含み、前記伝送する過程は、リモートホストへのネットワーク転送のために、前記送信データの前にいくつかのプロトコルヘッダを同時に付加する過程を含むことを特徴とする請求項 5 に記載の方法。

【請求項 9】

プロセッサとプロトコル層の順次式スタックとを有し、ネットワークインタフェースによりネットワークに接続されたホストコンピュータにより通信を行うための方法であって、該方法は、

前記ネットワークから前記ネットワークインタフェースによって、それぞれにデータ部分と、前記プロトコル層の順次式スタックにより処理した場合には、前記データのための前記ホストコンピュータ内の上位層の行先に前記データ部分を格納させる、関連するヘッダのシーケンスとを含む多数のパケットを有するメッセージを受信する過程と、

前記ネットワークインタフェースによって、前記関連するヘッダなしで、かつ任意のホスト CPU への割り込みを発生することなく、複数の前記データ部分を前記行先に送出する過程とを有することを特徴とする方法。

【請求項 10】

前記データ部分を前記行先に送出する過程の前に、前記プロトコル層のスタックによって前記パケットを処理するか否かを選択する過程をさらに含むことを特徴とする請求項9に記載の方法。

【請求項 1 1】

前記ヘッダを用いることなく前記データ部分を前記行先に送出する前に、前記ネットワークインタフェースで前記ヘッダを処理する過程をさらに含むことを特徴とする請求項9に記載の方法。

【請求項 1 2】

前記ネットワークインタフェースを介して前記ホストから前記ネットワークにデータファイルを伝送する過程をさらに含み、前記伝送する過程は、前記データファイルを一連のデータユニットに分割する過程と、前記データユニットの前に前記ヘッダを付加し、それにより一連のネットワークフレームを作成する過程と、任意のホストCPUに割り込みを発生することなく、前記ネットワークフレームを前記ネットワークに配置する過程とを含むことを特徴とする請求項9に記載の方法。

10

【請求項 1 3】

CPUと、メモリと、プロトコル層の順次式スタックとを備え、ネットワークインタフェース装置によりネットワークに接続されたホストコンピュータによって通信を行うための方法であって、

前記ネットワークインタフェース装置によって前記ネットワークから第1のパケットを受信する過程と、

20

前記ホストコンピュータによって、前記第1のパケットに対応する伝送制御プロトコル(TCP)状態を含む通信制御ブロックを作成する過程を含めて、前記第1のメッセージを処理する過程と、

前記通信制御ブロックを前記装置に転送する過程と、

前記ネットワークインタフェース装置によって前記ネットワークから、データとヘッダとを含み、かつ前記通信制御ブロックに対応する第2のパケットを受信する過程と、

前記ネットワークインタフェース装置によって前記ヘッダのTCP部分を処理することを含めて、前記ネットワークインタフェース装置によって前記ヘッダを処理する過程と、

前記ネットワークインタフェース装置によって前記データを、前記ホストコンピュータにより前記ヘッダを処理することなく、前記プロトコル層のTCP層より上位にある上位層に、該上位層に適した形で送出し、前記通信制御ブロックを用いて前記送出を誘導する過程とを有することを特徴とする方法。

30

【請求項 1 4】

前記通信制御ブロックに対応する第3のパケットを前記装置によって受信する過程と、前記装置から前記メモリに前記通信制御ブロックの制御を渡す過程と、前記CPUにより前記第3のパケットを処理する過程とをさらに含むことを特徴とする請求項1 3に記載の方法。

【請求項 1 5】

前記データを前記上位層に送出する前に、前記第2のパケットのIPアドレスとTCPポートとを前記通信制御ブロックに整合させる過程をさらに含むことを特徴とする請求項1 3に記載の方法。

40

【請求項 1 6】

前記ホストコンピュータから前記ネットワークに第3のパケットを伝送する過程をさらに含み、該過程は、前記通信制御ブロックを参照し、かつ複数のプロトコル層ヘッダを含む送信ヘッダをホストソースから取得されるデータに付加することにより、前記装置を経由して前記第3のパケットを送出する過程を含むことを特徴とする請求項1 3に記載の方法。

【請求項 1 7】

プロトコル処理装置によりネットワークに接続され、プロトコル処理スタックを有するローカルホストによりネットワーク通信を行うための方法であって、

50

前記プロトコル処理スタックにより、前記ローカルホストとリモートホストとの間の伝送制御プロトコル（TCP）コネクションを定義する通信制御ブロックを作成する過程と、

前記通信制御ブロックを前記プロトコル処理装置に渡し、それにより前記コネクションに関連しかつ前記ネットワークと前記ローカルホストとの間で転送されるメッセージパケットのTCP処理の制御を渡し、それにより、前記メッセージパケットを前記プロトコル処理スタックではなく、前記プロトコル処理装置によって処理する過程とを有することを特徴とする方法。

【請求項 18】

前記ネットワークと前記ローカルホストとの間で転送されかつ前記コネクションに関連する第2のメッセージパケットが全般に前記プロトコル処理スタックによって処理されるように、前記通信制御ブロックの制御を前記ローカルホストに渡す過程をさらに含むことを特徴とする請求項 17 に記載の方法。

【請求項 19】

前記プロトコル処理装置によって、前記ネットワークからメッセージフレームを受信する過程と、

前記プロトコル処理装置によって、前記メッセージフレームのヘッダを解析し、それにより前記メッセージフレームのTCP発信元ポート及び行先ポートを決定する過程と、

前記TCP発信元ポート及び行先ポートを前記通信制御ブロックと比較する過程とをさらに含むことを特徴とする請求項 17 に記載の方法。

【請求項 20】

前記プロトコル処理装置によって、前記ネットワークにメッセージフレームを伝送する過程をさらに含み、前記伝送する過程は、前記通信制御ブロックに基づいてヘッダを形成する過程と、前記メッセージフレームの前に前記ヘッダを付加する過程とを含むことを特徴とする請求項 17 に記載の方法。

【請求項 21】

ネットワークを経由してリモートホストに接続することができるローカルホストで用いるための装置であって、前記ローカルホストは、前記ローカルホストのアプリケーションと前記リモートホストのアプリケーションとの間のTCPコネクションを定義する通信制御ブロックを作成するプロトコル処理層のスタックを動作させるCPUを含み、

該装置は、

前記CPUと前記ネットワークとに接続され、前記プロトコル処理層によってネットワークメッセージを処理するか、前記プロトコル処理層を回避し、前記ネットワークと前記ローカルホストアプリケーションとの間で前記ネットワークメッセージに含まれるデータを転送するために前記通信制御ブロックを用いるかを選択するために構成されるプロセッサを有する通信処理機構を備えることを特徴とする装置。

【請求項 22】

前記プロセッサは、前記通信制御ブロックをその前記メッセージの処理に基づいて更新するように構成されることを特徴とする請求項 21 に記載の装置。

【請求項 23】

前記通信処理機構は、前記プロセッサに接続され、前記ネットワークから受信されるメッセージパケットの要約を生成するために構成される受信シーケンサを備え、前記メッセージパケットは、前記プロトコル層のスタックに対応する制御情報を含み、前記プロセッサは前記要約を前記通信制御ブロックと比較するように構成されていることを特徴とする請求項 21 に記載の装置。

【請求項 24】

前記プロセッサは、前記通信制御ブロックに対応するヘッダを作成し、前記プロトコル処理層のいくつかに対応する制御情報を含み、前記ホストから前記リモートヘッダに前記メッセージを伝送するために前記データの前に前記ヘッダを付加し、前記通信制御ブロックを更新するように構成されることを特徴とする請求項 21 に記載の装置。

【請求項 2 5】

前記通信処理機構は、前記通信制御ブロックに基づいて、前記データに付随する全ゆるヘッダを用いることなく、前記通信処理機構から前記ホストアプリケーションに前記データを送出するためのダイレクトメモリアクセスユニットを備えることを特徴とする請求項 2 1 に記載の装置。

【請求項 2 6】

前記プロセッサは複数のマイクロプロセッサを備え、前記マイクロプロセッサのうちの少なくとも 1 つは主に、前記ネットワークから前記ホストによって受信されるメッセージを処理するように構成され、前記マイクロプロセッサのうちの第 2 のマイクロプロセッサは前記ホストから前記ネットワークに伝送されるメッセージを処理するように構成されることを特徴とする請求項 2 2 に記載の装置。

10

【請求項 2 7】

前記マイクロプロセッサは、フェーズをローテーションする際に共有されるハードウェア機能を利用することを特徴とする請求項 2 6 に記載の装置。

【請求項 2 8】

前記通信処理機構は、複数のキュー記憶ユニットにおいて情報をキューに入れるために構成されるキューマネージャを備え、前記キュー記憶ユニットのうちの少なくとも 1 つは S R A M および D R A M を含むことを特徴とする請求項 2 1 に記載の装置。

【発明の詳細な説明】

【0001】

20

(技術分野)

本発明は、一般的にコンピュータ又は他のネットワークに関するものであり、得にネットワークに接続されたコンピュータのようなホスト間で通信される情報の処理に関するものである。

【0002】

(背景技術)

ネットワークコンピューティングの利点はますます明らかとなってきた。パーソナルコンピュータや他のエンドユーザ装置において情報、通信、又は計算能力を個人に提供することがもたらす利便性及び能率のために、イントラネット装置やアプリケーション及びインターネット等のネットワークコンピューティングが急速に成長することになった。

30

【0003】

周知のように、大部分のネットワークコンピュータ通信は、ネットワークに接続されたホストコンピュータ間での情報の転送のための階層化ソフトウェアアーキテクチャの支援によって達成される。その層が、情報を処理可能なセグメントに分割することを助け、各層の一般的な機能は、多くの場合開放型システム間相互接続 (O S I) と呼ばれる国際標準に準拠している。O S I は 7 種の処理の層を規定しており、情報がホストによって受信されると、それはエンドユーザに提示できるようにするためにそれらの層の間をわたされる。同様に、ホストからネットワークへの情報の送信では、情報がそれらの 7 種の処理層を逆の順に通過し得る。1 つの層による処理及びサービスの各ステップでは、処理された情報のコピーが行われる。T C P / I P (T C P は伝送制御プロトコルを表し、I P はインターネットプロトコルを意味する) と呼ばれる、広く実現されている別の参照モデルは、基本的に O S I の 7 層のうち 5 層を利用する。

40

【0004】

ネットワークは、例えば、それぞれが複数のホストを有する異なるローカルエリアネットワーク (L A N) 間のイーサネット接続又はインターネット接続のような高速パス、又はホスト間のデータ伝送のための様々な他の既知の手段の何れかを含み得る。O S I 標準によれば、物理層は各ホストにおいてネットワークに接続され、その物理層がネットワークを介しての生のデータビットの送信及び受信を提供する。データリンク層は、各ホストの物理層によりサービスされ、物理層から受け取ったデータのフレーム分割及び誤り訂正を提供するとともに、受信したホストによって送られた肯定応答フレームの処理を行う。各

50

ホストのネットワーク層は、各データリンク層によりサービスされ、主としてデータのパケットのサブネットのサイズ及び統合を制御する。

【 0 0 0 5 】

トランスポート層は各ネットワーク層からサービスされ、セッション層は各ホスト内の各トランスポート層によりサービスされる。トランスポート層はそれらの各セッション層からデータを受け取り、そのデータを他のホストのトランスポート層へ送信するためのより小さいユニットに分割する。他のトランスポート層は、各プレゼンテーション層への提示のためそのデータを連結する。セッション層によってホスト間の機能強化された通信の制御が可能となる。プレゼンテーション層はそれらの各セッション層によりサービスされ、各ホスト固有であり得るデータセマンティクス及び構文と、データ表現の標準化された構造との間での変換を行う。データの圧縮及び/または暗号化も、プレゼンテーション層レベルで達成され得る。アプリケーション層は、各プレゼンテーション層によりサービスされ、個々のホストに特有のプログラムと、アプリケーション又はエンドユーザの何れかに対して提示するための標準化されたプログラムとの間の変換を行う。TCP/IP標準は、下位の4層及びアプリケーション層を含むが、セッション層及びプレゼンテーション層の機能は隣接する層に統合している。一般的に、アプリケーション層、プレゼンテーション層、及びセッション層は上位層と定義され、トランスポート層、ネットワーク層、及びデータリンク層は下位層と定義される。

10

【 0 0 0 6 】

各層の規則及び規定はその層のプロトコルと呼ばれ、かつ各層のプロトコル及び一般的な機能は様々なホストにおいて概ね等しいことから、たとえ同位の層同士が情報を下位の各層を順次転送させることなく直接やり取りしないとしても、異なるホストの同一の層同士で直接行われる通信というものを考えることは有益である。各下位層は、通信される情報を処理を助けるべくその1つ上位の層にサービスする。各層は、処理及び次の層へのサービスの提供のため情報を保存する。一般的に用いられるハードウェア及びソフトウェアアーキテクチャや装置及びプログラムの多様性のために、各層は、介入し得るハードウェア及びソフトウェアの相違とは無関係に、データが適切な形式で目的の行先に到達し得ることを確実なものとする必要がある。

20

【 0 0 0 7 】

第1のホストから第2のホストに送信するためのデータの準備において、その層のプロトコルに関連する第1のホストの各層で幾つかの制御データが付加されるが、その制御データはそのホストの全ての下位の層において元のデータ(ペイロードデータ)から識別可能なものがある。従って、アプリケーション層はペイロードデータにアプリケーションヘッダを付加し、結合されたデータを送信側のホストのプレゼンテーション層に送る。プレゼンテーション層はその結合データを受け取り、それを処理してそのデータにプレゼンテーションヘッダを付加し、別の結合されたデータパケットが生成される。次に、ペイロードデータ、アプリケーションヘッダ、及びプレゼンテーションヘッダの結合によって生成されたデータをセッション層に渡し、セッション層は、そのデータへのセッションヘッダの付加や、生成されたデータの組み合わせのトランスポート層への提示等の必要な処理を行う。この処理は、情報が下位の層に移動する間継続され、それらの各層においてデータにトランスポートヘッダ、ネットワークヘッダ、及びデータリンクヘッダ、及びトレーラが付加され、各処理過程ではデータの転送及びコピーを行い、その後、データをビットパケットとしてネットワークを通して第2のホストに送信する。

30

40

【 0 0 0 8 】

受信側のホストは一般的に上述のプロセスの逆を行い、ネットワークからのデータビットの受信に始まり、最下位の(物理)層から最上位の(アプリケーション)層まで順番にヘッダが除去され、データが処理されて、受信側ホストの行先への送信が行われる。受信側ホストの各層は、その層に関連するヘッダのみを認識し、処理する。その層からみて、より上位の層の制御データはペイロードデータに含まれており、そこから識別できないからである。受信側ホストがそのデータをその目的の行先に適切な形式で渡すためには、複

50

数の割り込み、有効な中央処理装置（ＣＰＵ）による処理時間、及び反復されるデータコピー処理も必要であり得る。

【０００９】

上述の階層化プロトコル処理の説明は、主としてこの主題について詳述された大学レベルの教科書、例えばAndrew S. TanenbaumによるComputer Networks, Third Edition (1996)

が入手可能であることから、単純化されている。上記文献はこの引用により本明細書の一部とする。その文献に定義されているように、コンピュータネットワークとは、インターネットや、例えばローカルエリアネットワーク（ＬＡＮ）、ワイドエリアネットワーク（ＷＡＮ）、非同期転送モード（ＡＴＭ）、リング又はトークンリング、配線式、無線式、衛星、又は他の異なるプロセッサ間の通信能力を提供するための手段等のインターネットデバイスのような、相互接続された独立して操作されるコンピュータの集合体である。本明細書において、コンピュータは、データを処理するための論理機能とメモリ機能の両方を有する装置を含むものと定義されており、ネットワークに接続されたコンピュータ即ちホストは、それらが異なるオペレーティング装置に従って機能する場合、即ち異なるアーキテクチャを介して通信を行う場合、ヘテロジニアス（異種）であるという。

【００１０】

ネットワークがますます一般的となり、それによってやりとりされる情報がますます複雑化し内容豊富になるにつれて、そのようなプロトコル処理の必要性が高まってきた。ホストのＣＰＵの処理能力の大部分が、プロトコル処理の制御に充てられて、そのＣＰＵが他のタスクを行う能力は低下していると推定される。物理層やデータリンク層のような下位の層の補助のため、ネットワークインタフェースカードが開発されてきた。単に従来型の構成のＣＰＵにより大きい処理能力を付与することによって、プロトコル処理の速度を高めることも可能である。しかし、この解決方法は手際の良いものとは言えず、コストがかかる。しかし、様々なネットワーク、プロトコル、アーキテクチャ、オペレーティング装置、及びアプリケーションによって付与される複雑さによって、種々のネットワークホスト間での通信能力を与えるために大量の処理が要求されることになる。

【００１１】

（発明の開示）

本発明は、通信されるデータの転送の効率及びその処理速度を大幅に高める、ネットワーク通信の処理のための装置を提供する。本発明は、汎用プロセッサ上での階層化プロトコル処理を行うという、長年にわたって利用されてきた実施の形態を問い直すことによって達成されたものである。そのプロトコル処理方法と、その結果として形成されるアーキテクチャでは、例えばＴＣＰ／ＩＰのようなコネクションをベースにした階層化アーキテクチャの層を、所望の位置又はホスト上のバッファとの間で、より直接的にネットワークデータを送受信できる１層のより広い層に効果的にまとめる。この高速化処理はデータの送受信両方のためにホストに提供され、情報の交換に係するホストの両方がそのような特徴を有している場合も、その一方しかそのような特徴を有していない場合でも、処理能力を向上させる。

【００１２】

前記高速化処理では、その発信元において直接的にメッセージデータにアクセスしたりそれをその目的の行先に直接供給する高速パスを介してメッセージからのデータを処理できるようにする所定のメッセージに対応する制御命令を用いる。この高速パスは、データに付加されるヘッダの従来型のプロトコル処理をバイパスする。高速パスは、ネットワーク通信の処理のために設計された専用マイクロプロセッサを用いて、例えば反復されるコピー処理やＣＰＵへの割り込みのような、従来のソフトウェアの階層処理の遅れ及び問題点を回避している。要するに、前記高速パスは、従来型のネットワークスタックの複数の層において伝統的に見出される状態を、プロトコル層の厳格な区別及び分離を要求する従来の規定とは異なりこれらの層を全て包含する１個のステートマシンで置き換える。ホストは、高速パスコネクション又はメッセージ例外の処理をセットアップするために用いられ得る順次式プロトコル処理スタックを保持している。専用マイクロプロセッサ及びホスト

10

20

30

40

50

が、所定のメッセージまたはメッセージの一部をそのマイクロプロセッサによって処理するか、ホストのスタックによって処理するかを、インテリジェントに選択する。

【0013】

(発明の最良の実施の形態)

図1は、リモートホスト22にネットワーク25によって接続された本発明のホスト20を示す。本発明によって達成される処理速度の上昇は、既存のホストに容易に無理なく追加されるインテリジェントネットワークインタフェースカード(INIC)又はホストに統合される通信処理装置(CPD)によって提供され得、何れの場合もホストのCPUを大部分のプロトコル処理から解放し、そのCPUによって実行される他のタスクも改善することができる。第1の実施形態のホスト20は、PCIバス33によって接続されたCPU28及びCPD30を有する。CPD30はダイレクトメモリアクセス(DMA)ユニットによって制御されるメモリバッファ及び通信データの処理のために設計されたマイクロプロセッサを有する。例えば、半導体メモリ又はディスクドライブのような記憶機構35も、関連する制御機構とともにPCIバス33に接続されている。

10

【0014】

更に図2を参照すると、ホストのCPU28が記憶機構35に收容されたプロトコル処理スタック44を制御しており、前記スタックはデータリンク層36、ネットワーク層38、トランスポート層40、上位層46、及び上位層インタフェース42を有する。上位層46は、使用される特定のプロトコル及び通信されるメッセージに応じて、セッション層、プレゼンテーション層、及び/又はアプリケーション層を意味し得る。上位層インタフェース42は、CPU28及び関連するあらゆる制御機構と協働して、矢印48で示すように、上位層46又は記憶機構35へのファイルの転送及びそこからのファイルの取り出しを行える。コネクションコンテキスト50が生成されており、これは、後に説明するように、例えばプロトコルの種類及び各プロトコル層のための発信元及び行先アドレスのような様々なコネクションの特徴を要約する。そのコンテキストは、矢印52及び54で示すように、セッション層42のためのインタフェースとCPD30との間でやりとりされ得、CPD30又は記憶機構35の何れかに通信制御ブロック(CCB)として格納される。

20

【0015】

CPD30が特定のコネクションを確定するCCBを保持するとき、ネットワークからCPDによって受信されたそのコネクションに関連するデータはそのCCBによって参照され、次に高速バス58に従って記憶機構35に直接送信され得、データリンク層36、ネットワーク層38、及びトランスポート層40によって処理される順次式のプロトコルがバイパスされることになる。記憶機構35からリモートホスト22へのファイルの送信のようなメッセージの転送は高速バス58を介しても生じ得、その場合、トランスポート層40、ネットワーク層38、及びデータリンク層36による処理の間に、逐次ヘッダを付加するのではなく、CCBを参照するCPD30によってファイルデータのためのコンテキストが加えられる。CPD30のDMAコントローラが、これらのCPDと記憶機構35との間の転送を実行する。

30

【0016】

CPD30は、それぞれが可能な分離した状態を有する複数のプロトコルスタックを、高速バス処理のための1個のステートマシンにまとめる。この結果、その1個のステートマシンにおいては準備されていない例外条件が生じ得る。これようにしているのは、主としてそのような条件がまれにしか生じず、CPD上でそれを処理することはコストに対して殆ど処理能力上の利点を得られないからである。そのような例外は、CPD30又はCPU28で生じたものであり得る。本発明の利点として、高速バスCCB上で生ずる予期しない状況が取り扱われる方式がある。CPD30は、これらのまれな状況を制御ネゴシエーションを介してホストのプロトコルスタック44にCCB及び関連するあらゆるメッセージフレームを戻す、即ちフラッシュすることによって処理する。次にその例外条件は、従来の方式でホストのプロトコルスタック44によって処理される。しばらく後に、通常

40

50

は例外条件の処理が終了し、高速パス処理が再開し得るようになった直後に、ホストのスタック 44 は C C B を C P D に戻す。

【 0 0 1 7 】

このフォールバック能力によって、ホストのプロトコルの能力に影響を与える機能を C P D ネットワークマイクロプロセッサが取り扱うことができるようになるとともに、全体の処理能力に無視できる程度の影響しか与えない程度にまれにしか生じない例外がホストのスタックによって取り扱われることになる。カスタム設計のネットワークプロセッサは、ネットワーク情報の送受信のための独立したプロセッサと、補助及びキューのための別のプロセッサとを備え得る。好ましいマイクロプロセッサの実施形態は、パイプライン型の 3 個 1 組の受信、送信、及びユーティリティプロセッサを備える D M A コントローラがインプリメンテーションに組み入れられ、コントローラに隣接するバッファと長時間記憶機構のような他の位置との間でデータを迅速にやりとりさせるべく、ネットワークマイクロプロセッサと密接に協働する形で機能するものである。D M A コントローラに論理的に隣接するバッファを設けることによって、P C I バスに不必要な負荷が加わるのを避けられる。

【 0 0 1 8 】

図 3 は、本発明によって受信されたメッセージの概略流れ図である。ファイル転送のような大きい T C P / I P メッセージは、多数の分離された約 6 4 K B の転送でネットワークからホストに受信され得るが、それぞれが多数の約 1 . 5 K B のフレーム即ちネットワーク上の転送するためのパケットに分割されてもよい。Internet Packet Exchange (I P X) 上で Sequential Packet Protocol (S P X) 又は NetWare Core Protocol (N C P) を走らせる Novell 社の NetWare プロトコルスーツは、同様の方式で機能する。高速パスによって取り扱われ得るデータ通信の他の形態は、Transaction TCP (以下 T / T C P 又は T T C P という) であり、これは、幾つかのメッセージ初期化ダイアログを介してコネクションを開始させて後のメッセージでデータを転送するのではなく、初期のトランザクションリクエストでコネクションを開始させ、その後データを含む応答がそのコネクションに従って送られるような T C P のバージョンである。これらのプロトコルによって代表されるあらゆる転送において、各パケットは、従来通り転送されるデータの部分とともに、プロトコル層のそれぞれに対応するヘッダ及びそのメッセージの残りのパケット群に対するそのパケットを位置を決めるためのマーカを有する。

【 0 0 1 9 】

C P D によってネットワークからメッセージパケット即ちフレームが受信されたとき (4 7) 、ハードウェアの補助によりそれが初めに検査される。ここでは、様々な層のプロトコルの種類を決定し、関連するチェックサムを確認し、かつこれらの認識事項をステータスワードに要約する (5 7) 。これらのワードに含められるのは、そのフレームが高速パスデータフローの候補であるか否かの表示である。高速パス候補の選択 (5 9) は、そのホストが C P D によって取り扱われるそのメッセージコネクションから利益を得られるか否かに基いて行われ、そのパケットが例えば T C P / I P 又は S P X / I P X のような特定のプロトコルを表示するヘッダバイトを有しているか否かの判定が行われる。小さいパーセンテージで存在する高速パス候補でないフレームは、低速パスプロトコル処理のためにホストのプロトコルスタックに送られる (6 1) 。その後の各高速パス候補についてのネットワークプロセッサ作業は、T C P 又は S P X C C B のような高速パスコネクションがその候補に対して常に現存しているか否か、又はその候補が、例えば T T C P / I P トランザクションのためのもののような新たな高速パスコネクションをセットアップするために用いられ得るか否かを決定する。C P D によって提供される妥当性検査が、或るフレームが高速パスで処理されるか、低速パスで処理されるかに関わらず、処理の高速化をもたらし、誤りのない場合にのみ、検査されたフレームが、たとえ低速パス処理の場合でもホストの C P U によって処理される。

【 0 0 2 0 】

高速パスの候補となるものとしてC P Dハードウェアの補助によって決定された受信されたメッセージフレームは全て、ネットワークマイクロプロセッサ又はI N I C比較器回路によって調べられ(5 3)、それらがC P Dによって保持されたC C Bに一致するか否かが判定される。そのような一致を確認する際、C P Dは下位層のヘッダを取り除き、C P Dのダイレクトメモリアクセス(D M A)ユニットを用いて、残りのアプリケーションデータをフレームから直接その最終行先に送信する(6 9)。この操作は、例えばT C Pコネクションが既に存在しかつ行先バッファがネゴシエートされた場合には、メッセージパケットを受け取った直後に行われ得るが、初めに初期ヘッダを処理して、この転送のための最終行先アドレスの新たな組を取得することが必要なこともある。後者の場合には、C P Dは、行先アドレスを待つとともに後続のメッセージパケットをキューに入れ、次いでキューに入れられたアプリケーションデータをその行先にD M A転送する。

10

【 0 0 2 1 】

C C Bに一致しない高速パス候補を用いて、そのフレームを順次プロトコル処理のためにホストに送信(6 5)することによって、新たな高速パスコネクションをセットアップすることができる。この場合、ホストはそのフレームを用いてC C Bを生成し(5 1)、次いでC C BはC P Dにわたされて、そのコネクションにおける後のフレームを制御する。C C BはC P Dにキャッシュされ(6 7)、従来通りのソフトウェア階層処理を使用している場合には、処理された全てのプロトコルに対して適切な制御及び状態情報を有する。C C Bはまた、後の関連するメッセージパケット内に含められたアプリケーションレベルのデータを、直接的な使用に利用可能な形態でホストのアプリケーションに直接わたすのを容易にするために用いられる、転送毎の情報のための記憶空間も有している。C P Dは、ホストからそのコネクションのためのC C Bを受け取ったとき、コネクション処理の指揮を取る。

20

【 0 0 2 2 】

図4 Aにより明確に示されているように、リモートホスト2 2からメッセージパケットがネットワーク2 5を介して受け取られたとき、そのパケットはC P D3 0のハードウェア受信ロジック3 2に入り、これがヘッダ及びデータを検査し、ヘッダを構文解析してメッセージパケット及びステータスを特定するワードを生成し、ヘッダ、データ、及びワードをメモリ6 0に一時的に格納する。パケットの検査と同時に、受信ロジック3 2は、そのパケットが高速パス処理の候補であるか否かを前記ワードを用いて示す。図4 Aは、そのパケットが高速パス候補でない場合を示しており、この場合C P D3 0は、矢印5 6によって示されるように、ホストのC P Uによる処理のために、内部バスを通してメモリ6 0からデータリンク層3 6へ検査されたヘッダ及びデータを送信する。そのパケットは、データリンク層3 6、ネットワーク層3 8、トランスポート層4 0、及びセッション層4 2のホストのプロトコルスタック4 4によって処理され、次にパケットからのデータ(D) 6 3は、矢印6 5に示されるように、記憶機構3 5に送信され得る。

30

【 0 0 2 3 】

図4 Bは、C P Dの受信ロジック3 2が、例えばパケットのヘッダからそのパケットがT C P / I P、T T C P / I P、又はS P X / I P Xメッセージに属すること導き出すことによってそのメッセージパケットが高速パス処理の候補であることを決定する場合を示す。次に、C P D3 0におけるプロセッサ5 5は、高速パス候補を要約するワードがキャッシュ6 3に保持されたC C Bに一致するか否かを調べるためのチェックを行う。そのパケットに対して一致がみられないときは、C P Dは、ホストのプロトコルスタック4 4にメモリ6 0から検査されたパケットを処理のために送信する。ホストのスタック4 4は、そのパケットを用いて、パケットに関連するメッセージからのデータの行先を見出し、それを予約することを含むメッセージのコネクションコンテキストを生成することができ、そのコンテキストはC C Bの形態を取る。本実施例は、高速パス及び非高速パスの候補の両方の処理のために一個の専用ホストスタック4 4を用いるが、以下に説明する実施例では、高速パスの候補は、非高速パスの候補とは異なるホストスタックによって処理される。その初めのパケットから幾つかのデータ(D 1) 6 6が、矢印6 6で示すように、所望に

40

50

応じて記憶機構 35 における行先に送られ得る。次に C C D が、矢印 64 に示すように、キャッシュ 62 に保存される C P D 30 を送信する。例えば、T C P / I P に代表されるような伝統的なコネクションベースのメッセージの場合は、初めのパケットが、C C B が生成されて C P D 30 にわたされる前にホスト間で生ずるコネクション初期化ダイアログの一部であり得る。

【 0 0 2 4 】

ここで図 4 C を参照すると、初めのパケットと同一のコネクションからの後続のパケットが、C P D 30 によってネットワーク 25 から受信されたとき、パケットのヘッダ及びデータは、受信ロジック 32 によって検査され、そのヘッダは構文解析されて、メッセージパケットの要約及び対応する C C B を見つけるためのハッシュが生成され、その要約及びハッシュはワードに含められる。そのワードは、パケットとともにメモリ 60 に一時的に保存される。プロセッサ 55 は、キャッシュ 62 に保存された各 C C B とハッシュとの間の一致をチェックし、その一致を見出したときには、矢印 72 に示すように、高速パスを介して記憶機構 35 における行先に直接データ (D 2) 70 を送信し、セッション層 42、トランスポート層 40、ネットワーク層 38、及びデータリンク層 36 はバイパスさせる。メッセージからの残りのデータパケットは、D M A によって直接記憶機構に送信され得、C P U スタック 44 による比較的低速のプロトコル層処理及び反復されるコピー処理は回避される。

【 0 0 2 5 】

図 4 D は、例えば図 4 C に示されるような、それに対して高速パスがコネクションが確立されたメッセージが、C P D によって容易に取り扱われないパケットを有するときの稀な例を取扱うための手順を示す。この場合、そのパケットは、プロトコルスタックによって処理されるべく送信され、プロトコルスタックは、矢印 76 で示すように、C P D との制御ダイアログを介してキャッシュ 62 からのメッセージのために C C B にわたされ、そのメッセージの処理を引き継がせるために C P U に信号が送られる。次にプロトコルスタックによる低速パス処理によって、矢印 82 に示すように、そのパケットからのデータ (D 3) 80 が記憶機構 35 に送られることになる。一旦そのパケットが処理され、誤り状態が補正されると、C C B は制御ダイアログを介してキャッシュ 62 に戻され得、そのメッセージの後続のパケットからのペイロードデータは、同様に C P D 30 の高速パスを介して送信され得ることになる。従って、C P U 及び C P D は協働して、所定のメッセージが高速パスハードウェア処理に従って処理されるか、より一般的な C P U によるソフトウェア処理に従って処理されるのかを決定する。

【 0 0 2 6 】

リモートホスト 22 へ供給するためのホスト 20 からネットワーク 25 へのメッセージの送信も、図 5 に示すように、C P U を介した順次式プロトコルソフトウェア処理か、C P D 30 を介した高速化ハードウェア処理の何れかによって処理され得る。記憶機構 35 から C P U 28 によって選択されたメッセージ (M) 90 は、矢印 92 及び 96 によって示されるように、スタック 44 によって処理するためセッション層 42 に送信され得る。しかし、コネクションが存在し、かつ C P D 30 がそのメッセージのための適切な C C B を既に有している状況では、データパケットはホストスタック 44 をバイパスし、D M A によって直接メモリ 60 に送信され得る。このとき、プロセッサ 55 は、各データパケットに全ての適切なプロトコル層を含む 1 個のヘッダを加え、得られたパケットをリモートホスト 22 への送信のためネットワーク 25 に送る。この高速パス送信によって、たとえ 1 個のパケットのための処理でも非常に高速化することができ、より大きいメッセージの場合は大幅に高速化する。

【 0 0 2 7 】

高速パスコネクションのためのメッセージは現存しておらず、従って高速パス送信をガイドするための適切な制御及び状態情報と共に C C B を生成することによって利益が得られる。T C P / I P 又は S P X / I P X に代表される伝統的なコネクションベースのメッセージの場合、C C B は、コネクション初期化ダイアログの間に生成される。T T C P /

10

20

30

40

50

IPに代表されるような高速コネクションメッセージの場合、CCBはペイロードデータを送信する同じトランザクションで生成され得る。この場合、ペイロードデータの送信が高速バスコネクションをセットアップするために用いられたリクエストに対する応答であり得る。いずれの場合にも、CCBは、各プロトコル層に関連するプロトコル及びステータス情報を提供し、これには、ユーザが関与するもの及び1回の情報の転送のための記憶空間が含まれている。CCBは、プロトコルスタック44によって生成され、そのプロトコルスタックは、次に、矢印98によって示されるように、CPDのコマンドレジスタに書き込むことによってCCBをCPD30にわたす。CCBによってガイドされ、プロセッサ55は、ホストメモリ35における発信元からのデータのネットワークフレームサイズにされた一部分を、矢印99に示すように、DMAを用いてそれ自身のメモリ60に送る。次にプロセッサ55は、適切なヘッダ及びチェックサムをそのデータの一部分の先頭に付加し、得られたフレームを関連するプロトコルの規約に準拠させてネットワーク25に送信する。CPD30は、全てのデータがその行先に達したという確認応答を受信した後、CPDは応答バッファに書き込むことによってホスト35に通知する。

【0028】

従って、データ通信の高速バス送信はまた、ホストをCPUのフレーム毎の処理から解放させる。データ送信の大部分は、高速バスによってネットワークに送信され得る。入出力の両方が高速バスであることによって、上位層レベル即ちセッション層レベル以上における機能による割り込みが大幅な低減を達成され、かつネットワークマイクロプロセッサとホストとの間の相互作用は上位層が作ろうとする全転送サイズを用いて行われることになる。高速バス通信の場合、割り込みは(その大部分が)上位層メッセージトランザクション全体の始まりと終わりにおいてのみ生じ、各階層の部分及びそのトランザクションのパケットの送受信において割り込みは生じない。

【0029】

ホスト152に対するネットワークインタフェースを提供する、単純化されたインテリジェントネットワークインタフェースカード(INIC)150が図6に示されている。INIC150のハードウェアロジック171は、INICとホストとを接続する周辺バス(PCI)157でネットワーク155に接続されている。この実施例において、ホスト152はTCP/IPプロトコルスタックを有し、それはネットワーク155から受け取ったメッセージフレームの順次式ソフトウェア処理のための低速バス158を提供する。ホスト152のプロトコルスタックは、データリンク層160、ネットワーク層162、トランスポート層164、及びホスト152における通信データの発信元又は行先168を提供するアプリケーション層166を有する。図示されていない他の層、例えばセッション層やプレゼンテーション層もホストスタック152に含められ得、発信元又は行先はデータの性質に応じて変わり、実際にはアプリケーション層であり得る。

【0030】

INIC150は、ホストのプロトコルスタックを含む低速バス158に沿って処理するメッセージ、又はホストのプロトコルスタックをバイパスする高速159に沿って処理するメッセージとの間の選択を行うネットワークプロセッサ170を有する。受信されたパケットのそれぞれは、INIC150に含められたフライバイ(fly by)ハードウェアロジック171上で処理され、従って或るパケットの全てのプロトコルヘッダは、プロトコル層の間でデータのコピー、移動、記憶させることなく処理され得る。ハードウェアロジック171は、パケットバイトがそのハードウェアを通過するとき、選択されたヘッダバイトを分類することによって1回で所定のパケットのヘッダを処理する。選択されたバイトの処理の結果は、チェックサム検査を含むパケットの要約が生成されるまで、パケットのどの他のバイトを分類するかを決定することの助けとなる。次に受信されたパケットからのデータ及び処理されたヘッダが、そのパケットのヘッダ及びステータスを要約するワードとともにINICの記憶機構185に記憶される。ネットワーク記憶機構の構成のために、INIC150はIDE、SCSI、又は類似のインタフェースを有するディスクドライブのような周辺記憶装置に接続され得、その記憶装置のためのファイルキャッシ

ュは I N I C 1 5 0 のメモリ 1 8 5 上に存在する。そのようなネットワークインタフェースの幾つかは 1 つのホストのために存在し得、それぞれが関連する記憶装置を有する。

【 0 0 3 1 】

ネットワーク 1 5 5 から I N I C 1 5 0 によって受信されたメッセージパケットのハードウェア処理は、図 7 により詳細に示されている。受信されたメッセージパケットは、初めにメディアアクセスコントローラ 1 7 2 に入り、そのメディアアクセスコントローラはネットワークへの I N I C のアクセス及びパケットの受信を制御し、かつネットワークプロトコル管理のための統計的情報を提供する。そこから 1 回に 1 バイトのデータがアセンブリレジスタ 1 7 4 に流れる。アセンブリレジスタ 1 7 4 は、この実施例では 1 2 8 ビットの幅を有する。そのデータはフライバイシーケンサ 1 7 8 によって分類されるが、後に図 8 を参照してより詳細に説明するように、フライバイシーケンサは、パケットのバイトを、それらが移動している (fly by) ときに調べ、パケットを要約するために用いられるこれらのバイトからステータスを生成する。生成されたそのステータスはマルチプレクサ 1 8 0 によってデータと併合され、得られたデータは S R A M 1 8 2 に格納される。パケット制御シーケンサ 1 7 6 は、フライバイシーケンサ 1 7 8 を監視し、メディアアクセスコントローラ 1 7 2 からの情報を調べ、データのバイトをカウントし、アドレスを生成し、ステータスを送り、かつアセンブリレジスタ 1 7 4 から S R A M 1 8 2 、更には D R A M 1 8 8 へのデータの移動を管理する。パケット制御シーケンサ 1 7 6 は、S R A M コントローラ 1 8 3 を介して S R A M 1 8 2 におけるバッファを管理し、かつデータを S R A M 1 8 2 から D R A M 1 8 8 のバッファへ移動させる必要があるときに D R A M コントローラ 1 8 6 への表示も行う。一旦パケットのためのデータの移動が終了し、全てのデータが D R A M 1 8 8 のバッファに移動すると、パケット制御シーケンサ 1 7 6 は、フライバイシーケンサ 1 7 8 において生成されたステータスを、パケットデータの先頭に付加するべく D R A M 1 8 8 のバッファの先頭部及び S R A M 1 8 2 に動かす。次にパケット制御シーケンサ 1 7 6 は、受信バッファディスクリプタを受信キューに入れるようにキューマネージャに要求し、受信キューはプロセッサ 1 7 0 にパケットがハードウェアロジック 1 7 1 によって処理されたこと、及びその要約されたステータスを通知する。

【 0 0 3 2 】

図 8 は、フライバイシーケンサ 1 7 8 が幾つかの段を有していることを示し、この段のそれぞれは通常パケットヘッダの特定の部分を集中して取扱い、従ってその段に関連するステータスを生成するため特定のプロトコル層を集中して取り扱う。この実施形態では、フライバイシーケンサ 1 7 8 が、メディアアクセス制御シーケンサ 1 9 1 、ネットワークシーケンサ 1 9 2 、トランスポートシーケンサ 1 9 4 、及びセッションシーケンサ 1 9 5 を有する。より高位のプロトコル層に関連するシーケンサを追加的に提供することもできる。フライバイシーケンサ 1 7 8 は、パケット制御シーケンサ 1 7 6 及びフライバイシーケンサに所定のバイトがアセンブリレジスタ 1 7 4 から利用可能であるか否かを伝えるパケット制御シーケンサによる所定のポイントによってリセットされる。メディアアクセス制御シーケンサ 1 9 1 は、バイト 0 - 5 を見ることにより、或るパケットが、別のホストではなくホスト 1 5 2 を、又は別のホストに加えてホスト 1 5 2 をアドレス指定されているかを決定する。パケットのオフセット 1 2 及び 1 3 もメディアアクセス制御シーケンサ 1 9 1 によって処理され、例えばそのパケットがイーサネット又は 8 0 2 . 3 であるか否かのタイプフィールドを決定する。そのタイプフィールドがイーサネットである場合、それらのバイトはメディアアクセス制御シーケンサ 1 9 1 にそのパケットのネットワークプロトコルの種類も伝える。8 0 2 . 3 の場合は、それらのバイトは全フレームの長さを表示し、メディアアクセス制御シーケンサ 1 9 1 は更にパケット内の 8 個のバイトをチェックし、ネットワーク層の種類を決定する。

【 0 0 3 3 】

大部分のパケットについて、ネットワークシーケンサ 1 9 2 は、受信されたヘッダ長が正しい長さを有していることを検査し、ネットワーク層のヘッダを検査する。高速パス候補について、ネットワーク層のヘッダは、メディアアクセス制御シーケンサ 1 9 1 によって

10

20

30

40

50

行われた解析からIP又はIPXであることが判っている。例えば、タイプフィールドが802.3であり、ネットワークプロトコルがIPであると仮定すると、ネットワークシーケンス192はIPの種類を決定するためにバイト22で始まるネットワークヘッダの初めのバイトを解析する。IPヘッダの初めのバイトは、ネットワークシーケンス192によって処理され、そのパケットがどのIPの種類に関係しているかを決定する。そのパケットが、例えばIPバージョン4に関係するものであると決定することにより、ネットワークシーケンス192による更なる処理が方向付けられ、ネットワークシーケンスはパケットのトランスポートヘッダプロトコルの表示のためIPヘッダ内の10個のバイトに配置されたプロトコルのタイプを調べる。例えば、イーサネット上のIPの場合、そのIPヘッダはオフセット14から始まり、プロトコルタイプのバイトは、オフセット23であり、それはネットワークロジックによって処理され、トランスポート層プロトコルが例えばTCPであるか否かが決定される。通常20 - 40バイトであるネットワークヘッダの長さから、ネットワークシーケンス192が、トランスポート層ヘッダの検査のため、パケットのトランスポート層ヘッダの始まりを決定する。トランスポートシーケンス194は、トランスポート層ヘッダ及びデータのためのチェックサムを生成してもよく、トランスポート層ヘッダ及びデータは、少なくともTCPの場合にIPヘッダからの情報を含み得る。

【0034】

TCPパケットの例について続けると、トランスポートシーケンス194は、1つにはそのメッセージのためのTCP発信元及び行先ポートを決定するため、例えばそのパケットがNet Biosか他のプロトコルであるかを決定するために、ヘッダのトランスポート層部分における初めの数バイトの解析も行う。TCPヘッダのバイト12が、トランスポートシーケンス194によって処理され、TCPヘッダの長さが決定され、検査される。TCPヘッダのバイト13は、肯定応答(ack)フラグ及びプッシュフラグを除いて、例えばリセットや終了のような、プロセッサにこのパケットを例外として分類させ得る予期しないオプションを表示し得るフラグを含む。TCPオフセットバイト16及び17はチェックサムであり、それはハードウェアロジック171によって引き出されたり、格納されたりし、フレームの残りの部分はそのチェックサムに対して検査される。

【0035】

セッションシーケンス195は、セッション層ヘッダの長さを決定し、セッション層ヘッダは、Net Biosの場合には4バイトに過ぎず、そのうちの2バイトはNet Biosペイロードデータの長さを通知するが、他のプロトコルの場合にはより長いこともあり得る。セッションシーケンス195は、例えば読み出し又は書き込みとしてメッセージのタイプを分類するためにも用いられ得、そのために高速パスが特に有益なものとなり得る。メッセージの種類に応じて更なる上位層ロジックの処理が、フライバイシーケンス178及びパケット制御シーケンス176のハードウェアロジック171によって行われ得る。従って、ハードウェアロジック171はインテリジェントに1つのバイトストリームから選択されたバイトを分類することによってヘッダのハードウェア処理を誘導し、そのパケットのステータスが処理の進行中に決定された分類から構築される。一旦パケット制御シーケンス176が全てのパケットがフライバイシーケンス178によって処理されたことを検出すると、パケット制御シーケンス176は、フライバイシーケンス178によって生成されたステータス情報にパケット制御シーケンス176によって生成されたステータス情報を加え、プロセッサ170によるパケットの取扱いにおける便宜のため、そのパケットにそのステータス情報を先頭付加(パケットの先頭部分への付加)する。パケット制御シーケンス176によって生成された追加のステータス情報は、メディアアクセスコントローラ172のステータス情報及び発見された全ての誤り、又はアセンブリレジスタ若しくはDRAMバッファの何れかにおけるデータオーバーフロー、又はそのパケットに関連する他の種々雑多な情報を有する。パケット制御シーケンス176は、キューマネージャ184を介して受信バッファキュー及び受信統計的キューにエントリを入れることも行う。

10

20

30

40

50

【 0 0 3 6 】

ハードウェアロジック 171 によりパケットを処理することの利点は、従来のような順次ソフトウェアプロトコル処理とは異なり、そのパケットが各プロトコル層ヘッダの処理のために記憶機構に格納されたり、転送されたり、コピーされたり、又は引き出されたりする必要がなく、処理効率を劇的に向上させ、各パケットの処理時間を短縮できる点にある。そのパケットはビットがネットワークから受信されときのビットレート、例えば 100 BaseT 接続の場合 100 メガビット / 秒で処理され得る。従って、この速度で受信された 60 バイトの長さを有するパケットを分類する時間は、約 5 マイクロ秒である。ハードウェアロジック 171 でそのパケットを処理し、パケットデータを高速バスを介してそのホスト行先に送信するためにかかる全時間は、66 MHz の PCI バスを用いた場合、約 16 マイクロ秒以下であり得るが、300 MHz の PentiumII (登録商標) プロセッサによる従来通りのソフトウェアプロトコル処理では、使用中の装置において 200 マイクロ秒もの時間がかかり得る。従って、従来型の順次式ソフトウェアプロトコル処理を用いる高速 CPU の場合と比較して、高速バス 159 によって 1 桁以上の処理時間の短縮が達成され得、CPU の割り込みの低減やホストのバスの大域幅の節約によってもたらされる追加の時間短縮を考慮しなくても、ハードウェアロジック 171 及びプロセッサ 170 でプロトコルヘッダを処理することによって劇的な高速化が提供されることを立証している。

10

【 0 0 3 7 】

プロセッサ 170 は、記憶機構 185 に保持された受信されたメッセージパケットのそれぞれについて、そのパケットが高速バス 159 の候補であるか否かを選択し、そうである場合には、高速バスが既にそのパケットが属するコネクションに対してセットアップされているか否かを調べるためのチェックを行う。これを行うため、プロセッサ 170 は、初めにそのパケットヘッダが高速バス候補のために定義されたプロトコルのものであるか否かを決定するためにヘッダステータスの要約をチェックする。そうでない場合には、プロセッサ 170 は、INIC 150 における DMA コントローラに、そのパケットを低速バス 158 処理のためにホストに送信する命令を送る。メッセージが低速バス 158 処理される場合でも、INIC 150 は検査及びメッセージの種類の決定のような初期手順を実行し、検査されたメッセージを少なくともホストのデータリンク層 160 に送る。

20

【 0 0 3 8 】

高速バス 159 の候補の場合、プロセッサ 170 は、そのヘッダステータスの要約が INIC に保持された CCB に一致するか否かを調べるためのチェックを行う。そうである場合には、パケットからのデータは高速バス 159 を通してホストの行先 168 に送信される。高速バス 159 の候補のパケット要約が INIC に保持された CCB に一致しない場合には、そのパケットは低速バス処理のためにホスト 152 に送信されて、そのメッセージのための CCB が生成される。高速バス 159 を用いることは、断片化したメッセージの場合や他の複雑さのために必要でない、又は望ましいものではないこともある。しかし、大部分のメッセージについては、INIC の高速バス 159 は、メッセージ処理を非常に高速化し得る。従って INIC 150 は、従来のように所定のパケットの運命を決定するために幾つかのプロトコル層のそれぞれにおいてステートマシンを用いるのとは異なり、処理進行中に収集される情報に基いてデータをその行先に直接送信するか否かを決定する 1 個のステートマシンプロセッサ 170 を提供するものである。

30

40

【 0 0 3 9 】

ホスト 152 において受信された表示即ちパケットの処理において、ホストのプロトコルドライバは、その表示が高速バス又は低速バスの何れであるかに基いて処理経路を選択する。TCP/IP 又は SPX/IPX メッセージは、そのドライバによって形成され、高速バスパケットとの一致をとの一致をとり、かつそれをコネクション行先 168 に導くために、INIC に送られる。TCP/IP メッセージの場合、ドライバは、初期リクエストパケットの処理から、そのトランザクションのためのコネクションコンテキストを生成することができる。この処理では、メッセージ行先 168 の位置指定を行い、次にその

50

行先からの応答のための高速バスを提供するために C C B の形態で I N I C にそのコンテキストを送る。C C B は、プロトコル層及びメッセージのパケットに関するコネクション及び状態情報を含む。従って C C B は、発信元及び行先メディアアクセス制御 (M A C) アドレス、発信元及び行先 I P 又は I P X アドレス、発信元及び行先 T C P 又は S P X ポート、タイマーのような T C P 変数、ウィンドウプロトコルをスライドさせるための受信及び送信ウィンドウ、及びセッション層プロトコルを示す情報を有し得る。

【 0 0 4 0 】

I N I C におけるハッシュテーブルに C C B をキャッシュすることにより、そのパケットが高速バス 1 5 9 を介して処理され得るか否かを決定するべく入力されるパケットを要約するワードとの迅速な比較処理が行われるようになり、一方全ての C C B が処理のために I N I C にも保持される。この比較を高速化する他の方法としては、コンテンツのアドレス指定可能なメモリ (C A M) のようなハードウェア補助機構及び B ツリーのようなソフトウェア処理等が挙げられる。I N I C マイクロコード又は比較器回路が C C B との一致を検出すると、D M A コントローラは C P U による割り込み、プロトコル処理又はコピー処理を行うことなく、行先 1 6 8 にパケットからのデータをおく。受信したメッセージの種類に応じて、そのデータの行先は、ホスト 1 5 2 におけるセッション層、プレゼンテーション層、又はアプリケーション層、又はファイルバッファキャッシュであり得る。

【 0 0 4 1 】

図 9 は、ファイルサーバとして用いられるホスト 2 0 2 に接続された I N I C 2 0 0 を示す。この I N I C は、一般にファスト・イーサネットとして知られる 8 0 2 . 3 u 標準を用いる幾つかのネットワークコネクションのためのネットワークインタフェースを提供する。I N I C 2 0 0 は、P C I バス 2 0 5 によってサーバ 2 0 2 に接続されており、サーバ 2 0 2 は M A C 層 2 1 2、ネットワーク層 2 1 5、トランスポート層 2 1 7、及びアプリケーション層 2 2 0 を含む T C P / I P 又は S P X / I P X プロトコルスタックを有し、アプリケーション層の上に発信元 / 行先 2 2 2 が設けられているが、前に述べたように、アプリケーション層が発信元 / 行先であり得る。I N I C は、ネットワークライン 2 1 0、2 4 0、2 4 2、及び 2 4 4 にも接続されており、これらのネットワークラインは好ましくはファスト・イーサネット、ツイストペア、光ファイバー、同軸ケーブル、又は他の種類のラインであって、それぞれ 1 0 0 M b / s のデータ転送が可能なものであるが、より高速又はより低速のデータレートも可能である。ネットワークライン 2 1 0、2 4 0、2 4 2、及び 2 4 4 は、それぞれハードウェア回路の専用の行に接続されており、前記ハードウェア回路はそれぞれのネットワークラインから受信されたメッセージパケットの検査及び要約を行うことができる。従って、ライン 2 1 0 は、シーケンサ 2 5 0 の第 1 の水平方向の行に接続され、ライン 2 4 0 は、シーケンサ 2 6 0 の第 2 の水平方向の行に接続され、ライン 2 4 2 は、シーケンサ 2 6 2 の第 3 の水平方向の行に接続され、ライン 2 4 4 は、シーケンサ 2 6 4 の第 4 の水平方向の行に接続される。パケットがその水平方向のハードウェアの行の 1 つによって検査、要約された後、そのパケットはそのステータス要約とともに記憶機構 2 7 0 に格納される。

【 0 0 4 2 】

ネットワークプロセッサ 2 3 0 は、I N I C 2 0 0 に格納されたその要約及び任意の C C B との比較に基いて、ホストによって処理される低速バス 2 3 1 を通してそのパケットを送信するか否かを決定する。大部分のパケットは、そのような順次処理を避けることができ、一致する C C B に従ってサーバ内のデータ行先に高速バス 2 3 7 を通して D M A により送信されるそのデータ部分を有する。同様に、高速バス 2 3 7 は、データをネットワーク送信のためにパケットに分割し、全てのヘッダを付加してプロセッサ 2 3 0 によって発信元 2 2 2 から任意のネットワークラインへ直接データを送信する経路を提供し、同様に C P U による処理及び割り込みを最小限にする。明示のため、水平シーケンサ 2 5 0 のみがアクティブ状態として示されているが、実際にはシーケンサ行 2 5 0、2 6 0、2 6 2、及び 2 6 4 のそれぞれが、他の全てのシーケンサ行と同時に全二重通信を提供する。専用 I N I C 2 0 0 は、メッセージパケットとともに動作するとき、ソフトウェアプロトコ

ルスタックに従ってそれらのヘッダを順次処理する高度な汎用ホストCPUよりも非常に高速である。

【0043】

ファイル転送のような大きいメッセージのために最も一般的に用いられるネットワークプロトコルの1つは、TCP/IPを通してのサーバメッセージブロック(SMB)である。SMBはファイルが書き込まれるディスク又はプリンタのような特定の操作のために必要なリソースが、その操作が生成されるホストに存在するか、或いは結び付けられているか、又はファイルサーバのようなネットワークに接続された他のホストに存在するかを決定するリダイレクタソフトウェアとともに動作し得る。SMB及びサーバ/リダイレクタは、従来はトランスポート層によりサービスされるが、本発明ではSMB及びリダイレクタが、INICによりサービスされ得る。この場合には、大きいSMBトランザクションを受け取ったとき、INICバッファからDMAコントローラによってデータを送信することにより、ホストが取り扱わなければならない割り込みが著しく低減し得る。更に、このDMAは一般的にファイルデバースキャッシュにあるその最終的な行先にデータを転送する。本発明のSMB送信は、基本的に上述のSMBの受信を逆方向にたどるものであり、ホストからINICにデータが転送され、バッファに格納されるとともに、関連するプロトコルヘッダがネットワークラインを介してのリモートホストへの転送のためにINICにおけるデータの先頭に付加される。ホストの反復される割り込み無しに、カスタム設計のハードウェアを介しての複数のTCP/IP、NetBios、及びSMBプロトコル層及び複数のパケットのINICによる処理によって、SMBメッセージのネットワークラインへの送信速度を著しく高めることができる。

【0044】

図10に示すように、所定のメッセージがホスト202によって処理されるか、又はINIC200によって処理されるかを制御するために、メッセージコマンドドライバ300がホスト202にインストールされ、ホストプロトコルスタック310と協働して機能を果たす。コマンドドライバ300は、メッセージの受信又は送信に介入し、CCBを生成し、INIC200からCCBを送受信することができ、従って処理能力の向上を除いて、INICの機能はユーザからは透過的なものとなる。また、図面にはINICメモリ304及びINICミニポートドライバ306が示されており、INICミニポートドライバはネットワーク210から受信したメッセージパケットを、そのパケットが高速パス候補としてのラベル付けされているか否かに応じて、従来通りのプロトコルスタック310又はコマンドプロトコルスタック300の何れかへ誘導することができる。従来通りのプロトコルスタック310は、高速パス候補としてラベル付けされておらず、従ってコマンドスタック300によって処理されないメッセージの従来通りの下位層の処理のために、データリンク層312、ネットワーク層314、及びトランスポート層316を有する。下位層スタック310の上に存在するのは上位層318であり、これはやりとりされるメッセージに応じて、セッション層、プレゼンテーション層、及び/又はアプリケーション層を表す。コマンドドライバ300は、同様にデータリンク層320、ネットワーク層322、及びトランスポート層325を有する。

【0045】

ドライバ300は、メッセージをネットワーク210に送信するために、上位層318から送信されたメッセージが、コマンドスタック300、それに続けてINIC高速パスによって処理されるか、従来通りのスタック310によって処理されるかを決定する上位層インタフェース330を有する。その上位層インタフェース330が、従来通りホストのプロトコルスタックによるプロトコル処理の後のネットワークへの送信を目的とする上位層318から適切なメッセージを受け取ったとき、そのメッセージはドライバ300に送られる。次にINICが、INIC DMAユニットを介しての送信のために、メッセージデータのネットワークサイズにされた一部分を得て、そのデータの一部分の先頭にヘッダを付加し、形成されたメッセージパケットをワイヤを通して送信する。逆に、TCP、TTCP、SPX、又は類似のメッセージパケットを、高速パスコネクションのセットア

ップにおいて用いられるようにネットワーク 210 から受け取ったとき、ミニポートドライバ 306 は、そのメッセージパケットをコマンドドライバ 300 に処理するために転送する。ドライバ 300 は、そのメッセージパケットを処理してそのメッセージのためのコンテキストを生成し、ドライバ 302 は、そのコンテキスト及びコマンド命令を、高速バスを通しての同一のコネクションにおいて後のメッセージのデータを送信するための CCB として、INIC 200 に戻す。INIC によって、数百の TCP、TCP、SPX、又は類似の CCB コネクションは、その数に限界なく維持され得るが、INIC キャッシュが一杯になった場合のため、最小使用頻度 (LRU) アルゴリズムが用いられる。ドライバ 300 は、CCB として INIC 200 に送られる TCP リクエストのためのコネクションコンテキストも生成することができ、これによってそのリクエストに対する TCP 応答の高速バス転送が可能となる。高速化されていないプロトコルを有するメッセージは、プロトコルスタック 310 によって従来通りに処理され得る。

10

【0046】

図 11 は、Microsoft (登録商標) プロトコルメッセージのためのコマンドドライバソフトウェアの TCP/IP インプリメンテーションを示す。従来通りのホストプロトコルスタック 350 は、MAC 層 353、IP 層 355、及び TCP 層 358 を有する。コマンドドライバ 360 は、ホストスタック 350 と協働して働き、ネットワークメッセージを処理する。コマンドドライバ 360 は、MAC 層 363、IP 層 366、及び Alacritec TCP (ATCP) 層 373 を有する。従来型スタック 350 及びコマンドドライバ 360 は、ネットワークドライバインタフェース仕様 (NDIS) 層 375 を共有しており、これは INIC ミニポートドライバ 306 を相互作用する。INIC ミニポートドライバ 306 は、従来通りのホストスタック 350 又は ATCP ドライバ 360 の何れかによって処理するため受信した表示をソートする。TDI フィルタドライバ及び上位層インタフェース 380 は、同様に、TDI ユーザ 382 からネットワークに送られたメッセージが、コマンドドライバに送られて INIC の高速バスに転送されたか、或いはホストスタックによって処理されたかを判定する。

20

【0047】

図 12 は、典型的なクライアント 190 とサーバ 290 との間の SMB 交換を示す図であり、クライアントとサーバの両方は本発明の通信装置を有し、その通信装置はそれぞれデータの高速バス転送のためのそれらのコネクションを確定する CCB を有する。クライアント 190 は、INIC 150、802.3 準拠のデータリンク層 160、IP 層 162、TCP 層 164、Net Bios 層 166、及び SMB 層 168 を有する。クライアントは、通信処理のための低速バス 157 及び高速バス 159 を有する。同様にサーバ 290 は、INIC 200、802.3 準拠のデータリンク層 212、IP 層 215、TCP 層 217、Net Bios 層 220、及び SMB 層 222 を有する。サーバはネットワークライン 240、242、及び 244 に接続されているとともに、クライアント 190 に接続されたライン 210 にも接続されている。サーバも通信処理のための低速バス 231 及び高速バス 237 を有する。

30

【0048】

クライアント 190 がサーバ 290 上の 100 KB のファイルの読み出しを行おうとしていると仮定すると、クライアントは、ネットワーク 210 を介して読み出しブロック行 (RBR) SMB を送信することによって、サーバ 290 上のそのファイルの初めの 64 KB を要求することを開始し得る。RBR コマンドは、例えば 76 バイトに過ぎない大きさであり得、従ってサーバ上の INIC 200 は、メッセージタイプ (SMB) 及び比較的小さいメッセージサイズを認識し、かつその 76 バイトを高速バスを介してサーバの Net Bios に直接送信する。Net Bios は、SMB にデータを渡し、SMB は読み出しリクエストを処理し、データの 64 KB を取り出して、サーバのデータバッファに入れる。次に SMB は、Net Bios を呼出してそのデータを送信し、Net Bios はそのデータをクライアントのために出力する。従来型のホストでは、Net Bios が TCP 出力を呼出し、64 KB を TCP に渡し、TCP はそのデータを 1460 バイトの

40

50

セグメントに分割し、各セグメントをIP及び最終的にMAC（低速パス231）を介して出力する。本発明の場合には、その64KBのデータは、クライアント-サーバSMBコネクションに関する表示とともにATCPドライバに向かい、前記表示はINICによって保持されているCCBを示す。次にINIC200は、ホストバッファから1460バイトのセグメントをDMA送信し、TCP、IP及びMACのための適切なヘッダを一度に付加し、かつ完成したパケットをネットワーク210（高速パス237）上に送信する。INIC200は、この処理を64KBの転送の全てが送信されるまで反復する。通常、クライアントからその64KBのを受信したという確認応答を受信した後、INICは残りの36KBも高速パス237によって送信する。

【0049】

INIC150は、その応答が着信したとき、クライアント190上で動作しており、INIC150は、受信された第1のフレームからそのコネクションが受信高速パス159処理（TCP/IP、NetBios、CCBとの一致）であることを認識し、ATCPは、その第1フレームを用いてそのメッセージのためのバッファ空間を確保し得る。後者の場合には、フレームのNetBios部分の初めの128バイトをATCP高速パスを介して、ホストのNetBiosに直接わたすことによって達成され、これによりNetBios/SMBに全てのフレームのヘッダがわたされる。NetBios/SMBは、これらのヘッダを解析し、リクエストIDとの一致によってこれが元の行読み出しコネクションに対する応答であることを認識し、データを置くべき64KのバッファのリストをATCPにわたす。この段階において、ただ1個のフレームが着信しているが、この処理が生じている間に複数のフレームが着信し得る。クライアントバッファリストがATCPにわたされると直ぐに、それはその転送情報をINIC150にわたし、INIC150はそれらのバッファに蓄積された任意のフレームデータのDMA送信を開始する。

【0050】

図13は、INIC200を単純化して示した図であり、INIC200はネットワークインタフェースコントローラとプロトコルプロセッサの機能を一個のASICチップ400に結合したものである。この実施形態におけるINIC200は、全二重4チャンネル10/100メガビット/秒（Mbps）のインテリジェントネットワークインタフェースコントローラを提供し、このコントローラはサーバアプリケーションのための高速プロトコル処理のために設計されている。特にサーバの用途のために設計されているが、INIC200は、パーソナルコンピュータ、ワークステーション、ルータ、又は他のTCP/IP、TTCP/IP、又はSPX/IPXプロトコルが用いられているあらゆる他のホストに接続され得る。

【0051】

INIC200は、4つのネットワークライン210、240、242、及び244に接続されており、これらのラインは例えばツイストペア、同軸ケーブル、又は光ファイバーのような複数の異なる経路に沿ってデータを転送することができ、そのコネクションはそれぞれ市販の物理層チップ、例えばSEEQ Technology Incorporated, 47200 Bayside Parkway, Fremont, CA94538製のmodel 80220/80221 Ethernet Media Interface Adapterを介してのメディア独立インタフェース（MII）を提供する。そのラインは、好ましくは802.3準拠であり、INICとともに4つの完全なイーサネットノードを構成し、INICは10Base-T、10Base-T2、100Base-TX、100Base-FX、及び100Base-T4とともに将来のインタフェース標準をサポートしている。物理層の特定及び初期化は、ホストドライバ初期化ルーチンによって達成される。ネットワークライン210、240、242、及び244とINIC200との間のコネクションは、MACユニットのMAC-A 402、MAC-B 404、MAC-C 406、及びMAC-D 408によって制御され、これらのMACユニットはINICがネットワークライン210、240、242、及び244にアクセスするときに必要な制御を行うMACサブレイヤの基本的な機能を実行するための論理回路を含む。MACユニット402～408は、無差別でマルチキャスト又はユニキャストモードで動作すること

10

20

30

40

50

ができ、これによってINICがネットワークモニタとして機能し、ブロードキャスト及びマルチキャストパケットを受信し、かつ各ノードに対して複数のMACアドレスをインプリメントすることが可能になる。MACユニット402～408は、シンプルネットワーク管理プロトコル(SNMP)のために用いられ得る統計的情報も提供する。

【0052】

MACユニット402、404、406、及び408は、それぞれ送受信シーケンサ、XMT&RCV-A418、XMT&RCV-B420、XMT&RCV-C422、及びXMT&RCV-D424にそれぞれ配線410、412、414、及び416によって接続されている。各送受信シーケンサは、メッセージフレームがそのシーケンサを通過するとき、処理進行中に幾つかのプロトコル処理ステップを実行し得る。MACユニットと組み合わせ、この送受信シーケンサ418～422は、ハードウェア上のデータリンク層、ネットワーク層、トランスポート層、セッション層、及び使用する場合にはプレゼンテーション層及びアプリケーション層のプロトコルのためにパケットステータスをコンパイルして、従来型の順次式ソフトウェアエンジンと比較してそのようなプロトコル処理のための時間を非常に短縮させることができる。送受信シーケンサ410～414は、ライン426、428、430、及び432によってSRAM及びDMAコントローラに444に接続されており、DMAコントローラ444はDMAコントローラ438及びSRAMコントローラ442を有する。スタティックランダムアクセスメモリ(SRAM)バッファ440は、ライン441によってSRAMコントローラ442に接続されている。SRAM及びDMAコントローラ444は、ライン446を介して外部メモリ制御機構450と相互作用し、外部メモリバス455を介してダイナミックランダムアクセスメモリ(DRAM)バッファ460とフレームの送受信を行う。DRAMバッファ460は、ICチップ400に隣接して配置される。DRAMバッファ460は、4MB、8MB、16MB、又は32MBとして構成され得、かつ所望に応じてチップ上に配置され得る。SRAM及びDMAコントローラ444はライン464を介してPCIバスインタフェースユニット(BIU)468に接続されており、BIU468はINIC200とPCIインタフェースバス257との間のインタフェースを管理する。64ビットの多重化BIU468は、スレーブ及びマスター機能の両方についてPCIバス257に対する直接のインタフェースを提供する。INIC200は、64ビット又は32ビット何れかのPCI環境において動作し得、何れの構成においても64ビットのアドレス指定をサポートしている。

【0053】

マイクロプロセッサ470はライン472によってSRAM及びDMAコントローラ444に接続され、かつライン475を介してPCI BIU468に接続されている。マイクロプロセッサ470はチップ制御ストア480に存在するレジスタファイルに命令を送り、チップ制御ストアはSRAMの書き込み可能なオンチップ制御ストア(WCS)及び読み出し専用メモリ(ROM)を有し、かつライン477によってマイクロプロセッサに接続されている。マイクロプロセッサ470は入力されるフレームを処理し、ホストのコマンドを処理し、ネットワークトラフィックを誘導し、かつPCIバストラフィックを誘導し得る、プログラム可能なステートマシンを提供する。3つのプロセッサは、各クロックサイクル毎に1個の命令を発して完了させる3つのレベルのパイプライン型アーキテクチャに、共通のハードウェアを用いてインプリメントされている。受信プロセッサ482は主として通信の受信のために用いられ、送信プロセッサ484は主として通信の送信のために用いられて、全二重通信を容易に行えるようにしており、一方ユーティリティプロセッサ486はPCIレジスタアクセスの監視及び制御を含む様々な機能を提供する。

【0054】

3つのプロセッサ482、484、及び486に対する命令は、オンチップ制御ストア480に存在する。このため3つのプロセッサの機能は容易に定義しなおすことができ、従ってマイクロプロセッサ470は所定の環境に適合させることができる。例えば受信機能のために必要な処理の量は、送信機能及びユーティリティ機能の何れかのために要求され

10

20

30

40

50

るものを上回っていることがある。この状況では、幾つかの受信機能が送信プロセッサ 484 及び / 又はユーティリティプロセッサ 486 によって実行され得る。或いは、追加のレベルのパイプラインを生成して、3 つではなく 4 以上の仮想プロセッサを生成し、その追加のレベルを受信機能専用にすることができる。

【0055】

INIC200 は、この実施形態では DRAM460 におけるテーブルに維持されている最大 256 個の CCB をサポートし得る。しかし、SRAM440 において、順次検索を保存するためのハッシュ順の CCB インデックスも存在する。一旦ハッシュが生成されると CCB は SRAM にキャッシュされ、この実施例では最大 16 個の CCB が SRAM にキャッシュされる。SRAM にキャッシュされた 16 の CCB の割り当ては、以下に説明するように、最低使用頻度レジスタによって取り扱われる。これらのキャッシュ位置は、送信プロセッサ 484 と受信プロセッサ 486 との間で共有されており、従ってより重い負荷のかかるプロセッサがより多くのキャッシュバッファを使用することができる。また、シーケンサ間で共有される 8 個のコマンドバッファ及び 8 個のヘッダバッファも存在する。所定のヘッダ又はコマンドバッファは、所定の CCB バッファに静的にリンクされていない。そのリンクはフレーム毎に動的だからである。

【0056】

図 14 は、パイプライン型マイクロプロセッサ 470 の概要を示しており、このマイクロプロセッサでは受信、送信、及びユーティリティプロセッサのための命令が、クロックのインクリメント I、II、III に従って 3 つの交互のフェーズで実行され、そのフェーズはパイプライン段階のそれぞれに対応している。各フェーズは異なる機能を担っており、3 つのプロセッサのそれぞれは、各クロックインクリメントの間に異なるフェーズを占める。通常各プロセッサは、制御ストア 480 からの異なる命令ストリームに基いて動作し、それぞれそのフェーズを通してそれ自身のプログラムカウンタ及びステータスを有する。

【0057】

一般に、パイプライン型マイクロプロセッサの第 1 命令フェーズ 500 は、命令を終了し、その結果を行先のオペランドに格納し、次の命令を取り出し、かつその次の命令を命令レジスタに格納する。第 1 レジスタセット 490 は、命令レジスタを含む複数のレジスタを提供し、第 1 レジスタセットのための制御機構セット 492 は第 1 レジスタセット 490 に対する記憶機構のための制御を提供する。幾つかの項目が制御機構 492 によって変更されることなく第 1 フェーズを通過し、代わりに第 1 レジスタセット 490 又は RAM ファイルレジスタ 533 に単にコピーされる。第 2 命令フェーズ 560 は、第 1 レジスタセット 490 の命令レジスタに格納された命令をデコードし、生成された全てのオペランドを集める命令デコーダ及びオペランドマルチプレクサ 498 を有し、その収集されたオペランドは次に第 2 レジスタセット 496 のデコードレジスタに格納される。第 1 レジスタセット 490、第 2 レジスタセット 496、及び第 3 命令フェーズ 600 で用いられる第 3 レジスタセット 501 は、多くの同一のレジスタを有し、これは図 15A 乃至図 15C により詳細に示されている。命令デコーダ及びオペランドマルチプレクサ 498 は、RAM ファイルレジスタ 533 のデータポート及び 2 つのアドレスから読み出しを行うことができ、RAM ファイルレジスタは、第 1 フェーズ 500 及び第 2 フェーズ 560 の両方において動作する。プロセッサ 470 の第 3 フェーズ 600 は、第 2 レジスタセットからのオペランドに対する ALU 処理を行う論理演算装置 (ALU) 602 を有し、ALU は処理の結果を第 3 レジスタセット 501 に含められた結果レジスタに格納する。スタック交換機構 608 は、レジスタスタックを順序付けしなおすことができ、キューマネージャ 503 はプロセッサ 470 のためのキューを編成することができ、その結果は第 3 レジスタセットに格納される。

【0058】

命令は、環状パイプライン 505 によって示されるように第 1 フェーズ、次にそれに続く第 3 フェーズをまで継続する。所定のフェーズ内での組み合わせによる遅れを最小限する

10

20

30

40

50

ために、命令実行の3つのフェーズにわたって様々な機能が分散されていることに注意されたい。この実施形態の66MHzの周波数で、各クロックインクリメントが終了するまでに15ナノ秒の時間がかかり、3つのプロセッサのそれぞれに対して1個の命令が完了するまで全体で45ナノ秒かかることになる。循環する命令フェーズは図15A乃至図15Cにより詳細に示されており、各フェーズが異なる図面において示されている。

【0059】

詳述すると、図15Aは第1レジスタセット490及び関連する制御機構492を有する第1フェーズ500の幾つかの特定のハードウェア機能を示す。第2レジスタセット492のための制御機構には、SRAM制御機構502が含まれ、これはアドレスをロードし、データをSRAMアドレス及びデータレジスタ520に書き込むための論理制御機構である。従って第3フェーズ600からのALU602の出力は、SRAM制御機構502によってSRAMアドレス及びデータレジスタ520のアドレスレジスタ又はデータレジスタに置かれ得る。ロード制御機構504は、同様にファイルコンテキストレジスタ522に或るファイルに対するコンテキストの書き込みのための制御を提供し、別のロード制御機構506は様々なデータをフリップフロップレジスタ525に格納するための制御を提供する。例えば、送られたビットがセットされたか否かを示すALU条件コードは、第1フェーズ500において実行される処理なしで、ALU条件コードレジスタ528にクロックに応じて入力される。フラグデコード508は、例えばロックの設定のようなフラグレジスタ530に格納された様々な機能を実行することができる。

【0060】

RAMファイルレジスタ533は、アドレス及びデータのための1つの書き込みポート及びアドレス及びデータのための2つの読み出しポートを有し、従って2以上のレジスタが一度に読み出され得る。上述のように、RAMファイルレジスタ533は基本的に第1及び第2フェーズをまたいでおり、即ち第1フェーズ500において書き込みが行われ、第2フェーズ560において読み出しが行われる。制御ストア命令510によって制御ストア480から新たなデータが入ってくるため、図面には示されていないが、命令レジスタ535に格納された命令をプロセッサがプログラムしなおすことが可能となる。このためのアドレスは、フェッチアドレスレジスタ538に格納されたアドレスのうちどのアドレスを取り出すかを決定するフェッチ制御レジスタにおいて生成される。ロード制御機構515は、プログラムカウンタのための命令を供給し、プログラムカウンタは制御ストアのためのフェッチアドレスの場合と概ね同様に動作する。3つのレジスタの後入先出スタック544は、このフェーズにおいて他の操作を受けずに第1レジスタセットにコピーされる。最後に、デバッグアドレス548のためのロード制御機構517が所望に応じて含められるが、これによって生じ得る誤りの訂正が可能となる。

【0061】

図15Bは、RAMファイルレジスタ533からのアドレス及びデータの読み出しを含む、第2マイクロプロセッサフェーズ560を示す。スクラッチSRAM565には、第3フェーズにインクリメントされるまで初めの2つのフェーズを通過するレジスタを有する第1レジスタセットのSRAMアドレス及びデータレジスタ520から書き込みが行われる。スクラッチSRAM565は、命令デコーダ及びオペランドマルチプレクサ498によって読み出される。これは、上述のように、スタック544、デバッグアドレス548、及びSRAMアドレス及びデータレジスタのような例外を除いて、第1レジスタセットの大部分のレジスタと同様である。命令デコーダ及びオペランドマルチプレクサ498は、セット490の様々なレジスタ及びSRAM565を参照して命令をデコードし、次のフェーズにおける処理のためのオペランドを収集する。特に、以下のALU602に提供するオペランドを決定する。命令デコーダ及びオペランドマルチプレクサ498の出力は、ALUオペランド579及び582、ALU条件コードレジスタ580、及びこの実施形態では32のキューを制御し得るキューチャネル及びコマンド587レジスタを含む第2レジスタセット496における多数のレジスタに格納される。セット496における幾つかのレジスタは、プログラム制御機構590、リテラルフィールド589、テスト選択5

10

20

30

40

50

84、及びフラグ選択585を含む、デコーダ498によって実質的にデコードされることなく命令レジスタ535から直接ロードされる。第1フェーズ500のファイルコンテキスト522のような他のレジスタは、第2フェーズ560のファイルコンテキスト577に常に格納されるが、マルチプレクサ572によって集められたオペランドとしても取り扱われる。スタックレジスタ544は、スタックレジスタ594に単にコピーされる。プログラムカウンタ540は、このフェーズにおいてインクリメントされ(568)、レジスタ592に格納される。同様に、オプションのデバッグアドレス548もインクリメントされ(570)、ロード制御機構575は、この時点で、各フェーズにおける誤り制御を可能にするためにパイプライン505から供給され得、その結果はデバッグアドレス598に格納される。

10

【0062】

図15Cは、ALU及びキュー操作を含む、第3のマイクロプロセッサフェーズ600を示す。ALU602は、加算器、優先エンコーダ、及び他の標準的な論理機能を有する。ALUの結果は、レジスタであるALU出力618、ALU条件コード620、及び先行オペランド結果622に格納される。ファイルコンテキストレジスタ616、フラグ選択レジスタ626、及びリテラルフィールドレジスタ630は、単に前のフェーズ560からコピーされる。テストマルチプレクサ604は、条件ジャンプがジャンプを生ずるか否かを決定するために設けられ、その結果はテスト結果レジスタ624に格納される。テストマルチプレクサ604は、フェッチ制御機構511のような類似の判断とともに、第1フェーズ500において実行され得る。スタック交換608は、スタック594からプログラムカウンタをフェッチすることにより、又はプログラムカウンタをそのスタックに入れることによってスタックを上限に移動させ、その結果はプログラム制御機構634、プログラムカウンタ638、及びスタック640レジスタに格納される。SRAMアドレスは所望に応じてこのフェーズ600においてインクリメントされ得る。別のデバッグアドレス642のための別のロード制御機構610は、このフェーズにおいても誤り制御を可能にするために、この時点でパイプライン505から強制的にわたされる。図面では一体に示されているQRAM及びQALU606は、キューチャンネル及びコマンドレジスタ587から読み出され、SRAMに格納され、キューを再編成し、必要に応じてデータ及びポインタを加えたり取り除いたりして、データのキューを管理し、その結果をテストマルチプレクサ604、及びキューフラグ及びキューアドレスレジスタ628に送る。従って、QRAM及びQALU606は、3つのプロセッサのためのキューを管理する任務、ソフトウェアによってCPU上で順次実行される従来通りのタスクを想定しており、キューマネージャ606は、代わりに高速化した実質的に同時並行のハードウェアキューイングを提供する。

20

30

【0063】

図16は、キューマネージャ606によって管理される32のハードウェアキューの2つを示しており、そのキューのそれぞれは、SRAMヘッド、SRAMテール、及びDRAMボディにおいて情報をキューに入れる能力を有しており、これによって各キューの拡張及び個別の編成が可能となる。従ってFIFO700は、SRAM記憶ユニット705、707、709、及び711を有し、それぞれが全部で32個のバイトに対して8個のバイトを含んでいるが、これらのユニットの数及び処理能力は、実施形態によって変えることができる。同様に、FIFO702はSRAM記憶ユニット713、715、717、及び719を有している。SRAMユニット705及び707は、FIFO700のヘッドであり、ユニット709及び711はFIFOのテールであり、ユニット713及び715はFIFO702のヘッドであり、ユニット717及び719はそのFIFOのテールである。FIFO700の情報は、矢印722で示すように、ヘッドユニット705又は707に書き込まれ得、矢印725で示すように、テールユニット711又は709から読み出され得る。しかし特定のエントリは、ヘッドユニット705又は707から読み書きの両方が行われ得、又はテールユニット709又は741から読み書きの両方が行われ得、データ転送及びラテンシーが最小限にされている。同様に、FIFO702のため

40

50

の情報は、一般的に、矢印 7 3 3 で示すように、ヘッドユニット 7 1 3 又は 7 1 5 に書き込まれ、矢印 7 3 9 で示すように、テールユニット 7 1 7 又は 7 1 9 から読み出されるが、それが書き込まれた同一のヘッド又はテールユニットから読み出されることもある。

【 0 0 6 4 】

S R A M F I F O 7 0 0 及び 7 0 2 は、共に D R A M 4 6 0 に接続されており、この D R A M 4 6 0 によって、これらの F I F O を無限に拡張し S R A M ヘッド及びテールが一杯になった状況を取り扱うことが可能となる。例えば、3 2 個のキューのうちの、Q (0) とラベル付けされた第 1 のキューは、矢印 7 2 7 で示すように、F I F O 7 0 0 のヘッド又はテールにキューを入れる代わりにキューマネージャの指示の下で動作する D M A ユニットによって D R A M 4 2 0 においてエントリをキューに入れることができる。D R A M 4 6 0 に格納されたエントリは、矢印 7 3 0 で示されるように、S R A M ユニット 7 0 9 に戻され、その F I F O の長さ及び終了時間が延長される。S R A M から D R A M への転換は、一般的に S R A M が一杯になったときのために準備されている。D R A M はより低速であり、D M A 転送によって更なるラテンシーが生ずるからである。従って Q (0) は、F I F O 7 0 0 及び D R A M 4 6 0 の両方におけるキューマネージャ 6 0 0 によって格納されたエントリを含み得る。同様に、F I F O 7 0 2 に向かう情報は、例えば Q (2 7) に対応するものであり得、矢印 7 3 5 で示すように、D M A によって D R A M 4 6 0 に転送され得る。より低速の D R A M 4 6 0 にも関わらず、費用的に有効にキューを作成する能力は、初期化の間にユーザが決定でき、これにより所望に応じてキューのサイズを変更することが可能となる。D R A M 4 6 0 にキューイングされた情報は、矢印 7 3 7 に示すように、S R A M ユニット 7 1 7 に戻される。

【 0 0 6 5 】

3 2 のハードウェアキューのそれぞれのステータスは、従来通り図 1 7 に示すように、4 つの 3 2 ビットのレジスタのセット 7 4 0 に維持され、そこからアクセスされる。各レジスタにおける特定のビットは特定のキューに対応する。そのレジスタは、Q - O u t _ R e a d y 7 4 5、Q - I n _ R e a d y 7 5 0、Q - E m p t y 7 5 5、及び Q - F u l l 7 6 0 とラベル付けされている。特定のビットが Q - O u t _ R e a d y レジスタ 7 5 0 にセットされた場合、そのビットに対応するキューは読み出される状態となった情報を含み、Q - I n _ R e a d y 7 5 2 レジスタにおける同一のビットがセットされることは、そのキューが書き込み可能な状態となったことを意味する。同様に Q - E m p t y レジスタ 7 5 5 における特定のビットを正値にセットすることは、そのビットが空になったことに対応するキューを意味し、Q - F u l l レジスタ 7 6 0 の特定のビットが正値にセットされることは、そのビットが一杯になったことに対応するキューを意味する。従って Q - O u t _ R e a d y 7 4 5 は、ビット (2 7) 7 5 2、ビット (2 8) 7 5 4、ビット (2 9) 7 5 6、及びビット (3 0) を包含するビット (0) 7 4 6 ~ ビット (3 1) 7 4 8 を含む。Q - I n _ R e a d y 7 5 2 は、ビット (2 7) 7 6 6、ビット (2 8) 7 6 8、ビット (2 9) 7 7 0、及びビット (3 0) 7 7 2 を含むビット (0) 7 6 2 ~ ビット (3 1) 7 6 4 を含む。Q - E m p t y 7 5 5 は、ビット (2 7) 7 7 8、ビット (2 8) 7 8 0、ビット (2 9) 7 8 2、及びビット (3 0) 7 8 4 を含むビット (0) 7 7 4 ~ ビット (3 1) 7 7 6 を含む。Q - F u l l 7 6 0 は、ビット (2 7) 7 9 0、ビット (2 8) 7 9 2、ビット (2 9) 7 9 4、ビット (3 0) 7 9 6 を含むビット (0) 7 8 6 ~ ビット (3 1) 7 8 8 を含む。

【 0 0 6 6 】

F I F O 7 0 0 に対応する Q (0) は、フリーのバッファキューであり、全ての利用可能なバッファであるアドレスのリストを保持する。このキューは、マイクロプロセッサ又は他のデバイスがフリーバッファアドレスを必要とするときにアドレス指定され、従って共通に感知できる D R A M 4 6 0 を含む。従って、フリーバッファアドレスを必要とするデバイスは、Q (0) をチェックし、そのアドレスを得る。F I F O 7 0 2 に対応する Q (2 7) は、受信バッファディスクリプタキューである。受信シーケンサによって受信されたフレームが処理された後、そのシーケンサはそのフレームのためのディスクリプタを Q

(27)に格納しようとする。このようなディスクリプタの位置がSRAMにおいて即座に利用可能である場合は、Q-In_Ready 750のビット(27)766がセットされる。そうでない場合は、そのシーケンサはキューマネージャがSRAMからDRAMへDMA転送を開始させて、受信ディスクリプタを格納するための空間を空けるのを待たなければならない。

【0067】

SRAMとプロセッサとの間でのキューエントリの移動、受信及び送信シーケンサ間のキューエントリの移動、及びSRAMとDRAM間のキューエントリの移動を管理するキューマネージャの動作は、図18により詳細に示されている。キューを利用するリクエストとしては、プロセッサリクエスト802、送信シーケンサリクエスト804、及び受信シーケンサリクエスト806がある。そのキューのための他のリクエストは、DRAMからSRAMへのリクエスト808及びSRAMからDRAMへのリクエスト810であり、これはキューのDRAM及びSRAMヘッド又はテール間でデータをやりとりする際にキューマネージャの代わりに作用する。これらの様々なリクエストの何れが次のサイクルにおいてキューマネージャを用いるかの決定は、優先ロジックアービタ815によって取り扱われる。高い周波数での動作を可能にするため、キューマネージャはパイプライン型で、レジスタA818及びレジスタB820が一時的な記憶機構を提供するとともに、ステータスレジスタ822が次の更新までステータスを維持する。キューマネージャはDMA、受信及び送信シーケンサリクエストのために偶数番目のサイクルを予約し、プロセッサリクエストのために奇数番目のサイクルを予約する。二重ポートQRAM825は、各キューに関連する変数を格納し、各キューのための変数としては、キューのSRAM条件に対応するヘッド書き込みポインタ、ヘッド読み出しポインタ、テール書き込みポインタ、及びテール読み出しポインタと、キューのDRAM条件及びキューのサイズに対応するボディ書き込みポインタ及びボディ読み出しポインタとが挙げられる。

【0068】

アービタ815が実行される次の処理を選択した後、QRAM825の変数がフェッチされ、QALU828によって選択された処理に従って改変され、SRAM読み出しリクエスト830又はSRAM書き込みリクエスト840が生成され得る。その変数は更新され、更新されたステータスは、ステータスレジスタ822及びQRAM825に格納される。ステータスはアービタ815にも供給されて、以前にリクエストされた処理が完了したことを示す信号を供給して、リクエストの重複を防止する。ステータスレジスタ822は、4つのキューレジスタQ-Out_Ready 745、Q-In_Ready 750、Q-Empty 755、及びQ-Full 760を更新して、アクセスされたキューの新たなステータスを反映するようにする。同様に更新されるのは、SRAMアドレス833、ボディ書き込みリクエスト835及びボディ読み出しリクエスト838であり、これらはDMAを介してそのキューのためにSRAMヘッド及びテールとやりとりする。或いは、様々なプロセスが、Q書き込みデータ844によって示されるように、キューへの書き込みを行おうとすることもあり得、そのプロセスはマルチプレクサ846によって選択され、SRAM書き込みリクエスト840にパイプライン処理される。SRAMコントローラは、アクセスされたキューのテールの書き込み又はヘッドの読み出し及び確認応答を戻すことによって、読み出し及び書き込みリクエストをサービスする。このようにして、様々なキューが利用され、それらのステータスが更新される。

【0069】

図19A乃至図19Cは、INICキャッシュメモリに維持するコンテキスト又はCCBを選択するために用いられる最低使用頻度レジスタ900を示す。この実施形態におけるINICは、所定の時間においてSRAMに最大16のCCBをキャッシュでき、新たなCCBがキャッシュされるときには、古いものが破棄されなければならないことが多く、破棄されるCCBは、通常、このレジスタ900に従って最小使用頻度のCCBとして選択されたものである。この実施形態では、最大256のCCBのためのハッシュテーブルがSRAMに維持されており、一方最大256の全CCBがDRAMに保持されている。

最小使用頻度レジスタ900は、R0 - R15とラベル付けされた16個の4ビットブロックを有し、各ブロックはSRAMのキャッシュユニットに対応する。初期化の際、そのブロックは0 - 15に番号付けされ、番号0は、最小使用頻度(LRU)キャッシュユニットを表すブロックに格納され、番号15は、最大使用頻度(MRU)キャッシュユニットを表すブロックに格納される。図19Aは、LRUブロックR0が番号9を保持しており、MRUブロックR15が番号6を保持している、任意の時点におけるレジスタ900を示す。

【0070】

現在SRAMに保持されているものとは異なるCCBがキャッシュされたとき、図19Aでは、番号9を保持しているLRUブロックR0が読み出され、新たなCCBが番号9に対応するSRAMキャッシュユニットに格納される。番号9に対応する新たなCCBが現時点で最も使用頻度の高いCCBであることから、番号9は、図19Bに示すように、MRUブロックに格納される。他の番号は全て1個左のレジスタブロックにシフトされ、LRUブロックに番号1が残る。以前に番号9に対応するSRAMユニットにキャッシュされたCCBは、よりゆっくりとではあるが、よりコストパフォーマンスの良い形でDRAMに移されている。

【0071】

図19Cは、次に使用されるCCBが常にSRAMにキャッシュされたときの結果を示す。この例では、CCBが番号10に対応するSRAMユニットにキャッシュされており、そのCCBを使用した後、番号10にMPUBロックに格納される。番号10以外の最近使用された番号のみ(レジスタブロックR9 - R15)が左側にシフトされ、LRUブロックには番号1が残る。このようにして、INICはSRAMキャッシュに最もアクティブなCCBを維持する。

【0072】

場合によっては、使用されているCCBが、限定されたキャッシュメモリに保持することが望ましくないものであることがある。例えば、終了したことが判っているコンテキストのためのCCBは、キャッシュせずに、他のキャッシュされたCCBがSRAMに長く残るようにすることが好ましい。この場合、デキャッシュ可能なCCBを保持するキャッシュユニットを表す番号が、MRUブロックR15でないMRUブロックR0に格納され、従ってデキャッシュ可能なCCBがLRUブロックR0に保持された番号に対応するSRAMユニットにキャッシュされる新たなCCBが使用されるとすぐに置き換えられる。図19Dは、番号8(図19CにおけるブロックR9にあったもの)が、使用され次に終了されるCCBに対応する場合を示す。この場合では、番号8がブロックR6から取り除かれ、LRUブロックR0に格納される。次いでブロックR9の左側(R1 ~ R8)に以前に格納された全ての番号は、1個右側のブロックにシフトされる。

【0073】

図20は、最低使用頻度レジスタ900を操作するために用いられる論理ユニットをいくつか示す。16個の3又は4入力マルチプレクサ910のアレイ(明示のため、そのうちマルチプレクサMUX0、MUX7、MUX8、MUX9、及びMUX15のみが示されている)は、最低使用頻度レジスタ900の対応する16個のブロックに供給される出力を有する。例えば、MUX0の出力は、ブロックR0に格納され、MUX7の出力はブロックR7に格納される等である。レジスタブロックのそれぞれの値は、対応するマルチプレクサのための入力に接続されており、かつブロック番号のシフトの際に使用するために両隣のマルチプレクサの入力にも接続されている。例えば、R8に格納された番号は、MUX7、MUX8、及びMUX9の入力に供給される。MUX0及びMUX15のそれぞれは、ただ1つの隣接ブロックしか有しておらず、それらのマルチプレクサのための別の入力がLRU及びMRUブロックの選択のためにそれぞれ用いられる。MUX15は、4入力マルチプレクサとして図示されており、入力915はR0に格納された番号を供給する。

【0074】

10

20

30

40

50

16個の比較器920のアレイのそれぞれが、最低使用頻度レジスタ900の対応するブロックに格納された値を受け取る。各比較器は、ライン935を介してプロセッサ470からの信号を受け取り、従ってプロセッサ470によって送られた一致する番号を有するレジスタブロックが、論理回路930に真の値を出力し、一方他の15個の比較器は偽の値を出力する。論理回路930は、マルチプレクサへの入力を選択するためにマルチプレクサのそれぞれに接続された一対の選択ラインを制御し、従ってレジスタブロック番号のシフトを制御する。従って選択ライン939はMUX0を制御し、選択ライン944はMUX7を制御し、選択ライン949はMUX8を制御し、選択ライン954はMUX9を制御し、選択ライン959はMUX15を制御する。

【0075】

CCBが使用されるとき、プロセッサ470は、そのCCBが16個のキャッシュユニットの1つを現在保持されているCCBに一致するか否かを調べるためのチェックを行う。一致が見出された場合は、そのプロセッサは、ライン935を通してそのキャッシュユニットに対応するブロック番号、例えば番号12を有する信号を送る。比較器920は、そのライン935からの信号とブロック番号とを比較し、比較器C8はその信号に一致するブロックR8に真の出力を供給し、他の全ての比較器は偽の値を出力する。論理回路930はプロセッサ470からの制御の下で、選択ライン959を用いてMUX15のためライン935からの入ロクを選択し、MRUブロックR15に番号12を格納する。論理回路930も、MUX15を除いて、MUX8及びより高位のマルチプレクサのための対になった選択ラインを通して信号を送り、各マルチプレクサMUX8及びそれより上位のマルチプレクサへの入力として、1個右側のレジスタブロック(R9 - R15)に格納された値を選択することにより、それらの出力を1ブロック左にシフトさせる。MUX8の左側のマルチプレクサの出力は一定に選択される。

【0076】

一方プロセッサ470が16個のキャッシュユニットのなかでCCBについて一致を見出せなかった場合は、プロセッサは、ライン966を通してLRUブロックR6から読み出しを行い、LRUブロックに対応するキャッシュを特定し、そのキャッシュに格納されたデータをDRAMに書き込む。この場合では番号3であるR0に格納された番号は、MRUブロックR5における記憶機構のためにMUX15への入力915としてライン959を選択することによって選出される。他の15個のマルチプレクサは、それらの各レジスタブロックにそれぞれの右隣のレジスタブロックに格納されていた番号を出力する。

【0077】

プロセッサが、使用後にキャッシュからCCBを取り除こうとしている状況では、MRUブロックR5ではなく、LRUブロックR0がそのCCBを維持するキャッシュユニットに対応する番号の置き換えのために選択される。SRAMから除去するためにLRUブロックR0に置かれるCCBに対応する番号(例えばブロックR9に保持されている番号1)は、ライン935を通してプロセッサ470に送られ、それは比較器C9によって一致を取られる。プロセッサは、論理回路930にMUX0への入力935としてライン939を選択することにより、番号1をR0に入力するように命令する。MUX9への選択ライン954は、入力としてレジスタブロックR8に保持された番号を選択し、従ってR8からの番号はR9に格納される。R0~R9の間の他のレジスタブロックによって保持された番号は、同様に右側にシフトされ、R9の右側のレジスタブロックにおける番号は一定に維持される。これによって、多くのサイクルにわたって使用されなくなったCCBが限りのあるキャッシュメモリに維持されることがなくなるとともに、それらの特定する番号がMRUからLRUブロックにレジスタブロックを通して動くことになる。

【0078】

既に述べたように、上述のデータ通信の処理のための装置及びシステムは、大きなコネクションベースのメッセージを処理するために必要な時間を劇的に短縮させる。プロトコル処理速度は、専用に設計されたプロトコル処理ハードウェアによって汎用CPUが従来のプロトコルソフトウェアを実行する場合と比較して非常に高速化され、かつホストのCP

10

20

30

40

50

Uへの割り込みも実質的に低減する。これらの利点は、インテリジェントネットワークインタフェースカード（INIC）を追加することによって既存のホストにもたすこともできるが、或いはプロトコル処理ハードウェアをCPUに組み込んでもよい。何れの場合も、プロトコル処理ハードウェア及びCPUは、インテリジェントに与えられたメッセージを処理するデバイスを決定し、メッセージの条件に基いてその処理の割り当てを変更することができる。

【図面の簡単な説明】

【図1】 ネットワーク通信を高速化するための通信処理装置を有するホストコンピュータを備えた、本発明の装置の平面図。

【図2】 高速パス、低速パス、及び高速パスと低速パスとの間のコネクションコンテキストの転送を含む、ネットワーク通信の処理における図1のホストの情報の流れを示す図。

【図3】 本発明によるメッセージ受信の流れ図。

【図4】 A図は、低速パスによって処理される初めのメッセージパケットを受信した、図1のホストにおける情報の流れを示す図、B図は、高速パスによって処理される初めのメッセージパケットを受信した、図1のホストにおける情報の流れを示す図、C図は、高速パスによって処理される順次式のメッセージパケットを受信した、図4Bのホストにおける情報の流れを示す図、D図は、低速パスに戻すための処理を生じさせる誤りを有するメッセージパケットを受信した、図4Cのホストにおける情報の流れを示す図。

【図5】 高速パス又は低速パスの何れかによってメッセージを転送する、図1のホストにおける情報の流れを示す図。

【図6】 TCP/IP処理スタックを有するクライアントに関連するインテリジェントネットワークインタフェースカード（INIC）の第1の実施形態における情報の流れを示す図。

【図7】 パケット制御シーケンサ、フライバイシーケンサを含む、図6に示すINICの実施形態のためのハードウェアロジックを示す図。

【図8】 INICによって受信されたときにヘッダバイトを解析する、図7のフライバイシーケンサを示す図。

【図9】 TCP/IP処理スタックを有するサーバに関連するINICの第2の実施形態における情報の流れを示す図。

【図10】 高速パスのための通信制御ブロックを生成し制御するための、図9にインストールされたコマンドドライバを示す図。

【図11】 NetBios通信のために構成された図10のコマンドドライバ及びTCP/IPスタックを示す図。

【図12】 図6のクライアントと図9のサーバとの間の通信の交換を示す図。

【図13】 図9のINICに含められたハードウェア機構を示す図。

【図14】 各フェーズに1個のプロセッサがある3つのフェーズを含む、図13のINICに含められた3個1組のパイプライン型マイクロプロセッサを示す図。

【図15】 A図は図14のパイプライン型マイクロプロセッサの第1フェーズを示す図、B図は図14のパイプライン型マイクロプロセッサの第2フェーズを示す図、C図は図14のパイプライン型マイクロプロセッサの第3フェーズを示す図。

【図16】 図14のマイクロプロセッサと相互作用し、かつSRAM及びDRAMを備えた複数のキュー記憶ユニットを示す図。

【図17】 図16のキュー記憶ユニットのための、1組のステータスレジスタを示す図。

【図18】 図16及び図17のキュー記憶ユニット及びステータスレジスタと相互作用する、キューマネージャを示す図。

【図19】 A図はキャッシュメモリの割り当てのために用いられる最低使用頻度レジスタの様々な動作段階の1つを示す図、B図はキャッシュメモリの割り当てのために用いられる最低使用頻度レジスタの様々な動作段階の1つを示す図、C図はキャッシュメモリの

10

20

30

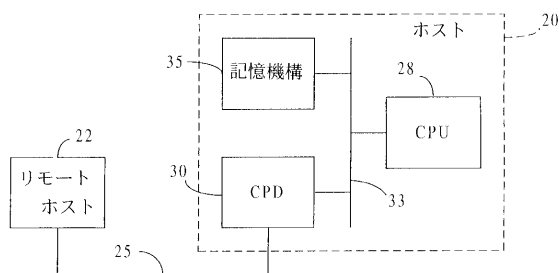
40

50

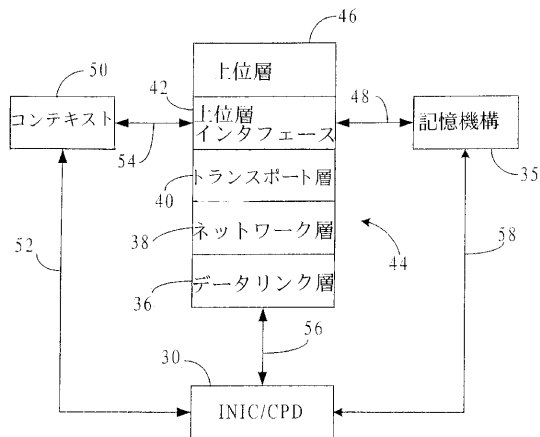
割り当てのために用いられる最低使用頻度レジスタの様々な動作段階の1つを示す図、D図はキャッシュメモリの割り当てのために用いられる最低使用頻度レジスタの様々な動作段階の1つを示す図。

【図20】 図19A乃至図19Dの最低使用頻度レジスタを操作するための用いられるデバイスを示す図。

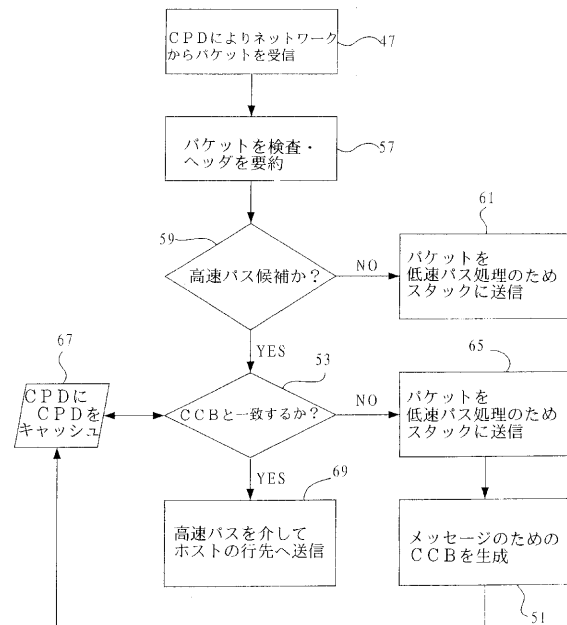
【図1】



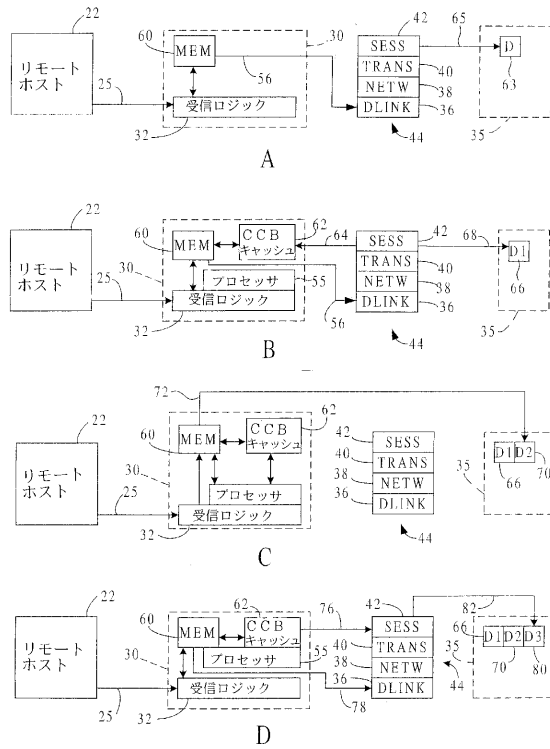
【図2】



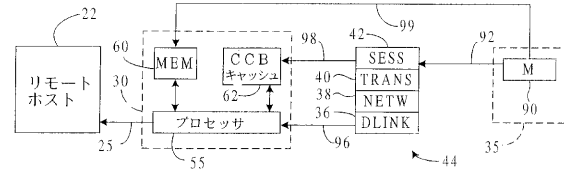
【図3】



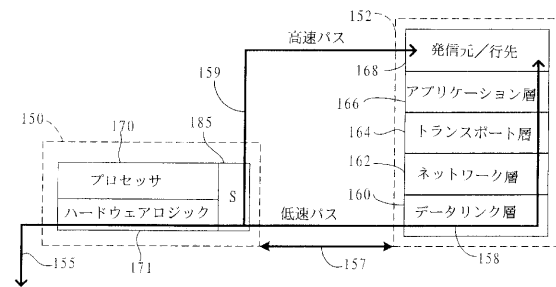
【図 4】



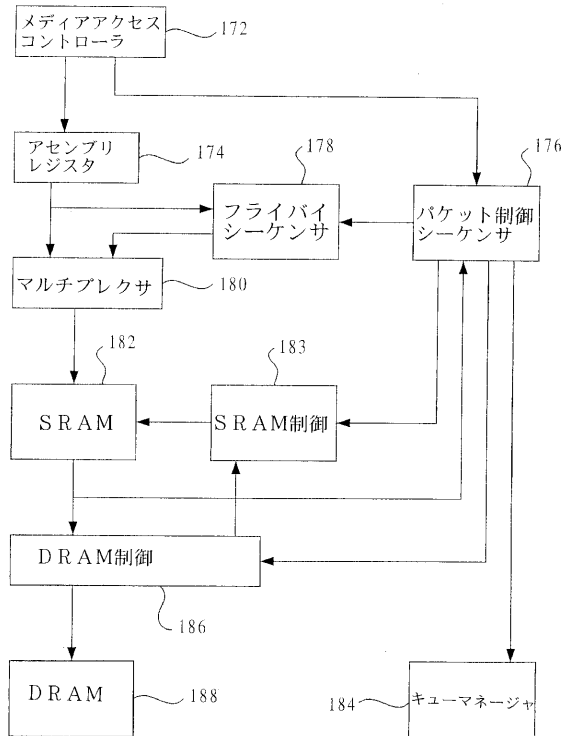
【図 5】



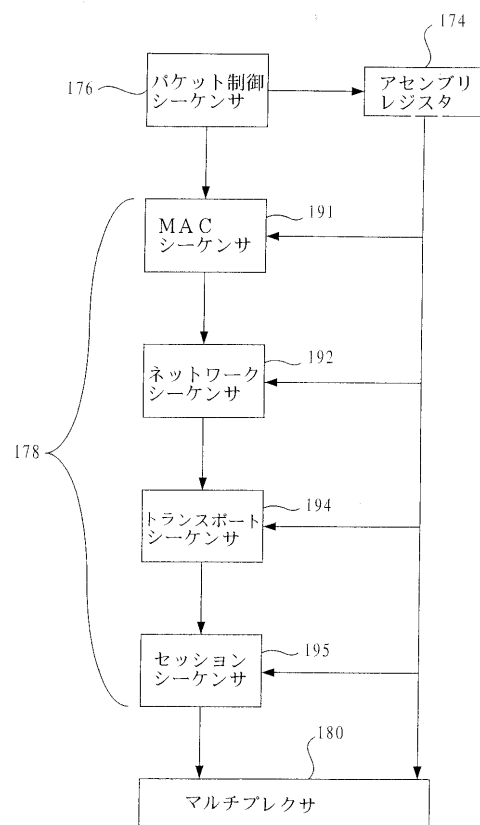
【図 6】



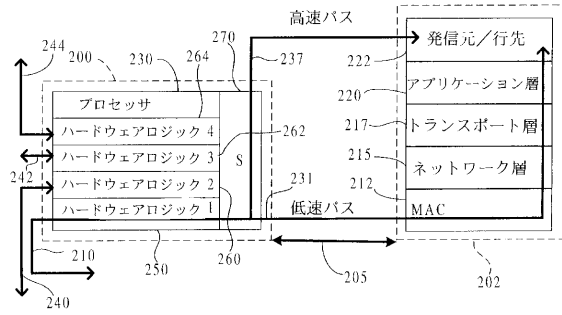
【図 7】



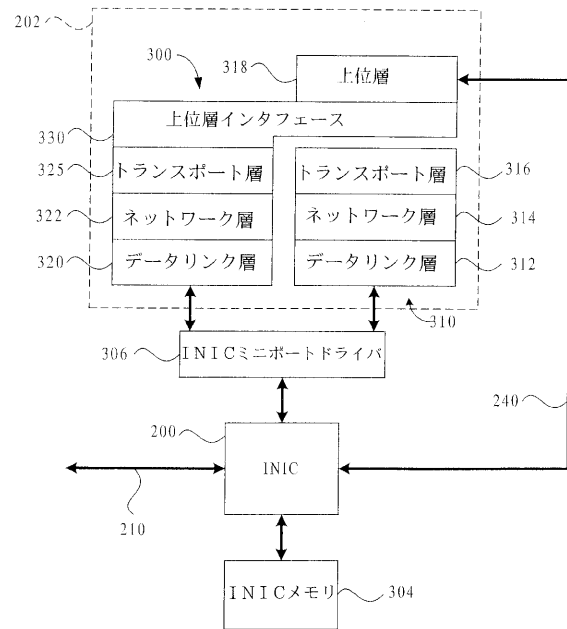
【図 8】



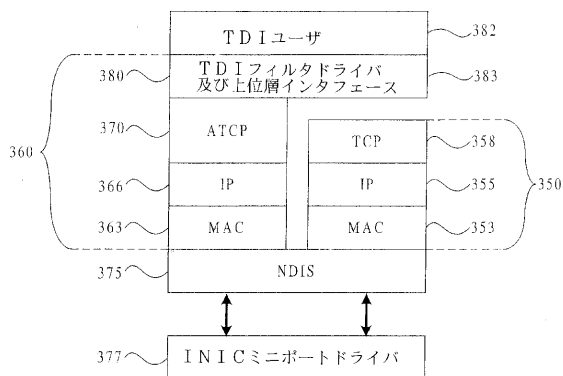
【図 9】



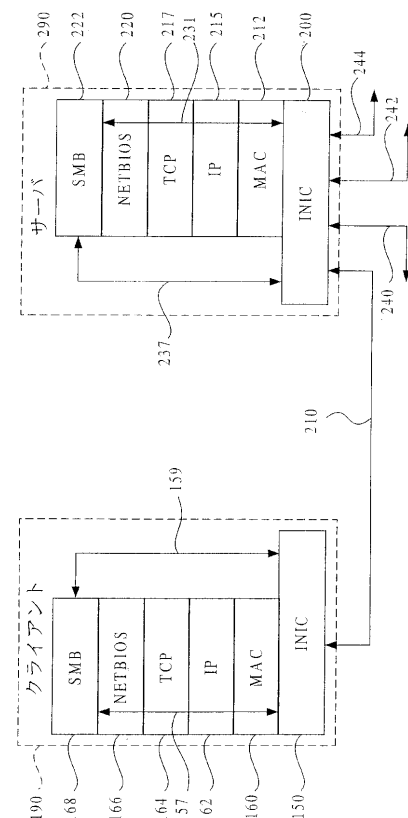
【図 10】



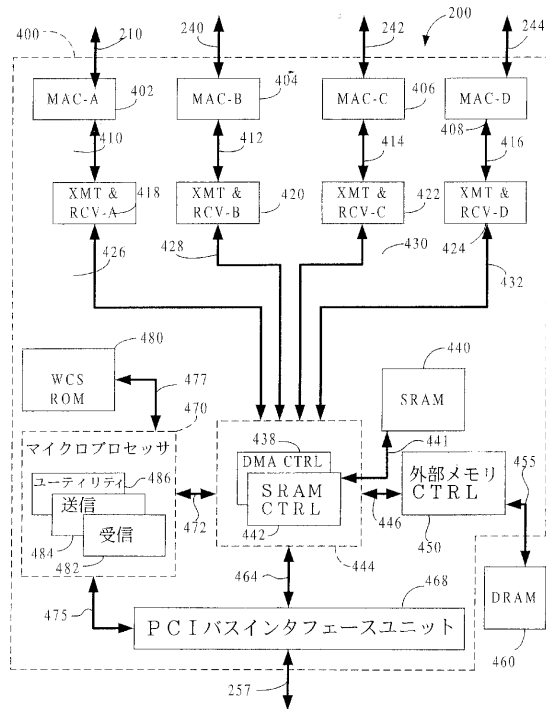
【図 11】



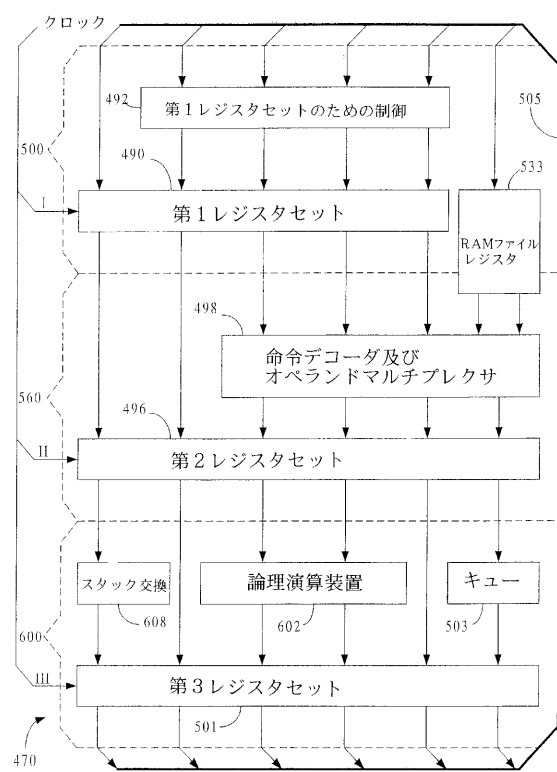
【図 12】



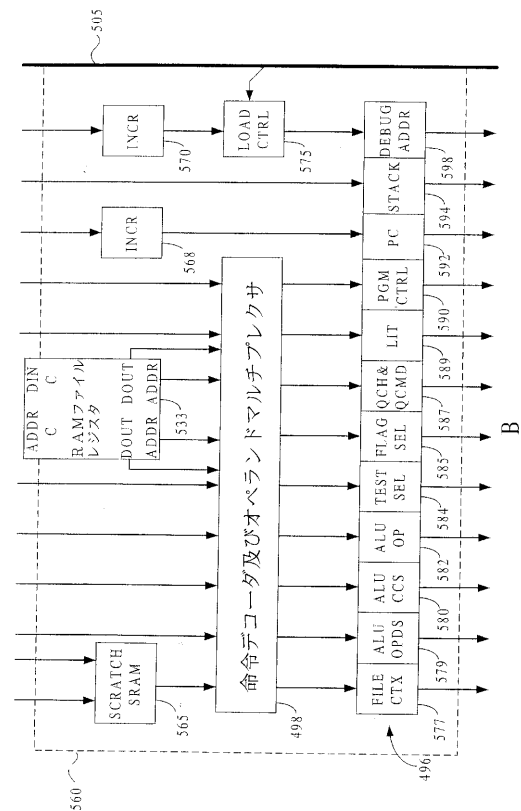
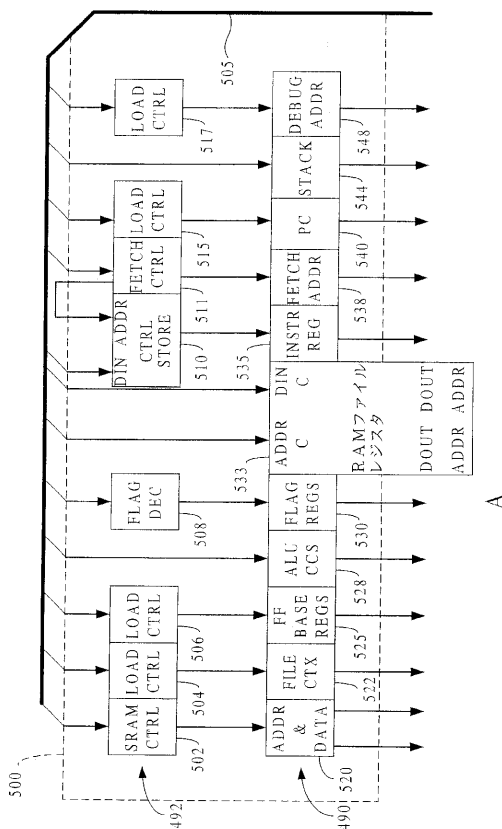
【図13】

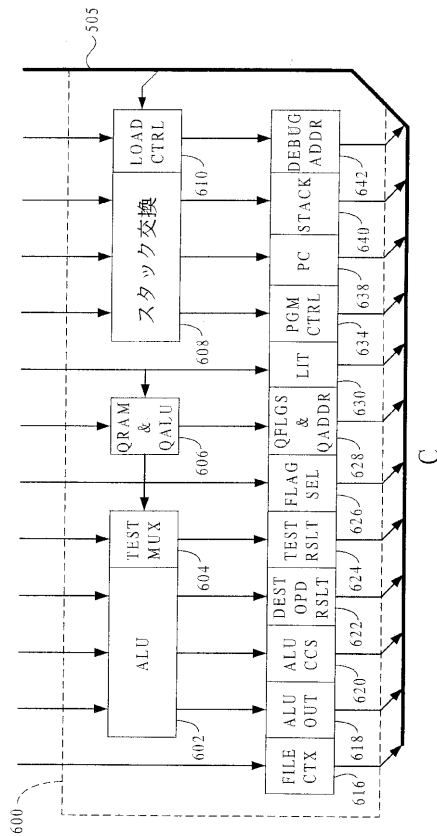


【図14】

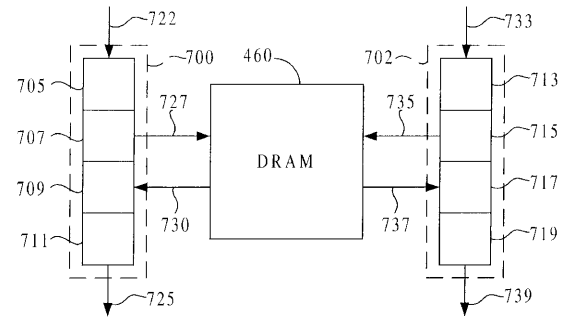


【図15】

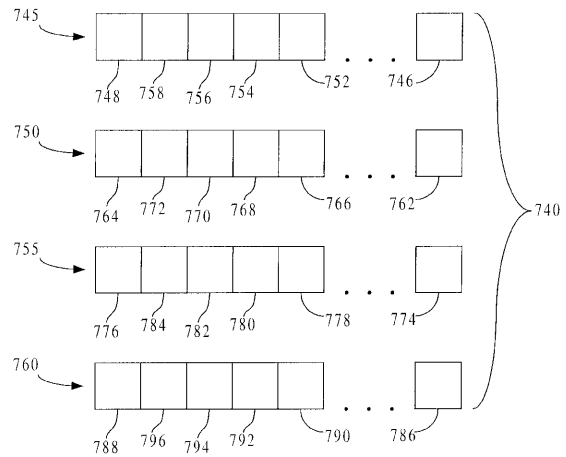




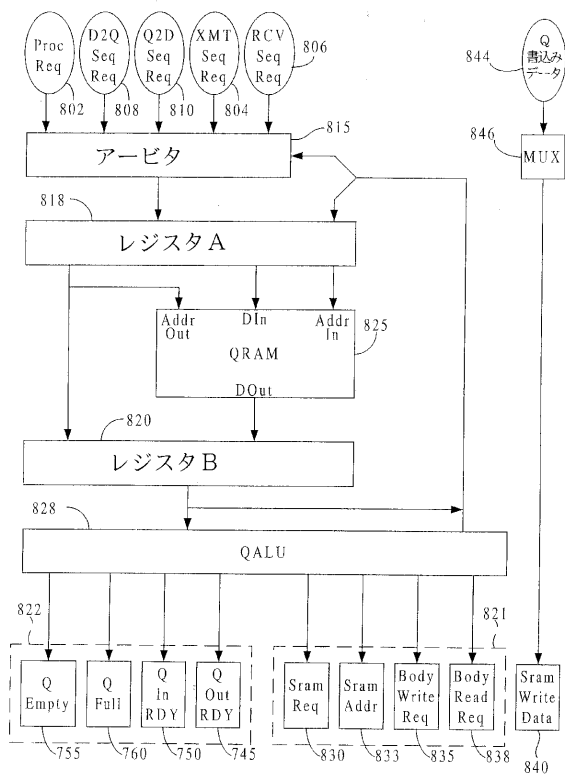
【図 16】



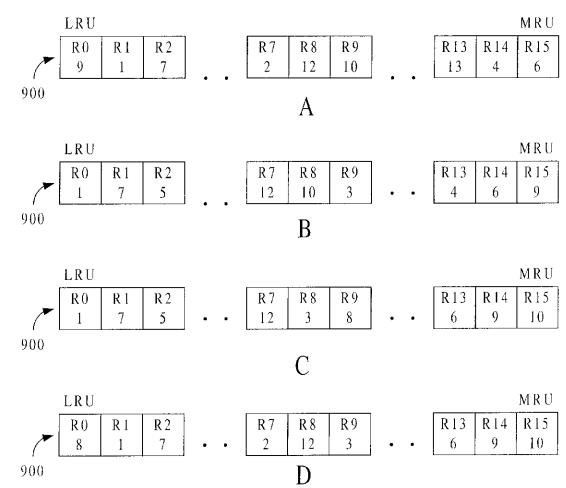
【図 17】



【図 18】



【図 19】



フロントページの続き

- (72)発明者 クラフト, ピーター, ケイ
アメリカ合衆国, カリフォルニア州・ 9 4 1 1 4 , サンフランシスコ, ヘンリー・ストリート・ 1
5 6
- (72)発明者 ヒギン, デヴィッド, エー
アメリカ合衆国, カリフォルニア州・ 9 5 0 7 0 , サラトガ, ロス・アラモス・ドライブ・ 1 7 8
8 0
- (72)発明者 フィルブリック, クライヴ, エム
アメリカ合衆国, カリフォルニア州・ 9 5 1 2 5 , サンノゼ, ロイコット・ウェイ・ 1 1 7 0
- (72)発明者 スター, ダリル
アメリカ合衆国, カリフォルニア州・ 9 5 0 3 5 , ミルピタス, フォルソム・コート・ 4 4 6

審査官 小曳 満昭

- (56)参考文献 特開平 0 2 - 2 3 8 5 4 4 (J P , A)
特開平 0 9 - 1 2 8 3 1 4 (J P , A)
特開平 0 3 - 2 5 0 9 4 6 (J P , A)
特開平 0 2 - 1 3 7 5 5 5 (J P , A)
特開昭 6 2 - 0 3 8 0 6 0 (J P , A)
特開平 0 1 - 2 5 6 2 4 8 (J P , A)
向井伸一郎, L A Nカード周辺ハードウェア動作概要と機能評価, インターフェース, C Q出版
株式会社, 1 9 9 2年 4月 1日, 第18巻, 第4号, p.223~237
陣崎 明 外 3 名, 大規模広域並列分散システムの実現を目指す超高速インターネットの構想,
電子情報通信学会技術研究報告 C P S Y 9 7 - 6 1 , 社団法人電子情報通信学会, 1 9 9 7年
8月 2 0日, 第97巻, 第226号, p.75~82

(58)調査した分野(Int.Cl. , D B 名)

G06F 13/00、
H04L 12/00-12/26、12/50-13/18、
29/00-29/12