

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5066702号  
(P5066702)

(45) 発行日 平成24年11月7日 (2012. 11. 7)

(24) 登録日 平成24年8月24日 (2012. 8. 24)

(51) Int. Cl.

F I

G 0 6 F 13/38 (2006. 01)

G 0 6 F 13/38 3 5 0

G 0 6 F 12/00 (2006. 01)

G 0 6 F 12/00 5 4 5 M

G 0 6 F 13/10 (2006. 01)

G 0 6 F 13/10 3 4 0 B

請求項の数 35 (全 50 頁)

(21) 出願番号 特願2002-531029 (P2002-531029)  
 (86) (22) 出願日 平成13年9月24日 (2001. 9. 24)  
 (65) 公表番号 特表2004-510252 (P2004-510252A)  
 (43) 公表日 平成16年4月2日 (2004. 4. 2)  
 (86) 国際出願番号 PCT/US2001/030150  
 (87) 国際公開番号 W02002/027519  
 (87) 国際公開日 平成14年4月4日 (2002. 4. 4)  
 審査請求日 平成20年9月24日 (2008. 9. 24)  
 (31) 優先権主張番号 09/675, 700  
 (32) 優先日 平成12年9月29日 (2000. 9. 29)  
 (33) 優先権主張国 米国 (US)  
 (31) 優先権主張番号 09/675, 484  
 (32) 優先日 平成12年9月29日 (2000. 9. 29)  
 (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・9 5  
 0 1 2, サン・ノゼ, イースト・ギッシュ  
 ・ロード・2 3 4  
 (74) 代理人 100098062  
 弁理士 梅田 明彦  
 (72) 発明者 スター, ダリル, ディー  
 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  
 (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

最終頁に続く

(54) 【発明の名称】 インテリジェントネットワークストレージインタフェースシステム及びデバイス

(57) 【特許請求の範囲】

【請求項 1】

ネットワークと記憶装置間で情報を転送するための装置であって、  
 ファイルシステムを動作させる C P U と、ホストメモリバスにより前記 C P U に接続され  
 たホストメモリとを有するホストコンピュータと、  
 前記ホストコンピュータと前記ネットワークと前記記憶装置とに接続され、前記ネット  
 ワークと記憶装置間で通信されるデータを格納するのに適したインタフェースファイルキ  
 ャッシュを含むインタフェースメモリを有するインタフェース装置とを備え、  
 前記インタフェースファイルキャッシュが前記ファイルシステムにより制御され、  
 前記ホストコンピュータが、通信制御ブロックを作成しかつ前記通信制御ブロックを前  
 記インタフェース装置へ送るように構成され、  
 前記インタフェース装置が、前記通信制御ブロックに従って前記ネットワークと前記フ  
 ァイルキャッシュ間で前記データを通信するように構成されていることを特徴とする装置  
 。

【請求項 2】

前記ホストコンピュータが、通信制御ブロックを作成しかつ前記通信制御ブロックを前  
 記インタフェース装置へ送るように構成され、かつ  
 前記インタフェース装置が、前記通信制御ブロックに従って前記記憶装置を前記ファイ  
 ルキャッシュ間で前記データを通信するように構成されていることを特徴とする請求項 1  
 に記載の装置。

10

20

**【請求項 3】**

前記データが前記ホストコンピュータに入力しないことを特徴とする請求項 1 に記載の装置。

**【請求項 4】**

前記データが、前記ネットワークと前記インタフェース装置間をファイル形式で通信され、かつ

前記データが前記インタフェース装置と前記記憶装置間をブロック形式で通信されることを特徴とする請求項 1 に記載の装置。

**【請求項 5】**

前記データがトランスポートレイヤヘッダを含むヘッダに添付され、かつ前記インタフェース装置が前記トランスポートレイヤヘッダを処理するように構成されたメカニズムを有することを特徴とする請求項 1 に記載の装置。

10

**【請求項 6】**

前記インタフェース装置が、前記ホストを有する第 1 パスを介して、または前記ホストを有しない第 2 パスを介して前記ネットワークと前記記憶装置間で前記データを転送するかどうかを選択するように構成されたメカニズムを有することを特徴とする請求項 1 に記載の装置。

**【請求項 7】**

前記インタフェース装置が、前記記憶装置に接続された S C S I コントローラを有することを特徴とする請求項 1 に記載の装置。

20

**【請求項 8】**

前記インタフェース装置が、前記記憶装置に接続されたファイバチャネルコントローラを有することを特徴とする請求項 1 に記載の装置。

**【請求項 9】**

前記インタフェース装置が、前記記憶装置に接続された R A I D コントローラを有することを特徴とする請求項 1 に記載の装置。

**【請求項 10】**

前記インタフェース装置が、前記ネットワーク及び前記記憶装置の少なくとも一方に接続されたネットワークポートを有することを特徴とする請求項 1 に記載の装置。

**【請求項 11】**

30

前記ホストコンピュータ、前記ネットワーク及び第 2 記憶装置に接続された第 2 インタフェース装置を更に備え、

前記第 2 インタフェース装置が、前記ネットワークと前記第 2 記憶装置間で通信されるデータを格納するのに適した第 2 インタフェースファイルキャッシュを含む第 2 インタフェースメモリを有し、

前記第 2 インタフェースファイルキャッシュが前記ファイルシステムにより制御されることを特徴とする請求項 1 に記載の装置。

**【請求項 12】**

ネットワークと記憶装置間で情報を転送するための装置であって、

ホストメモリバスによりホストメモリに接続されたプロセッサを有し、前記ホストメモリが、ネットワーク接続を定義する通信制御ブロックを作成するように前記プロセッサにより動作可能なプロトコルスタックを有するホストコンピュータと、

40

前記ホストコンピュータに接続されかつ前記ネットワークと前記記憶装置間に結合されたインタフェース装置とを備え、前記インタフェース装置が、前記通信制御ブロックを格納するのに適したインタフェースメモリと、データを前記通信制御ブロックと関連付け、それにより前記データが前記ホストコンピュータに出会うことなく前記ネットワークと前記記憶装置間で通信されるように構成されたメカニズムとを有することを特徴とする装置。

**【請求項 13】**

前記ホストコンピュータがファイルシステムを有し、かつ前記インタフェースメモリが

50

前記データを格納するのに適したファイルキャッシュを有し、前記ファイルシステムが前記データの前記ファイルキャッシュへの格納を管理することを特徴とする請求項 1 2 に記載の装置。

【請求項 1 4】

前記データが、トランスポートレイヤヘッダを含む少なくとも 1 つのパケットで前記ネットワーク上を移動し、かつ前記インタフェース装置が前記ヘッダ処理するための回路を有することを特徴とする請求項 1 2 に記載の装置。

【請求項 1 5】

前記トランスポートレイヤヘッダが T C P ヘッダであることを特徴とする請求項 1 4 に記載の装置。

【請求項 1 6】

前記トランスポートレイヤヘッダが U D P ヘッダであることを特徴とする請求項 1 4 に記載の装置。

【請求項 1 7】

前記インタフェース装置が入力 / 出力バスにより前記ホストコンピュータに接続され、かつ前記通信コントロールブロックが前記入力 / 出力バス上で前記インタフェース装置と前記ホストコンピュータ間を移動することを特徴とする請求項 1 2 に記載の装置。

【請求項 1 8】

前記インタフェース装置が更に、前記記憶装置に接続された S C S I コントローラを有することを特徴とする請求項 1 2 に記載の装置。

【請求項 1 9】

前記インタフェース装置が更に、前記記憶装置に接続されたファイバチャネルコントローラを有することを特徴とする請求項 1 2 に記載の装置。

【請求項 2 0】

前記インタフェース装置が更に、前記記憶装置に接続された R A I D コントローラを有することを特徴とする請求項 1 2 に記載の装置。

【請求項 2 1】

前記インタフェース装置が前記記憶装置と、少なくとも 1 つのネットワークポートにより前記ネットワークに接続されていることを特徴とする請求項 1 2 に記載の装置。

【請求項 2 2】

前記通信制御ブロックに対応する先のまたは後のデータが前記ホストコンピュータを通過することを特徴とする請求項 1 2 に記載の装置。

【請求項 2 3】

前記ホストコンピュータが、前記インタフェース装置によってアクセス可能な U D P ソケットを指定するように構成され、かつ前記インタフェース装置が、前記 U D P ソケットに従って前記ネットワークと前記ファイルキャッシュ間で前記データを通信するように構成されていることを特徴とする請求項 1 2 に記載の装置。

【請求項 2 4】

ネットワークと記憶装置間で情報を転送するための装置であって、

ファイルシステムを動作させる C P U と、ホストバスにより前記 C P U に接続されたホストメモリとを有し、前記ファイルシステムにより情報がファイルの階層構造として論理的に編成されるホストコンピュータと、

前記ホストコンピュータ、前記ネットワーク及び前記記憶装置に接続され、前記ファイルシステムの制御下で前記ネットワークと前記記憶装置間で通信されるデータを格納するのに適したインタフェースファイルキャッシュを含むインタフェースメモリを有し、前記ファイルシステムによって前記ファイルキャッシュ内の情報が前記ファイルの階層構造の一部として論理的に編成されるようにしたインタフェース装置とを備え、

前記ホストコンピュータが、前記インタフェース装置によりアクセス可能な U D P ソケットを指定するように構成され、かつ、

前記インタフェース装置が、前記 U D P ソケットに従って前記ネットワークと前記ファ

10

20

30

40

50

イルキャッシュ間で前記データを通信するように構成されていることを特徴とする装置。

【請求項 2 5】

前記ホストコンピュータが、前記インタフェース装置によりアクセス可能なアプリケーションレイヤヘッダを作成するように構成され、かつ前記インタフェース装置が前記アプリケーションレイヤヘッダを前記データにプリペンドするように構成されていることを特徴とする請求項 2 4 に記載の装置。

【請求項 2 6】

前記ホストコンピュータが、前記インタフェース装置によりアクセス可能なリアルタイムトランスポートプロトコルヘッダを作成するように構成され、かつ前記インタフェース装置が前記リアルタイムトランスポートプロトコルヘッダを前記データにプリペンドするように構成されていることを特徴とする請求項 2 4 に記載の装置。

【請求項 2 7】

前記データが関連する UDP ヘッダと共に格納され、かつ前記インタフェース装置が前記 UDP ヘッダを処理するように構成されたメカニズムを有することを特徴とする請求項 2 4 に記載の装置。

【請求項 2 8】

前記データが前記インタフェース装置により UDP ヘッダでプリペンドされて UDP データグラムを作成し、かつ前記インタフェース装置が前記データグラムを複数のフラグメントに分割するように構成されたメカニズムを有することを特徴とする請求項 2 4 に記載の装置。

【請求項 2 9】

前記データが複数のフラグメントに配置され、かつ前記インタフェース装置が、UDP ヘッダに対応する前記フラグメントを連結するように構成されたメカニズムを有することを特徴とする請求項 2 4 に記載の装置。

【請求項 3 0】

前記データが前記ホストコンピュータに入力しないことを特徴とする請求項 2 4 に記載の装置。

【請求項 3 1】

前記データが音声データを含むことを特徴とする請求項 2 4 に記載の装置。

【請求項 3 2】

前記データが映像データを含むことを特徴とする請求項 2 4 に記載の装置。

【請求項 3 3】

前記データがリアルタイム通信の一部であることを特徴とする請求項 2 4 に記載の装置。

【請求項 3 4】

ネットワークと記憶装置間で情報を転送するための装置であって、

ファイルシステムを動作させる CPU と、メモリバスにより前記 CPU に接続されたメモリとを有するコンピュータと、

前記コンピュータ、前記ネットワーク及び前記記憶装置に接続され、前記ネットワークと記憶装置間で通信されるデータを格納するインタフェースファイルキャッシュを含むインタフェースメモリを有するインタフェース装置とを備え、

前記インタフェースファイルキャッシュが前記ファイルシステムにより制御され、

前記データが前記インタフェース装置と前記記憶装置間での前記データの転送の際に iSCSI ヘッダに添付され、

前記コンピュータが、通信制御ブロックを作成し、かつ前記通信制御ブロックを前記インタフェース装置へ送るように構成され、

前記インタフェース装置が、前記通信制御ブロックに従って前記ネットワークと前記ファイルキャッシュ間で前記データを通信するように構成されていることを特徴とする装置。

【請求項 3 5】

前記インタフェース装置がギガビットイーサネットネットワークにより前記記憶装置に接続されていることを特徴とする請求項 3 4 に記載の装置。

【発明の詳細な説明】

【0001】

【技術分野】

本発明は、ネットワーク通信及びストレージに関する。

【0002】

【背景技術】

過去数 10 年以上に亘って、ネットワークコンピューティングの利益及び進歩がコンピュータネットワークの著しい成長を促しており、これがより一層の進歩、成長及び利益に拍車をかけている。しかしながら、この成長と共に従来のネットワークデバイスを用いる上で混乱及び障害が起こっている。例えば、ネットワークに接続されたコンピュータの CPU が、そのネットワーク通信を処理するのに使用する時間の部分が増加して、他の作業に使うことができる時間が少なくなっている場合がある。特に、ネットワークとディスクドライブのようなコンピュータの記憶装置との間でファイルデータを移動させることへの要求が加速度的に増えている。従来、このようなデータはネットワーク上での移動のためにパケットに分割され、各パケットは受信コンピュータの CPU によって一度に 1 つのレイヤが処理される制御情報のレイヤにカプセル化される。CPU の速度が一定の割合で増加しているにも拘わらず、このファイル転送のようなネットワークメッセージのプロトコル処理は、市販の最速 CPU の利用可能な処理能力の大部分を消費する場合がある。

【0003】

このような状況は、その主な機能がネットワーク上でファイルデータを転送することによってその接続型ディスクまたはテープドライブ上でファイルを格納しかつ検索することであるネットワークファイルサーバにとって、より一層挑戦的なものになる場合がある。ネットワーク及びデータベースが成長するにつれて、そのようなサーバに格納される情報の量が爆発的に増加し、そのようなサーバ接続型記憶装置の限界を露呈している。ホスト CPU により処理されるプロトコルの上述した問題に加えて、従来の小型コンピュータシステムインタフェース (SCSI) インタフェースのようなパラレルデータチャネルの限界が、ストレージの必要性が増加するにつれて明らかになっている。例えば、パラレル SCSI インタフェースは、サーバに接続することができるストレージデバイスの数及びストレージデバイスとサーバ間の距離を制限する。

【0004】

Tom Clark 著の書籍「Designing Storage Area Networks」(著作権 1999 年)に記載されるように、サーバ接続型のパラレル SCSI ストレージデバイスの限界に対する 1 つの解決策は、ネットワークサーバの正面にある既存のローカルエリアネットワーク (LAN) に別のファイルサーバを接続することを必要とする。このネットワーク接続型ストレージ (NAS) によって、ネットワーク上の他のサーバ及びクライアントから NAS ファイルサーバへのアクセスが可能になるが、元のネットワークサーバに専用の記憶容量が増加するものではない。逆に、NAS は、元のネットワークサーバが様々な NAS ファイルサーバと通信する必要があることから、そのサーバが必要とするプロトコル処理が増える場合がある。更に、NAS ファイルサーバのそれぞれが、次に、プロトコル処理の負担及びストレージインタフェースの限界を被ることになる。

【0005】

ストレージエリアネットワーキング (SAN) が、サーバの背後に接続されたストレージデバイスのネットワークでデジチェーンにより接続された SCSI ストレージデバイスを置き換えることによって、ネットワーク上でのファイル転送及び記憶に対する必要性の増加に対する別の解決策を提供している。イーサネットまたはファーストイーサネットのような従来のネットワーク規格の代わりに、SAN は、ファイバチャネル (FC) と呼ばれる新しいネットワーキング規格を使用する。しかしながら、それが比較的最近の導入であるため、多くの市販されている FC デバイスは互いに互換性がない。また、FC ネット

10

20

30

40

50

トワークはサーバと記憶装置のようなネットワーク上の２点間の通信のためにバンド幅を専用にすることがあり、前記２点が通信していない時はバンド幅が無駄になる。

【０００６】

今や公知のＮＡＳ及びＳＡＮは、転送されかつ格納されるデータの形によって区別することができる。ＮＡＳデバイスは一般にデータファイルを別のファイルサーバまたはクライアントへ及びそこから転送するのに対し、ＳＡＮ上ではデバイスレベルのデータのブロックが転送される。このため、ＮＡＳデバイスは通常記憶のためにファイルとブロック間で変換するためのファイルシステムを有するが、ＳＡＮはそのようなファイルシステムを持たないストレージデバイスを備えることができる。

【０００７】

これに代えて、ＮＡＳファイルサーバは、イーサネットＳＡＮの一部として、サーバに専用のイーサネットベースのネットワークに接続することができる。Marc Farleyは、その書籍「Building Storage Networks」（著作権２０００年）において、イーサネット上でストレージプロトコルを実行することが可能であり、それによってファイバチャネルの非互換性の問題を解消できると述べている。しかしながら、ＳＡＮのようなネットワークトポロジを使用することによってサーバに接続されたストレージデバイスの数の増加によって、そのサーバが実行しなければならないプロトコル処理の量が増加する。上述したように、このようなプロトコルの処理は既に最新型のサーバを乱用している。

【０００８】

ファイル転送のようなネットワークメッセージの従来の処理の一例は、ネットワークデータの格納を遅くするいくつかの処理ステップを説明している。ネットワークインタフェースカード（ＮＩＣ）は一般に、ホストがネットワークにアクセスできるようにする媒体アクセスコントロール（ＭＡＣ）機能を提供することに加えて、ホストとネットワーク間の物理的接続を提供する。ホストに送られたネットワークメッセージのパケットがＮＩＣに到着すると、そのパケットのＭＡＣレイヤヘッダは処理されかつ前記パケットはＮＩＣにおいて巡回冗長検査（ＣＲＣ）を受ける。次に前記パケットは、周辺装置相互接続（ＰＣＩ）のような入力／出力（Ｉ／Ｏ）バスを通してホストに送られかつホストメモリに格納される。次に、ＣＰＵが、プロトコルスタックからの命令を実行することによって、前記パケットのヘッダレイヤのそれぞれを順に処理する。これには、最初に前記パケットを格納するためにホストメモリバスを通る移動と、次に各ヘッダレイヤを順に処理するためにホストメモリバスを通る移動とが必要である。そのパケットの全ヘッダレイヤが処理された後、該パケットからのペイロードデータが、前記メッセージの他の同様に処理されたペイロードパケットと共にファイルキャッシュ内にグループ分けされる。前記データは、ディスク上に格納するためのファイルブロックとして前記ファイルシステムに従ってＣＰＵにより再アセンブリされる。全パケットが処理されかつメッセージがファイルキャッシュにファイルブロックとして再アセンブリされた後、前記ファイルは、それぞれに数個のペイロードパケットから構築することができるデータブロックにして、ホストメモリバス及びＩ／Ｏバス上をホストストレージヘディスク上で長期間格納するために、一般にＩ／ＯバスにブリッジされたＳＣＳＩバスを介して送られる。

【０００９】

これに代えて、ＳＡＮ上でファイルを格納するために、ファイルキャッシュの再アセンブリされたファイルがブロックでホストメモリバス及びＩ／Ｏバス上を前記ＳＡＮ用に構成されたＩ／Ｏコントローラへ送られる。ＳＡＮがＦＣネットワークである場合には、ファイバチャネルプロトコル（ＦＣＰ）に従ってファイルブロックをＳＡＮ上のストレージデバイスに送ることができる特別のＦＣコントローラが設けられる。ファイルがＮＡＳデバイス上で格納されることになっている場合には、該ファイルはＮＡＳデバイスへ送りまたは再送することができ、これがパケットを上述したとほぼ同様に処理するが、ＮＡＳデバイスのＣＰＵ、プロトコルスタック及びファイルシステムを使用し、ＮＡＳデバイスの記憶装置に前記ファイルのブロックを格納する。

【００１０】

このように、ホストに接続されたN A SまたはS A Nでの格納のためにネットワークから該ホストに送られるファイルは、一般に前記ファイルの各メッセージパケットについてI / Oバスを通る2つの移動が必要である。更に、各パケットのヘッダレイヤの制御情報は、一時的に格納され、一度に1つのレイヤが処理されかつ次にI / Oバスへ送り戻されるので、ホストメモリバスを繰り返し通過することになる。このようなファイルをクライアントからのリクエストに回答してS A N上のストレージから検索することはまた、従来ホストC P U及びファイルシステムによる相当量の処理を必要とする。

【0011】

【発明の開示】

データ転送の制御をホストに残しつつ、ネットワークと記憶装置間のデータ転送を加速するためのハードウェア及び処理メカニズムを提供する、ローカルホストのためのインテリジェントネットワークインタフェースカード( I N I C )のようなインタフェース装置が開示されている。このインタフェース装置は、ネットワークパケットヘッダを処理するためのハードウェア回路を有し、ネットワークと記憶装置間でのデータ転送のためのホストにより設定された専用の高速バスを使用することができる。ホストC P U及びプロトコルスタックにより、高速バス上でのデータ転送のためのプロトコル処理が回避され、ネットワーク及びストレージサブシステムの多くの要求からホストバスのバンド幅を解放する。独立ディスクの冗長アレイ( R A I D )または複数のドライブからなる他の構成を有することができる記憶装置が、S C S Iのようなパラレルチャネルによって、またはイーサネットもしくはファイバチャネルのようなシリアルチャネルによってインタフェース装置に接続することができ、かつインタフェース装置はP C IバスのようなI / Oバスによってローカルホストに接続することができる。追加の記憶装置を、S C S Iのようなパラレルインタフェースによりローカルホストに接続することができる。

【0012】

ファイルキャッシュがホストをバイパスすることができるデータを格納するためにインタフェース装置上に設けられ、インタフェース装置のファイルキャッシュにおけるデータの編成が、ホスト上のファイルシステムによって制御される。この構成によって、リモートホストと記憶装置間のデータ転送が、データがインタフェース装置とI / Oバス上のローカルホストの間を通過することなく、インタフェース装置の高速バス上で処理することができる。また、従来の通信プロトコル処理と対照的に、高速バスデータのための制御信号がホストメモリバス上を繰り返し移動して一時的に格納され、かつ次にホストC P Uにより一度に1レイヤが処理されるということがない。従って、ホストは、ホスト制御型記憶装置上でのファイル読みまたは書き込みのためのデータトラフィックの大部分に関与することから解放することができる。

【0013】

追加のインタフェース装置をI / Oバスを介してホストに接続することができ、追加の各インタフェース装置はホストファイルシステムにより制御されるファイルキャッシュを有し、かつ追加のネットワーク接続を提供しかつ/または追加の記憶装置に接続される。複数のインタフェース装置が単一のホストに接続されることによって、該ホストは複数のストレージネットワークを制御することができ、ホスト制御型ネットワークとの間のデータフローの大部分がホストプロトコル処理、I / Oバスを通る移動、ホストバスを通る移動、ホストメモリへの格納をバイパスする。ある実施例では、記憶装置がそのようなインタフェース装置にギガビットイーサネットネットワークにより接続することができ、前記欠陥なしでファイバチャネルの速度及びバンド幅を提供し、かつイーサネットベースのネットワークの大規模導入ベース及び互換性の利益が得られる。

【0014】

【発明を実施するための最良の形態】

本発明によるネットワークデータ通信システムの概要が図1に示されている。ホストコンピュータ20が、ローカルまたはワイドエリアネットワーク25もしくはインターネット28のようなネットワークに接続するための1つまたは複数のポートを有することができ

10

20

30

40

50

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）２２のようなインタフェース装置に接続されている。ホスト２０は、ホストバス３５によりホストメモリ３３に接続された中央処理装置（ＣＰＵ）３０のようなプロセッサを有し、図示しないオペレーティングシステムが、ファイルシステム２３を含む様々なタスク及び装置を監督するためにメモリ３３に存在する。同様にホストメモリ３３に格納されているのは、ネットワーク通信を処理するためのインストラクションのプロトコルスタック３８、及びＩＮＩＣ２２とプロトコルスタック３８間を通信するＩＮＩＣドライバ３９である。キャッシュマネージャ２６は、ファイルシステム２３及びＷｉｎｄｏｗｓ（登録商標）ＮＴまたは２０００の仮想メモリマネージャのような任意のメモリマネージャ２７の制御下で動作して、ファイルストリームと称されるファイル部分をホストファイルキャッシュ２４上に格納しかつ検索する。

10

#### 【００１５】

ホスト２０は、ＰＣＩバスのようなＩ／Ｏバス４０によってＩＮＩＣ２２に接続され、これはホストＩ／Ｏブリッジ４２によりホストバス３５に結合されている。前記ＩＮＩＣは、ＩＮＩＣバス４８により相互に接続されたインタフェースプロセッサ４４及びメモリ４６を有する。ＩＮＩＣバス４８は、ＩＮＩＣ・Ｉ／Ｏブリッジ５０を有するＩ／Ｏバス４０に結合されている。同様にＩＮＩＣバス４８に接続されているのは、ネットワークメッセージの上位レイヤプロセッシングを行うハードウェアシーケンサ５２のセット即ち集合である。ＬＡＮ／ＷＡＮ２５及びインターネット２８への物理的接続は、従来の物理的レイヤハードウェアＰＨＹ５８によって提供される。各ＰＨＹ５８装置は、媒体アクセスコントロール（ＭＡＣ）６０の対応する装置に接続され、各ＭＡＣ装置によって、前記ＩＮＩＣと前記ネットワークの１つとの間に従来型のデータリンクレイヤ接続が提供される。

20

#### 【００１６】

ディスクドライブまたはディスクドライブの集合と対応するコントローラのようなホスト記憶装置６６が、ＳＣＳＩアダプタのような従来のＩ／Ｏコントローラ６４によりＩ／Ｏバス４０と接続することができる。パラレルデータチャネル６２がコントローラ６４をホスト記憶装置６６に接続している。これに代えて、ホスト記憶装置６６は、独立ディスクの冗長アレイ（ＲＡＩＤ）とすることができ、かつＩ／Ｏコントローラ６４はＲＡＩＤコントローラとすることができ、ファイルシステム２３の命令下で動作するＩ／Ｏドライバ６７、例えばＳＣＳＩドライバモジュールがコントローラ６４と対話して、ホスト記憶装置６６へのデータの書込み及び読出しを行う。好ましくは、ホスト記憶装置６６が、ホストメモリ３３にキャッシュすることができる、ファイルシステム２３を含むホスト２０のオペレーティングシステムコードを有する。

30

#### 【００１７】

ディスクドライブまたはディスクドライブ及び対応するコントローラの集合のようなＩＮＩＣ記憶装置７０が、マッピングインタフェースコントローラ、ＩＮＩＣ・Ｉ／Ｏコントローラ７２を介してＩＮＩＣバス４８に接続され、これが次にパラレルデータチャネル７５により前記ＩＮＩＣ記憶装置に接続されている。ＩＮＩＣ・Ｉ／Ｏコントローラ７２は、ＳＣＳＩコントローラとすることができ、これはパラレルデータチャネル７５によりＩＮＩＣ記憶装置７０に接続される。これに代えて、ＩＮＩＣ記憶装置７０はＲＡＩＤシステムとすることができ、かつＩ／Ｏコントローラ７２をＲＡＩＤコントローラとし、多数のまたは分岐したデータチャネル７５とすることができ、同様に、Ｉ／Ｏコントローラ７２は、ＩＮＩＣ記憶装置７０のためのＲＡＩＤコントローラに接続されたＳＣＳＩコントローラとすることができ、別の実施例では、ＩＮＩＣ記憶装置７０は、ファイバチャネル（ＦＣ）ネットワーク７５に接続され、かつＩ／Ｏコントローラ７２はＦＣコントローラである。ＩＮＩＣ・Ｉ／Ｏコントローラ７２はＩＮＩＣバス４８に接続されたように図示されているが、Ｉ／Ｏコントローラ７２は、これに代えてＩ／Ｏバス４０に接続することができる。ＩＮＩＣ記憶装置７０は、任意により、そこからオペレーティングシステムの核がロードされる、ホスト２０のルートディスクを有することができる。ＩＮＩＣメモリ４６は、ＬＡＮ／ＷＡＮ２５のようなネットワークから受信しまたはそれに送信される

40

50



パケットを一時的に格納するためのフレームバッファ77を有する。また、I N I Cメモリ46は、I N I C記憶装置70に格納しまたはそこから検索するデータを一時的に格納するために、インタフェースファイルキャッシュ、I N I Cファイルキャッシュ80を有する。I N I Cメモリ46が図1においてそれを簡単にするために1個のブロックで記載されているが、メモリ46はI N I C22の様々な位置に配置された別個の装置で形成することができ、かつダイナミックランダムアクセスメモリ(D R A M)、スタティックランダムアクセスメモリ(S R A M)、リードオンリメモリ(R O M)及び他の形態のメモリで構成することができる。

#### 【0018】

ファイルシステム23は、記憶装置66、70及びファイルキャッシュ24、80の情報の編成に関する一般的な知識を有する高レベルのソフトウェアエンティティであり、ストレージアーキテクチャの特性及び性能を実行するアルゴリズムを提供する。ファイルシステム23は、ファイルの階層構造として、記憶装置66、70及び各ファイルキャッシュ24、80に格納されている情報を論理的に編成するが、このような論理ファイルは、記憶装置66または70の異なるディスクの本質的に異なるブロックに物理的に配置できるものである。また、ファイルシステム23は、記憶装置66、70及びファイルキャッシュ24、80上でのファイルデータの格納及び検索を管理する。前記ファイルシステム下でホスト20上で動作するI/Oドライバ67ソフトウェアが、各記憶装置66、70のコントローラ64、72と対話して、データのブロックを操作し、即ちそれらの記憶装置からデータブロックを読み出しまたはそれらに書き込む。ホストファイルキャッシュ24及びI N I Cファイルキャッシュ80によって、記憶装置66、70から読み出されまたはそれらに書き込まれるデータの記憶空間が提供され、前記データは記憶装置66、70の物理的ブロック形式とアプリケーションについて使用される論理的ファイル形式との間でファイルシステム23によりマッピングされる。ファイルに関連しかつホストファイルキャッシュ24及びI N I Cファイルキャッシュ80に格納されるバイトの線形ストリームをファイルストリームと称する。ホストファイルキャッシュ24及びI N I Cファイルキャッシュ80はそれぞれ、その対応するキャッシュに保持されるファイルストリームをリストするインデックスを有する。

#### 【0019】

ファイルシステム23は、記憶装置66、70上のファイルブロックのアドレスを決定するのに使用することができるメタデータを有し、最近アクセスされたファイルブロックのアドレスへのポインタがメタデータキャッシュにキャッシュされる。ファイルブロックへのアクセスが、例えばL A N / W A N 25上のリモートホストにより要求されると、ホストファイルキャッシュ24及びI N I Cファイルキャッシュ80のインデックスが最初に参照されて、前記ブロックに対応するファイルストリームがそれらのキャッシュに格納されているかどうかを見る。前記ファイルストリームがファイルキャッシュ24または80に見当たらない場合には、そのブロックへのリクエストが、前記メタデータにより示された適当な記憶装置のアドレスへ送られる。1つまたは複数の従来のキャッシングアルゴリズムをファイルキャッシュ24及び80のキャッシュマネージャ26が用いて、前記キャッシュが空で新しいデータをキャッシュすべき場合にどのデータを捨てるべきかを選択する。I N I Cファイルキャッシュ80上でのファイルストリームのキャッシングによって、I N I C記憶装置70に格納されたファイルブロックのI/Oパス40及びデータチャネル75双方におけるトラフィックが大幅に減少する。

#### 【0020】

ホスト20に送られたネットワークパケットがI N I C22に到着すると、そのパケットのヘッダがシーケンサ52により処理されて、前記パケットを有効にしかつ該パケットのサマリまたは記述子が作成され、該サマリは前記パケットにプリペンドされかつフレームバッファ77に格納され、前記パケットへのポインタが9に記憶される。前記サマリは、前記パケットヘッダのプロトコルタイプ及びチェックサムの結果を表す状態語である。この状態語に含まれるものは、前記フレームが高速パスデータフローの対象であるかどうか

10

20

30

40

50

を示すものである。従来の方法と異なり、トランスポート及びセッションレイヤの情報を  
含むプロトコル情報を有する上位レイヤヘッダが、シーケンサ52のハードウェア論理に  
より処理されて、サマリが作成される。前記シーケンサの専用論理回路によって、パケッ  
トヘッダを、該パケットがネットワークから到着するのと同様に高速で仮想的に処理する  
ことが可能になる。

#### 【0021】

次に、前記INICは、CPU30実行プロトコルスタック38によるヘッダの「低速パ  
ス」プロセッシングのために前記パケットをホストメモリ33に送るか、または「高速パ  
ス」に従って前記パケットデータを直接INICファイルキャッシュ80またはホストフ  
ァイルキャッシュ24のいずれかに送るかを選択する。前記高速パスは、シーケンシャル  
でエラーのないメッセージ毎に複数のパケットを有するデータトラフィックの相当大部分  
について選択することができ、データを繰り返しコピーしたり、ホストメモリバス35を  
繰り返し往来するような、前記CPUによる各パケットの時間のかかるプロトコルプロセ  
ッシングが回避される。パケットがINICファイルキャッシュ80内に直接移される高  
速パス状況の場合、ホストバス35とI/Oバス40における追加の往来が同様に回避  
される。低速パスプロセッシングは、INIC22の高速パスによって従来は転送されな  
いあらゆるパケットをホスト20により従来通り処理できるようにする。

#### 【0022】

ホスト20における高速パス能力を与えるためにリモートホストとの接続が最初に設定さ  
れ、それにはハンドシェーク、認証及び他の接続初期化の手順を含むことができる。通信  
制御ブロック(CCB)が、TCP/IPまたはSPX/IPXプロトコルにより代表さ  
れるような接続ベースのメッセージのための接続初期設定手順の際にプロトコルスタック  
38により作成される。前記CCBは、ソース及びデスティネーションアドレス並びにポ  
ートのような接続情報を有する。TCP接続の場合、CCBは、ソース及びデスティネ  
ーション媒体アクセスコントロール(MAC)アドレス、ソース及びデスティネーションI  
Pアドレス、ソース及びデスティネーションTCPポート、並びにスライディングウィン  
ドウプロトコルのための送信受信ウィンドウ及びタイマのようなTCP変数からなる。接  
続が設定された後、前記CCBはINICドライバ39により前記ホストからINICメ  
モリ46へ、そのメモリ46内にコマンドレジスタを書き込むことによって送られ、そこ  
では他のCCBと共にCCBキャッシュ74に格納することができる。また、前記INIC  
は、前記CCBをパケットサマリと加速して突き合わせるためにキャッシュしたCCB  
に対応するハッシュテーブルを作成する。

#### 【0023】

ファイル書込みのような前記CCBに対応するメッセージを前記INOCが受け取ると、  
該メッセージの最初のパケットのヘッダ部分がホスト20に送られて、CPU及びプロト  
コルスタック38により処理される。この前記ホストに送られたヘッダ部分は、前記パケ  
ットの一定のオフセットで始まることが知られている、前記メッセージのセッションレイ  
ヤヘッダを有し、かつ任意により前記パケットからいくつかのデータを含んでいる。プロ  
トコルスタック38のセッションレイヤによる前記セッションレイヤヘッダのプロセッシ  
ングによって、前記ファイルに属するものとして前記データが認識されかつ前記メッセ  
ージのサイズが表示され、これらが前記ファイルシステムにより使用されて、前記メッセ  
ージデータをホストファイルキャッシュ24またはINICファイルキャッシュ80にキャ  
ッシュするかどうか、及び選択されたファイルキャッシュ内の前記データのデスティネ  
ーションを確保するかどうか決定される。前記ホストに送られたヘッダ部分に何らかのデ  
ータが含まれていた場合には、それは前記デスティネーションに格納される。選択された  
前記ファイルキャッシュの前記デスティネーションのバッファアドレスのリストがINIC  
22に送られて、前記CCBにまたはそれと共に格納される。また、前記CCBは、メ  
ッセージの長さ及び処理されたパケットの数並びに順序のような、前記メッセージに関す  
る状態情報を保持し、どのユーザが関与しているか及び転送情報毎の記憶空間を含む、各  
プロトコルレイヤに関するプロトコル及びステータス情報を提供する。

## 【 0 0 2 4 】

前記 C C B が一旦デスティネーションを示すと、該 C C B に対応するパケットの高速バスプロセッシングが可能になる。次に受け取ったパケットをシーケンサ 5 2 により上述したようにプロセッシングしてパケットサマリを生成した後、前記パケットサマリのハッシュが前記ハッシュテーブルと比較され、かつ必要な場合には前記 C C B キャッシュ 7 4 に格納された C C B と比較され、高速バス接続が設定されているメッセージに前記パッケージが属するかどうかを決定する。パケットサマリが C C B に整合すると、他の例外的な条件が存在しないと仮定して、ネットワークまたはトランスポートレイヤヘッダ無しで、前記パケットのデータがダイレクトメモリアクセス ( D M A ) 装置 6 8 により、前記 C C B により示されたファイルキャッシュ 8 0 またはファイルキャッシュ 2 4 のデスティネーションに送られる。

10

## 【 0 0 2 5 】

前記メッセージからの全てのデータがファイルストリームとして I N I C ファイルキャッシュ 8 0 またはホストファイルキャッシュ 2 4 内にキャッシュされた後のいくつかの時点において、データのファイルストリームが、ファイルシステム 2 3 の制御下にある D M A 装置 6 8 により、そのファイルキャッシュから I N I C 記憶装置 7 0 またはホスト記憶装置 6 6 に、前記ファイルシステムの制御下で送られる。通例、ホストファイルキャッシュ 2 4 内にキャッシュされたファイルストリームは、I N I C 記憶装置 6 6 上に格納されるのに対し、I N I C ファイルキャッシュ 8 0 内にキャッシュされたファイルストリームは、I N I C 記憶装置 7 0 上に格納されるが、このような構成は必ずしも必要でない。ファイル転送に関する後のリクエストは、そのリクエストが同一のソース及びデスティネーション I P アドレス及びポートを含んでいると仮定して、同じ C C B により取り扱うことができ、書込みリクエストの最初のパケットは前記ホスト C P U により処理されて、前記メッセージを格納するためのホストファイルキャッシュ 2 4 または I N I C ファイルキャッシュ 8 0 内の位置を決定する。また、前記ファイルキャッシュをバイパスしてリモートホストから受け取ったメッセージからのデータを格納するためのデスティネーションとして I N I C 記憶装置 7 0 またはホスト記憶装置 6 6 上の位置を予定するように前記ファイルシステムを構成することが可能である。

20

## 【 0 0 2 6 】

本発明の基本的な理解を促すための概要が図 2 に記載されており、同図では、各パスについての情報についての基本タイプを示すことによって、図 1 のネットワークデータ記憶システムのための情報の流れに関する主要なパスを分離している。図 2 は、細い矢印を有する基本的に制御情報からなる情報のフローパス、太く白い矢印を有する基本的にデータからなる情報のフローパス、及び太く黒い矢印を有する制御情報及びデータ双方からなる情報のフローパスを示している。ホスト 2 0 が基本的に制御情報の流れに関係しているのに対して、I N I C 記憶装置 7 0 は基本的にデータ転送に関連していることに注意する。

30

## 【 0 0 2 7 】

L A N / W A N 2 5 のようなネットワークと I N I C 2 2 間の情報の流れは制御情報及びデータを含むことができ、従って太く黒い矢印 8 5 で示されている。L A N / W A N 2 5 のようなネットワークと I N I C 2 2 間の情報の流れ 8 1 の例は、ファイルの読出しまたは書込みに加えて、接続初期設定の対話及び肯定応答のような、制御情報に包含されるようなファイルデータを含むパケットとして送られる制御情報が含まれる。シーケンサ 5 2 がファイル書込みからの制御情報を処理し、かつデータ及び制御情報を I N I C フレームバッファ 7 7 との間でやりとりするが、従ってそれらの転送は太く黒い矢印 8 8 で表されている。フレームバッファ 7 7 に格納されるデータに関する制御情報は、細い矢印 9 0 で示すようにプロセッサ 4 4 により処理され、かつネットワーク接続初期設定パケット及びセッションレイヤヘッダのような制御情報は、太い矢印 9 2 で示すように、プロトコルスタック 3 8 に送られる。前記ホストにより接続が設定されると、その接続に関する、C C B のような制御情報を、細い矢印 9 4 で示すようにホストプロトコルスタック 3 8 と I N I C メモリ 4 6 との間で送ることができる。I N I C 記憶装置 7 0 から読み出されまた

40

50

は書き込まれるデータの一時的な格納が、太く白い矢印 96、98 で示されるように、I N I C ファイルキャッシュ 80 及びフレームバッファ 77 により行われる。I N I C ファイルキャッシュ 80 に格納されている全てのファイルストリームの制御及び知識が、細い矢印 91 で示すように、ファイルシステム 23 によって提供される。ホスト記憶装置 66 がネットワークアクセス可能なデータを格納しない実施例では、ファイルシステムの情報が、矢印 81 で示すように、ホストファイルキャッシュ 24 とホスト記憶装置 66 との間で送られる。同図に記載されていない別の実施例では、ホスト記憶装置を備えておらず、または基本的にネットワークファイル転送のためにホスト記憶装置及びホストファイルキャッシュを用いることができる。

#### 【0028】

図 2 からわかるように、ネットワークファイル読みまたは書き込みのデータは基本的に I N I C 22 を通過しかつホスト 20 を避けており、制御情報が基本的に前記ホストと I N I C 間で送られる。このネットワークと記憶装置間のファイル転送に関する制御情報の分離によって、前記ホストは、前記ネットワークと記憶装置との間で前記 I N I C を流れるファイル転送を管理することが可能になり、他方前記 I N I C がそれらデータ転送のための高速パスを提供し、それによってデータ処理能力が加速される。前記 I N I C のデータ高速パスにより得られる処理能力の増加によって、ホスト及び I N I C は、ファイルサーバとして機能することに加えて、例えば映像のような高バンド幅アプリケーションのデータベースサーバとして機能することができる。

#### 【0029】

図 3 は、ネットワークから受け取ったメッセージを格納するために図 1 のシステムが実行するいくつかのステップを示している。L A N / W A N 25 のようなネットワークから送られたパケットは、最初に P H Y 装置 58 により I N I C 22 において受け取られ 100、かつ M A C 装置 60 が前記パケットがホスト 20 にアドレス指定されていることを確認するようにリンクレイヤプロセッシングを実行する。次に、ネットワーク、前記パケットのトランスポート及び任意によりセッションレイヤヘッダがシーケンサ 52 により処理され 102、それにより前記パケットが有効になりかつそれらヘッダのサマリが作成される。次に、前記サマリをパケットに追加して、フレームバッファ 77 の 1 つに格納する。次にプロセッサ 44 が、前記パケットサマリを検査することにより、前記パケットが高速パスプロセッシングの対象であるかどうかを決定する 106。パケットが高速パスの対象であるかどうかは、単にサマリ内に示されているパケットのプロトコルによって決定することができる。パケットが高速パス対象でない場合には、そのパケットはプロトコルスタック 38 からの C P U 30 実行命令によって前記パケットのヘッダを処理するために、I / O バス 40 を介してホストメモリ 33 に送られる 108。

#### 【0030】

前記パケットは高速パス対象である場合には、前記パケットサマリが、それぞれに C C B として表される、前記カードによって取り扱われる高速パス接続の集合と前記サマリを C C B ハッシュ及び C C B キャッシュと突き合わせることによって比較する。前記サマリが前記 I N I C メモリに保持されている C C B と整合しない場合には、前記パケットは、プロトコルスタックからの C P U 実行命令により前記パケットのヘッダを処理するためにホストメモリに送られる 112。パケットが接続初期設定対話の一部である場合には、前記パケットを用いて前記メッセージについて C C B を作成することができる 115。前記パケットサマリがむしろ前記 I N I C メモリに保持されている C C B に整合する場合には、前記プロセッサは、例えば断片化されたまたは壊れたパケットを含む例外状態を検討し、そのような例外状態が見つかった場合には、プロトコルプロセッシングのために C C B 及び前記パケットをホストプロトコルスタック 38 にフラッシュする。パケットサマリは C C B に整合しているが該パケットのデスティネーションが C C B で示されていない場合には、前記パケットのセッションレイヤヘッダをホストプロトコルスタック 38 に送って、前記ファイルシステムに従って、ホストファイルキャッシュまたは I N I C ファイルキャッシュ内のデスティネーションを決定し 122、そのデスティネーションのキャッシュ

アドレスのリストをCCBと共にINICに格納する。また、前記INICはCCBをホスト116にフラッシュさせかつパケットをスタックによるプロセッシングのためにホストに送らせることになる例外状態について前記パケットサマリを検査する114。

#### 【0031】

パケットサマリがCCBに整合しかつ該パケットのデスティネーションが該CCBと共に格納され、かつ例外状態が存在しない場合には、前記パケットからのデータがDMAにより、CCBにより示されたホストファイルキャッシュまたはINICファイルキャッシュのデスティネーションに送られる125。この場合におけるメッセージパケットはホストプロトコルプロセッシングスタックによるヘッダのプロセッシングをバイパスし、高速パスデータ転送を提供する。パケットのデータが高速パスを介してINICファイルキャッシュ及びINICストレージに送られているような場合には、前記パケットはホストによるプロトコルプロセッシングを避けるだけでなく、前記I/Oバスまたはホストメモリバスを横切らないので、従来のネットワークストレージと比較して大幅に時間を節約するCPUプロセッシング及びバストラフィックが得られる。

#### 【0032】

##### 【発明の実施の形態】

図4は、LAN/WAN25のようなネットワークからのリクエスト200に回答して、ホスト記憶装置66またはINIC記憶装置70からファイルまたはファイルの一部を検索するために図1のシステムによって実行されるいくつかのステップを示している。まず、リクエストパケットをプロトコルスタックにより処理し、それにより前記リクエストをファイルシステムに送る。前記ファイルシステムは前記リクエストに示されたファイルを、該ファイルに対応するファイルストリームがINICファイルキャッシュまたはホストファイルキャッシュにキャッシュされているかどうかを決定すること、及び前記ファイルストリームが前記キャッシュのいずれかにも配置されていない場合には前記ファイルに対応するブロックが前記ホスト記憶装置またはINIC記憶装置に格納されているかどうかを決定することを含めて、位置付けする。ファイルストリームが前記キャッシュのいずれにも配置されていないと仮定した場合、前記ファイルブロックはホストファイルキャッシュに読み取られ204、またはINICファイルキャッシュに読み取られる206。大抵の場合には、前記ホスト記憶装置に格納されたファイルブロックはホストファイルキャッシュ内に読み取られ、かつINIC記憶装置に格納されたファイルブロックはINICファイルに読み取られることになるが、このマッピングは必然ではない。例えば、ホスト記憶装置に格納されたファイルブロックをINICファイルキャッシュ内に読み取り、それによってホストメモリバス上のトラフィックを減少させることが好ましい場合がある。

#### 【0033】

ファイルブロックがホストファイルキャッシュにキャッシュされていない場合には、前記ホストは前記リクエストに対応するCCBが前記INICにより保持されているかどうかを注意することによって、ファイルを高速パスプロセッシングにより送るかどうかを決定する210。前記ホストが高速パスを使うのではなく、低速パスによって前記ファイルをホストから送ることを選択した場合、前記ホストはプロトコルスタックを実行して、ホストファイルキャッシュに保持されているデータについてヘッダを作成し、かつ次に前記ヘッダ及びチェックサムをデータに追加して、従来と同様にINICによりネットワーク上を送信するためにネットワークフレームを作成する212。次に、INICはDMAを用いてホストから前記フレームを入手し214、かつ次にINICは前記フレームをネットワーク上に送る208。そうしないで前記ファイルを高速パスで送る場合には、INICプロセッサはCCBを用いてヘッダ及びチェックサムを作成し、ホストファイルキャッシュからのデータのフレームサイズのセグメントをDMAし、かつ次に前記ヘッダ及びチェックサムを前記データセグメントにプリペンドしてネットワークフレームを作成し218、前記ホストをプロトコルプロセッシングから解放する。

#### 【0034】

同様にして、ファイルブロックがINICファイルキャッシュ内にキャッシュされた場合

には、前記ホストがＣＣＢがＩＮＩＣにより保持されているかどうかを注意することによって、高速バスプロセッシングにより前記ファイルを送るかどうかを決定する２２０。前記ホストが高速バスを使用しないことを選択した場合２２２、ホストＣＰＵは、ファイルブロックデータについてヘッダ及びチェックサムを準備し、前記ヘッダをホストメモリに格納する。次に前記ホストは、ＩＮＩＣに指示して、ホストメモリからのヘッダをＩＮＩＣメモリ内のデータにプリペンドすることによってネットワークフレームをアSEMBリし、次にＩＮＩＣによりネットワークを介して送られるメッセージフレームを作成する。この高速バスを使用しない場合でさえ、データはＩ／Ｏバスを通してホストに送られたりＩＮＩＣから戻されることが無く、Ｉ／Ｏバスまたはネットワークによりホストに接続された記憶装置上に配置されたファイルブロックの送信という従来方法と比較して、Ｉ／Ｏトラフィックが少なくなる。その代わりに高速バスが選択された場合２２５、ＩＮＩＣプロセッサはＣＣＢに対応するヘッダ及びチェックサムを作成し、かつ前記ヘッダ及びチェックサムをＩＮＩＣファイルキャッシュからのデータセグメントにプリペンドしてネットワークフレームを作成し、これは次にＩＮＩＣによってネットワーク上を送信される。この高速バスの場合には、前記ホストが、Ｉ／Ｏバストラフィックから解放されることに加えて、プロトコルプロセッシング及びホストメモリバストラフィックから解放される。

#### 【００３５】

図５は、ホスト２０がＩ／Ｏバス４０を介して、第１ＩＮＩＣ２２に加えていくつかのＩ／Ｏ・ＩＮＩＣに接続されたネットワーク記憶システムを示している。本実施例の各ＩＮＩＣは少なくとも１つのネットワークに接続されかつ少なくとも１つの記憶装置が接続されている。かくして第２ＩＮＩＣ３０３はＩ／Ｏバス４０に接続され、第２ネットワーク３０５への前記ホストのインタフェースを提供する。第１ＩＮＩＣ２２は、上述したように実際はいくつかのネットワークに接続するためにいくつかのネットワークポートを持つことができ、同様に第２ＩＮＩＣ３０３も１つのネットワーク３０５以外のネットワークに接続することができるが、この図面では簡単化するために、各ＩＮＩＣについて１つのネットワーク接続のみが示されている。第２ＩＮＩＣ３０３は、第２ＩＮＩＣ記憶装置３０８に接続されたＳＣＳＩアダプタのようなＩ／Ｏコントローラを有する。これに代えて、第２ＩＮＩＣ記憶装置３０８はＲＡＩＤシステムとすることができ、かつ第２ＩＮＩＣ３０３はＲＡＩＤコントローラを有しまたはそれと接続することができる。別の実施例では、第２ＩＮＩＣ３０３が、ＦＣネットワークループまたはＦＣアダプタ及びネットワークラインによって第２ＩＮＩＣ記憶装置３０８に接続されたＦＣコントローラを有する。第Ｎ・ＩＮＩＣ３１０で示されるように、Ｎ個のＩＮＩＣをホスト２０にＩ／Ｏバスを介して接続することができる。第Ｎ・ＩＮＩＣ３１０は、第Ｎネットワーク３１３及び第Ｎ記憶装置３１５にそれぞれ接続されたネットワークインタフェース及びストレージインタフェースを提供する回路及び制御インストラクションを有する。第Ｎ・ＩＮＩＣ３１０はいくつかのネットワークに接続するためにいくつかのネットワークポートを持つことができ、かつ第２の第Ｎ・ＩＮＩＣ３１０が同様に２以上のネットワークに接続することができる。第Ｎ・ＩＮＩＣ３１０のストレージインタフェース回路及び制御インストラクションは、例えば、ＳＣＳＩケーブルによって第Ｎ・ＩＮＩＣ記憶装置３１５に接続されたＳＣＳＩコントローラを含むことができる。これに代えて、第Ｎ・ＩＮＩＣ記憶装置３１５はＲＡＩＤシステムとすることができ、かつ第Ｎ・ＩＮＩＣ３１５はＲＡＩＤコントローラを有しまたはそれに接続することができる。さらに別の実施例では、第Ｎ・ＩＮＩＣ３１０が、ＦＣアダプタ及びＦＣネットワークラインによって第Ｎ・ＩＮＩＣ記憶装置３１５に接続されたＦＣコントローラを有する。

#### 【００３６】

前記ファイルシステムは、ネットワークアクセス可能ファイルが、ホスト記憶装置６６ではなく、いずれか１つのネットワーク記憶装置６６、３０５、３１５に格納されるように構成することができ、その代わりにそれはファイルシステムコード及びプロトコルスタックを含み、それらの複製が前記ホストにキャッシュされる。この構成によって、ホスト２０はネットワークファイル転送を制御するが、それらファイル内のデータの相当大部分は

10

20

30

40

50

、前記ホストに入ることすらなく、INIC 22、203、310を通して高速バスにより転送することができる。ファイルブロックがネットワークと同じINICに接続された記憶装置との間で転送される場合には、ファイルデータは決してI/Oバスまたはホストメモリバスを横切ることがない。ファイルがネットワークと異なるINICに接続された記憶装置との間で転送される場合には、ファイルブロックはDMAによって、記憶装置に接続されたINIC上のファイルキャッシュとの間でI/Oバスを通して送ることができ、依然としてホストメモリバスを回避する。この通常は回避される最悪の場合では、データはあるINICから別のものに転送することができ、それには、2回のI/Oバス転送や従来通りの繰り返しのホストメモリバスへのアクセスではなく、それに伴うI/Oバスの転送が1回であり、同様にホストメモリバスが回避される。

10

#### 【0037】

図6は、上述したINIC・I/Oコントローラ72を必要とすることなく、ネットワーク通信接続及びネットワークストレージ接続を提供するINIC 400を含むネットワークストレージシステムを示している。簡潔にするために、ホスト20及び関連する要素は、図1から変更はないものとして図示されているが、これは必ずしも当てはまらない。本実施例のINICは、第1LAN 414、第2LAN 416、第1SAN 418及び第2SAN 420に接続されたネットワーク接続またはポートを有する。ネットワーク414、416、418、420のいずれかまたは全部は、イーサネット、ファーストイーサネットまたはギガビットイーサネット規格に従って動作することができる。例えば802.3z及び802.3ab規格で記述されるギガビットイーサネットは、1ギガビット/秒または10ギガビット/秒のデータ転送速度、もしくは将来的におそらくはより速い速度を提供する。SAN 418、420は、TCP/IPまたはSCSIカプセル化プロトコルにおいてSCSIのようなストレージプロトコルを実行することができる。そのようなストレージプロトコルの1つが、J. Satran他によりthe Internet Engineering Task Force (IETF)のthe Internet-Draftにおいて「iSCSI (Internet SCSI)」の表題を付して2000年6月に、それより早い時期のInternet-Draftでは「SCSI/TCP (SCSI over TCP)」の表題を付して2000年2月に記載されている。イーサネットストレージと称される、このような別のプロトコルがセッションレイヤにおいてSCSIカプセル化プロトコル(SEP)を使用し、かつ基本的にTCPが使用されているWANまたはインターネット上でデータが転送されているかどうか、またはSTPが使用されるLANまたはSAN上でデータが転送されているかによって、トランスポートレイヤにおいてTCPまたはSANTランスポートプロトコル(STP)を用いることを提案している。

20

30

#### 【0038】

ホスト20は、PCIバスのようなI/Oバス40によってINIC 400に接続され、これはPCIバスインタフェースのようなINIC・I/Oブリッジ406によってINICバス404に接続されている。INIC 400は、I/Oバス40によってINICメモリ410に接続された専用プロセッサ408を有する。INICメモリ410は、フレームバッファ430及びINICファイルキャッシュ433を有する。また、INICバス404には、ネットワーク、トランスポート及びセッションレイヤプロセッシングを含むネットワークメッセージのプロセッシングを提供するハードウェアシーケンサ412の集合が接続されている。LAN 414、416及びSAN 418、420への物理的接続は、従来の物理的レイヤハードウェアPHY 422によって提供される。各PHY 422装置は、媒体アクセスコントロール(MAC)424の対応する装置と接続され、MAC装置はそれぞれINIC 400と1つのネットワークとの間にデータリンクレイヤ接続を提供する。

40

#### 【0039】

図7は、INIC 400と第1イーサネット-SCSIアダプタ452、第2イーサネット-SCSIアダプタ454、及び第3イーサネット-SCSIアダプタ456との間に接続されたギガビットイーサネットライン450を有するSAN 418を示している。イ

50

ーサネット - S C S I アダプタ 4 5 2、4 5 4、4 5 6 は T C P 接続を形成しかつ遮断することができ、I N I C 4 0 0 に S C S I コマンドを送りまたはそれから受け取り、かつ 4 5 0 を介して I N I C 4 0 0 にデータを送りまたはそれから受け取ることができる。第 1 記憶装置 4 6 2 が第 1 S C S I ケーブル 4 5 8 によって第 1 イーサネット - S C S I アダプタ 4 5 2 に接続されている。同様に、第 2 記憶装置 4 6 4 が第 2 S C S I ケーブル 4 5 9 によって第 2 イーサネット - S C S I アダプタ 4 5 4 に接続され、かつ第 3 記憶装置 4 6 6 が第 3 S C S I ケーブル 4 6 0 によって第 3 イーサネット - S C S I アダプタ 4 5 6 に接続されている。記憶装置 4 6 2、4 6 4、4 6 6 は S C S I 規格に従って各アダプタ 4 5 2、4 5 4、4 5 6 によって動作する。各記憶装置は、それぞれのアダプタにディジーチェーン式に接続された多数のディスクドライブを含むことができる。

10

#### 【 0 0 4 0 】

図 8 は、第 1 イーサネット - S C S I アダプタ 4 5 2 を詳細に示しており、これは本実施例において図 1 に示すものと類似する I N I C である。アダプタ 4 5 2 は、従来の P H Y 4 7 0 によって提供されるネットワークライン 4 5 0 との物理的レイヤ接続及び従来の M A C 4 7 2 によって提供される媒体アクセスを有する 1 個のネットワークポートを有する。上位レイヤプロセッシングを含むネットワークメッセージのパケットヘッダのプロセッシングは、アダプタバス 4 7 7 を介してプロセッサ 4 8 0 及びアダプタメモリ 4 8 2 に接続されたシーケンサ 4 7 5 によって提供される。アダプタメモリ 4 8 2 はフレームバッファ 4 8 4 及びファイルキャッシュ 4 8 6 を有する。また、アダプタバス 4 7 7 には S C S I コントローラが接続され、かつこれは S C S I チャンネル 4 5 8 によって第 1 記憶装置 4 6 2 に接続されている。アダプタ 4 5 2 と I N I C 2 0 との 1 つの違いは、アダプタ 4 5 2 が低速バスメッセージを処理するために C P U 及びプロトコルスタックを有するホストに必ずしも接続されていないことである。本実施例において接続の設定は、例えば接続初期設定対話の際にアダプタ 4 5 2 に初期パケットを I N I C 4 0 0 が送ることによって、アダプタ 4 5 2 により行われ、前記パケットはシーケンサ 4 7 5 により処理されかつプロセッサ 4 8 0 に送られて C C B を作成する。ソフトウェアプロトコルスタックを実行する C P U による低速バスプロセッシングを必要とする一定の状態は、このようなアダプタ 4 5 2 と I N I C 4 0 0 間の通信環境ではなお一層起こりにくいものである。アダプタ 4 5 2 と I N I C 4 0 0 間で送られるメッセージは、S C S I / T C P 及び簡易ネットワーク管理プロトコル ( S N M P ) のような単一のまたは限られたプロトコルレイヤの集合に従って構成することができ、かつ単一のソースへまたはそれから単一のもしくは制限された数のデスティネーションに送られる。従来の通信ネットワークに複雑な状態を生じさせる変数の多くを減らすことによって、高速バスプロセッシングの使用を増やすことができ、エラープロセッシングのためのアダプタ 4 5 2 の必要性が少なくなる。アダプタ 4 5 2 は、I P 及び T C P 上でいくつかのタイプのストレージプロトコルを処理する能力を持たせることができるが、それはアダプタ 4 5 2 が、I N I C 4 0 0 に接続される代わりに、ネットワークストレージのためにそのようなプロトコルの 1 つを使用するホストに接続されるような場合である。ネットワーク 4 5 0 がストレージ転送に専用の S A N ではないが通信トラフィックを取り扱うような場合、低速バスパケットのためにプロトコルスタックを実行する C P U を有するホストに接続された I N I C は、アダプタ 4 5 2 の代わりに用い

20

30

40

#### 【 0 0 4 1 】

図 9 に示すように、I N I C 4 0 0 に類似する追加の I N I C が、I / O バス 4 0 を介してホスト 2 0 に接続され、追加の各 I N I C によって追加の L A N 接続が提供されかつ / または追加の S A N に接続されている。複数の I N I C は、第 N ・ I N I C 4 9 0 で表されており、これは第 N ・ S A N 4 9 2 及び第 N ・ L A N 4 9 4 に接続されている。ホスト 2 0 に接続された複数の I N I C によって、前記ホストは複数のストレージネットワークを制御することができ、ホスト制御ネットワークとの間で行き来するデータの流れの相当大部分がホストプロトコルプロセッシング、I / O バスを介しての移動、ホストバスを介しての移動及びホストメモリへの格納をバイパスする。

50



## 【 0 0 4 2 】

図 1 の I N I C 2 2 がネットワーク 2 5 のようなネットワークから受信するメッセージパケットのプロセッシングが、図 1 0 により詳細に示されている。受け取ったメッセージパケットは、まず媒体アクセスコントローラ 6 0 に入り、これがネットワークへの I N I C アクセス及びパケットの受け取り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ことができる。そこから、データは一度に 1 バイトづつ、本実施例では 1 2 8 ビットの幅を有するアセンブリレジスタ 5 0 0 内に流れる。このデータは、図 1 1 に関してより詳細に説明するように、フライバイシーケンサ 5 0 2 によって分類され、これはパケットのバイトをそれらが近くを飛ぶごとに調査し、それらのバイトから前記パケットを要約されるのに用いられることになるステータスを生成する。このように作成されたステータスは、マルチプレクサ 5 0 5 により前記データと併合され、その結果得られたデータが S R A M 5 0 8 に格納される。パケット制御シーケンサ 5 1 0 がフライバイシーケンサ 5 0 2 を監督し、媒体アクセスコントローラ 6 0 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、かつアセンブリレジスタ 5 0 0 から S R A M 5 0 8 及び最終的に D R A M 5 1 2 へのデータの移動を管理する。パケット制御シーケンサ 5 1 0 は、S R A M コントローラ 5 1 5 を介して S R A M 5 0 8 のバッファを管理し、またデータを S R A M 5 0 8 から D R A M 5 1 2 のバッファに移動させる必要がある時、D R A M コントローラ 5 1 8 に知らせる。前記パケットに関するデータの移動が一旦完了しかつ全データが D R A M 5 1 2 のバッファに移動すると、パケットコントロールシーケンサ 5 1 0 は、フライバイシーケンサ 5 0 2 によって生成されたステータスを、パケットデータにプリペンドされるように S R A M 5 0 8 に及び D R A M 5 1 2 バッファの開始部分に移動させる。次に、パケット制御シーケンサ 5 1 0 は、キューマネージャ 5 2 0 に受信バッファ記述子を受信 9 に入力するように要求し、これが次にプロセッサ 4 4 にパケットがハードウェア論理及びその要約されたステータスによって処理されたことを通知する。

## 【 0 0 4 3 】

図 1 1 は、フライバイシーケンサ 5 0 2 がいくつかの階層を有することを示しており、各階層は一般にパケットヘッダの特定部分及び特定のプロトコルレイヤに、そのレイヤに関連するステータスを生成するために、焦点を合わせている。本実施例のフライバイシーケンサ 5 0 2 は、媒体アクセス制御シーケンサ 5 4 0、ネットワークシーケンサ 5 4 2、トランスポートシーケンサ 5 4 6 及びセッションシーケンサ 5 4 8 を有する。より上位のプロトコルレイヤに関係するシーケンサを追加で設けることができる。フライバイシーケンサ 5 0 2 は、パケット制御シーケンサ 5 1 0 によってリセットされ、かつ一定のバイトがアセンブリレジスタ 5 0 0 から利用可能であるかをフライバイシーケンサに知らせるパケット制御シーケンサによってポイントが与えられる。媒体アクセス制御シーケンサ 5 4 0 は、バイト 0 ~ 5 を見ることによってパケットが別のホストではなくまたは別のホストに加えてホスト 2 0 にアドレス指定されていることを判断する。また、パケットのオフセット 1 2、1 3 が媒体アクセス制御シーケンサ 5 4 0 により処理されて、タイプフィールド、例えば前記パケットがイーサネットまたは 8 0 2 . 3 であるかどうかを決定する。タイプフィールドがイーサネットの場合、それらのバイトは更に前記パケットのネットワークプロトコルタイプを媒体アクセス制御シーケンサ 5 4 0 に知らせる。8 0 2 . 3 の場合には、そうではなく、それらのバイトはフレーム全体の長さを示し、媒体アクセス制御シーケンサ 5 4 0 が前記パケット内に更に 8 バイトを調べて、ネットワークレイヤのタイプを決定することになる。

## 【 0 0 4 4 】

大部分のパケットについて、ネットワークシーケンサ 5 4 2 が、受け取ったヘッダの長さが正しい長さであるかを検証し、かつネットワークレイヤのヘッダをチェックサムする。高速パスの対象について、ネットワークレイヤのヘッダは、媒体アクセス制御シーケンサ 5 4 0 により行われる分析から、I P または I P X であることがわかる。例えば、タイプフィールドが 8 0 2 . 3 でありかつネットワークプロトコルが I P であると仮定すると、ネットワークシーケンサ 5 4 2 は、I P タイプを決定するために、バイト 2 2 で始まるこ

10

20

30

40

50

とになるネットワークレイヤヘッダの最初のバイトを分析する。IPヘッダの最初のバイトは、ネットワークシーケンスにより処理されて、前記パケットが関連するIPタイプを決定することになる。前記パケットが例えばIPバージョン4に関するものであることを決定することで、ネットワークシーケンス542による処理が更に行われ、該ネットワークシーケンスは前記パケットのトランスポートヘッダプロトコルを表示するためにIPヘッダ内への10バイトに配置されたプロトコルタイプをも見る。例えば、イーサネット上でのIPの場合、IPヘッダはオフセット14で始まり、かつプロトコルタイプバイトはオフセット23であり、これはネットワーク論理により処理されて、トランスポートレイヤプロトコルが例えばTCPであるかどうかを決定する。一般に20～40バイトであるネットワークレイヤヘッダの長さから、ネットワークシーケンス542は、トランスポートレイヤヘッダを検証するためにパケットのトランスポートレイヤヘッダの開始部分を決定する。トランスポートシーケンス546は、トランスポートレイヤヘッダ及びデータのチェックサムを生成することができ、これには、少なくともTCPの場合にIPヘッダからの情報を含むことができる。

10

#### 【0045】

TCPパケットの例について続けて説明すると、トランスポートシーケンス546はまた、ヘッダのトランスポートレイヤ部分の最初のいくつかのバイトを分析して、前記パケットがNet Biosまたは別のプロトコルであるかどうかのようなメッセージについてTCPソース及びデスティネーションポートを部分的に決定する。TCPヘッダのバイト12をトランスポートシーケンス546により処理されて、TCPヘッダの長さを決定しつつ検証する。TCPヘッダのバイト13は、認証フラグ及びプッシュフラグは別として、プロセッサにこのパケットを例外として分類させることがある、リセット及び終了のような予期しないオプションを示すことがあるフラグを有する。TCPオフセットバイト16、17は、ハードウェア論理によって取り出されかつ格納されるチェックサムであり、前記フレームの残りの部分は該チェックサムに対して検証される。

20

#### 【0046】

セッションシーケンス548は、セッションレイヤヘッダの長さを決定し、これは、Net Biosの場合にはたった4つのバイトであり、その内の2つがNet Biosペイロードデータの長さを知らせるが、これは他のプロトコルについてより大きくすることができる。また、セッションシーケンス548を用いて、例えば、特に高速パスが有利であるような読出しまたは書込みとしてメッセージのタイプを分類することができる。メッセージのタイプによって、更に上位レイヤの論理プロセッシングが、パケット制御シーケンス510及びフライバイシーケンス502のハードウェア論理によって実行することができる。このようにして、シーケンス52は、バイトの単一のストリームから選択したバイトの分類によってヘッダのハードウェアプロセッシングを知的に行い、パケットのステータスがオンザフライで決定された分類から構築される。全てのパケットがフライバイシーケンス502により処理されたことをパケット制御シーケンス510が一旦検出すると、パケット制御シーケンス510は、フライバイシーケンス502により生成されたステータス情報及びパケット制御シーケンス510により生成されたあらゆるステータス情報を付け加え、かつプロセッサ44によるパケットの取り扱いに好都合なように、それらのステータス情報をパケットにプリペンド（前部に追加）する。パケット制御シーケンス510により生成された追加のステータス情報には、媒体アクセスコントローラ60のステータス情報及び発見されたあらゆるエラー、またはアセンブリレジスタまたはDRAMバッファにおけるデータのオーバフロー、または前記パケットに関する他の様々な情報が含まれる。また、パケット制御シーケンス510は、キューマネージャ520を介して受信バッファ9及び受信統計9への入力を記憶する。

30

40

#### 【0047】

ハードウェア論理によりパケットを処理することの利点は、従来のシーケンシャルソフトウェアプロトコルプロセッシングと対照的に、パケットを格納したり、移動したり、コピーしたり、各プロトコルレイヤヘッダを処理するためにストレージから引き抜いたりする

50

必要がなく、処理効率の大幅な増加及び各パケットについて処理時間の節約が得られることである。パケットは、ネットワークから受け取ったレートビットで、例えば100ベースT接続の場合に100メガビット/秒で処理することができる。このレートで受信されかつ60バイトの長さを有するパケットを分類するための時間は、従って約5ミリ秒である。このパケットをハードウェア論理で処理しかつパケットデータを高速バスを介してそのホストデスティネーションに送るための合計時間は、CPUの割込みの減少で得られる追加の時間の節約及びホストバスバンド幅の節約を考慮することさえ無く、従来のシーケンシャルソフトウェアプロトコルプロセッシングを用いる従来のCPUが必要とするものよりも、大幅に少なくすることができる。デスティネーションがINICキャッシュにあるような場合には、ホストバス35及びI/Oバス40について追加のバンド幅の節約が得られる。

10

#### 【0048】

プロセッサ44は、フレームバッファ77に保持される各受信メッセージパケットについて、そのパケットが高速バスの対象であるかどうかを選択し、かつそうである場合には、前記パケットが属する接続について高速バスが既に接続されているかどうかを調べる。これをするために、プロセッサ44はまずヘッダのステータスサマリを調べて、パケットのヘッダが高速バスの対象について定義されたプロトコルからなるものであるかどうかを決定する。そうでない場合には、プロセッサ44は、INIC22内のDMAコントローラに命令して、前記パケットを低速バスプロセッシングのためのホストに送る。メッセージの低速バスプロセッシングの場合でさえ、このようにINIC22は、メッセージタイプ

20

#### 【0049】

高速バス対象の場合、プロセッサ44は、ヘッダのステータスサマリがINICにより保持されてるCCBと一致するかどうかを調べる。そうである場合、前記パケットからのデータは、高速バスに沿って前記ホストのデスティネーション168に送られる。高速バス対象のパケットサマリがINICにより保持されるCCBと一致しない場合には、前記パケットは低速バスプロセッシングのために前記ホストに送られ、該メッセージのCCBを作成する。また、断片化されたメッセージや他の複雑さがある場合には高速バスが用いられないことがある。しかしながら、相当大部分のメッセージについては、INIC高速バスによってメッセージのプロセッシングを大幅に加速することができる。このようにINIC22は、与えられたパケットの運命を決定するためにいくつかのプロトコルレイヤのそれぞれにステートマシンを用いる従来のものと対照的に、オンザフライで集めた情報に基づいて、データをそのデスティネーションに直接送るかどうかを決定する単一のステートマシンプロセッサ44が提供される。

30

#### 【0050】

CCBをINICのハッシュテーブルにキャッシングすることによって、到来するパケットを要約する語との素早い比較が行われ、パケットを高速バスを介して処理できるかどうか決定される一方、一杯のCCBがプロセッシングのためにINICに保持される。この比較を加速する別の方法には、B-ツリーのようなソフトウェア処理またはコンテンツアドレス可能メモリ(CAM)のようなハードウェアアシストが含まれる。INICマイクロコードまたはコンパレータ回路がCCBとの整合を検出すると、DMAコントローラが前記パケットからのデータをホストメモリ33またはINICファイルキャッシュ80のデスティネーションにCPUによる割込み、プロトコルプロセッシングまたは複製なしで配置する。受信したメッセージのタイプによって、前記データのデスティネーションはホスト20、またはホストファイルキャッシュ24、INICファイルキャッシュ80のセッションレイヤ、プレゼンテーションレイヤまたはアプリケーションレイヤであったりする。

40

#### 【0051】

ファイル転送のような大きなメッセージについて最も一般的に使用されているネットワー

50

クプロトコルの1つは、TCP/IP上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プリンタまたはディスクのような特定の動作のために必要なリソースが前記動作が行われたホストに存在するかまたはそれに関連するかもしくはファイルサーバのようなネットワークに接続された別室のホストに配置されているかどうかを決定するリダイレクタソフトウェアと協力して動作することができる。SMB及びサーバフラッシュリダイレクタは、通常トランスポートレイヤからサービスを受ける。本発明では、SMB及びリダイレクタはそうではなく、INICからサービスを受けることができる。この場合、大きなSMBトランスアクションを受け取った時、INICバッファからDMAコントローラによってデータを送ることは、ホストが取り扱わなければならない割込みを大幅に少なくすることができる。更に、このDMAは、通常データをそのホストファイルキャッシュ24またはINICファイルキャッシュ80内のデスティネーションに移動させ、そこから次にブロックでホスト記憶装置66またはINIC記憶装置70へそれぞれフラッシュされる。

10

#### 【0052】

SMB高速バス電送は一般に上述したSMB高速バス受信を逆にし、データのブロックがホスト記憶装置66またはINIC記憶装置70からホストファイルキャッシュ24またはINICファイルキャッシュ80へそれぞれ読み取られる一方、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リモートホストへの伝送のために、INICにより前記データにプリペンドされる。カスタムハードウェア経由かつホストの繰り返し割込み無しにINICにより複数パケット及び複数TCP、IP、NetBios及びSMBプロトコルレイヤのプロセッシングによって、SMBメッセージをネットワークラインに伝送する速度を大幅に増加させることができる。図4に関連して上述したように、送信したファイルブロックがINIC記憶装置70上に格納される場合には、ホストバス35バンド幅及びI/Oバスバンド幅40における追加の節約を得ることができる。

20

#### 【0053】

図12は、そのいずれも同図には記載されていないが、ネットワークメッセージを処理するために、INICと共にホストにより使用されるAlacritechのプロトコルスタック38を示している。INICドライバ560はINICをホストのオペレーティングシステムにリンクさせ、INICとプロトコルスタック38間で通信を送ることができる。本実施例のプロトコルスタック38はデータリンクレイヤ562、ネットワークレイヤ564、トランスポートレイヤ566、上位レイヤインタフェース568及び上位レイヤ570を有する。上位レイヤ570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よって、セッションレイヤ、プレゼンテーションレイヤ及び/またはアプリケーションレイヤを表すことができる。プロトコルスタック38は低速バスのパケットヘッダを処理し、接続を作りかつ壊し、INICへの高速バス接続のためにCCBを分配し、かつINICからホスト20へフラッシュされる高速バス接続のためにCCBを受け取る。上位レイヤインタフェース568は一般に、データリンクレイヤ562、ネットワークレイヤ564及びトランスポートレイヤ566により作成された接続及びステータス情報に基づいてCCBを組み立てること、及びINICデバイスドライバ560経由CCBをINICに分配すること、またはINICデバイスドライバ560経由INICからフラッシュされたCCBを受け取ることに責任を有する。

30

40

#### 【0054】

図13は、マイクロソフト(登録商標)のオペレーティングシステムと協力してネットワーク通信を処理するための複数のプロトコルスタックを有するAlacritechのプロトコルスタック38の別の実施例を示している。従来のマイクロソフトのTCP/IPプロトコルスタック580はMACレイヤ582、IPレイヤ584及びTCPレイヤ586を有する。コマンドドライバ590がホストスタック580と協力して動作し、ネットワークメッセージを処理する。コマンドドライバ590はMACレイヤ592、IP594及びAlacritechのTCP(APCP)レイヤ596を有する。従来のスタック580とコマンドドライバ590とが、INICデバイスドライバ570と対話するネットワークドライバ

50

インタフェース仕様 ( N D I S ) レイヤ 5 9 8 を共有する。 I N I C デバイスドライバ 5 7 0 は、従来のホストスタック 5 8 0 または A T C P ドライバ 5 9 0 によって処理するために受信表示をソートする。 P D I フィルタドライバ及び上位レイヤインタフェース 5 7 2 は同様に、 T D I ユーザ 5 7 5 からネットワークに送られるメッセージがコマンドドライバにかつ恐らくは I N I C の高速バスへ向きを変えられるか、またはホストスタックにより処理されるかを決定する。

#### 【 0 0 5 5 】

図 1 4 は、共に I N I C を有するサーバ 6 0 0 とクライアント 6 0 2 との間の S M B エクステンションを示しており、各 I N I C は、ギガビットイーサネットコンプライアントで有り得るネットワーク 6 0 4 上でのデータの高速バス移動のためにその接続及びステータスを定義する C C B を有する。クライアント 6 0 2 は I N I C 6 0 6 、 8 0 2 . 3 コンプライアントデータリンクレイヤ 6 0 8 、 I P レイヤ 6 1 0 、 T C P レイヤ 6 1 1 、 A T C P レイヤ 6 1 2 、 N e t B i o s レイヤ 6 1 4 及び S M B レイヤ 6 1 6 を有する。前記クライアントは、通信処理のために低速バス 6 1 8 及び高速バス 6 2 0 を有する。同様に、サーバ 6 0 0 は I N I C 6 2 2 、 8 0 2 . 3 コンプライアントデータリンクレイヤ 6 2 4 、 I P レイヤ 6 2 6 、 T C P レイヤ 6 2 7 、 A T C P レイヤ 6 2 8 、 N e t B i o s レイヤ 6 3 0 及び S M B レイヤ 6 3 2 を有する。サーバ直接接続型記憶装置 6 3 4 が、 S C S I チャンネルのようなパラレルチャンネル 6 3 8 上でサーバと接続される。これは I N I C 6 2 2 にも接続された I / O バス 6 3 9 に接続されている。ネットワーク記憶装置 6 4 0 がネットワークライン 6 4 4 上で I N I C 6 2 2 に接続され、かつ N A S 記憶装置 6 4 2 が、ギガビットイーサネットコンプライアントで有り得る同じネットワーク 6 4 4 に接続されている。サーバ 6 0 0 は、図 1 4 に示されていないファイルキャッシュと I N I C 6 2 2 間で I / O バス 6 3 8 上を通過する通信プロセッシングのための低速バス 6 4 6 及び高速バス 6 4 8 を有する。

#### 【 0 0 5 6 】

ネットワーク記憶装置 6 4 0 または 6 4 2 とクライアント 6 0 2 との間で転送される、 I / O バスを通らないデータのために、サーバの制御下で I N I C 6 2 2 によりストレージ高速バスが提供される。データは S C S I / T C P または I S C S I のようなブロック形式に従って I N I C 6 2 2 とネットワーク記憶装置 6 4 0 との間で通信され、データは T C P / N e t B i o s / S M B のようなファイル形式に従って I N I C 6 2 2 と N A S 記憶装置 6 4 2 との間で送信される。いずれのストレージ高速バスについても、 I N I C 6 2 2 が、記憶装置 6 4 0 または 6 4 2 との接続を定義する別の C C B を保持することができる。以下の説明を簡便にするために、そのサーバ 6 0 0 とのネットワーク 6 0 4 上での接続を定義する I N I C 6 0 6 により保持される C C B をクライアント C C B と称し、そのクライアント 6 0 2 とのネットワーク 6 0 4 上での接続を定義する I N I C 6 2 2 により保持される C C B をサーバ C C B と称する。そのネットワーク記憶装置 6 4 0 とのネットワーク 6 4 4 上での接続を定義する I N I C 6 2 2 により保持された C C B を S A N ・ C C B と称し、その N A S 記憶装置 6 4 2 とのネットワーク 6 4 4 上での接続を定義する I N I C 6 2 2 により保持された C C B を N A S ・ C C B と称する。追加のネットワークライン 6 5 0 、 6 5 2 を別の通信及び / またはストレージネットワークと接続することができる。

#### 【 0 0 5 7 】

クライアント 6 0 2 が、ネットワーク記憶装置 6 4 0 上にブロックで格納されたサーバ 6 0 0 上の 1 0 0 K B ファイルを読み出したいと仮定すると、前記クライアントは、ネットワーク 6 0 4 を通って S M B 読出しリクエストを送ることによって開始し、前記サーバ上のそのファイルの最初の 6 4 K B を要求することができる。この要求は、例えば 7 6 バイトだけであってもよく、前記サーバ上の I N I C 6 2 2 はメッセージのタイプ ( S M B ) 及び比較的小さいメッセージのサイズを認識し、前記 7 6 バイトを直接 A T C P フィルタレイヤ 6 2 8 に送り、それが前記リクエストをサーバの N e t B i o s 6 3 0 に送信する。 N e t B i o s 6 3 0 はセッションヘッダを S M B 6 3 2 に送り、これが前記読出しリ

10

20

30

40

50

クエストを処理して、要求されたデータがホストまたはI N I Cファイルキャッシュ上に保持されているかどうかを決定する。要求されたデータがファイルキャッシュにより保持されていない場合、S M Bはファイルシステムに読出しリクエストを出して、ネットワーク記憶装置6 4 0からのデータをI N I C 6 2 2ファイルキャッシュ内に書き込む。

【0 0 5 8】

これを実行するために、ファイルシステムはI N I C 6 2 2に、ネットワーク記憶装置6 4 0から6 4 K BのデータをI N I C 6 2 2ファイルキャッシュ内にフェッチするように指示する。次にI N I C 6 2 2は、前記データへのリクエストをネットワーク記憶装置6 4 0へネットワーク6 4 4で送信する。このリクエストは前記ブロックを読取る記憶装置6 4 0への1つまたは複数のS C S Iコマンドの形を取ることができ、前記コマンドにはI S C S Iまたは類似のプロトコルに従って、T C P / I Pヘッダが付加されている。記憶装置6 4 0上のコントローラがそのディスクドライブから要求されたブロックを読出し、I S C S Iまたは類似のプロトコルヘッダを前記ブロックまたは該ブロックのフレームサイズの部分に付け加え、かつその結果得られるフレームをネットワーク6 4 4上でI N I C 6 2 2に送ることによって、前記コマンドに応答する。前記フレームはI N I C 6 2 2により受け取られ、I N I C 6 2 2シーケンサにより処理され、ストレージC C Bと整合され、かつ要求された1 0 0 K Bファイルの一部分を形成するI N I Cファイルキャッシュの6 4 K Bファイルストリームとして再び組立てられる。前記ファイルストリームがI N I C 6 2 2ファイルキャッシュ上に一旦格納されると、S M Bは読出し応答を組立てかつそのファイルストリームを示すスキッタギャザリストをI N I C 6 2 2に送り、かつ前記応答をI N I C 6 2 2に送って、サーバC C Bに従ってネットワーク上で前記データを送る。I N I C 6 2 2は前記スキッタギャザリストを用いて、そのファイルキャッシュからデータパケットを読出し、これは、サーバC C Bに基づいてI N I Cにより作成されたI P / T C P / N e t B i o s / S M Bヘッダとプリペンドされ、かつその結果のフレームをネットワーク6 0 4上に送る。前記ファイルの残りの3 6 K Bは、同様の手法で送られる。このようにして、ネットワーク記憶装置上のファイルは、サーバの制御下で、I / Oパスまたはサーバプロトコルスタックに出会うファイルからのデータ無しで転送することができる。

【0 0 5 9】

クライアント6 0 2により要求されたデータがN A S記憶装置6 4 2上に格納される場合には、リクエストはサーバ6 0 0からその記憶装置6 4 2へ送ることができ、これはクライアント6 0 2にアドレス指定されたヘッダを有するデータを送ることによって応答し、サーバ6 0 0がルータとして機能する。サーバ6 0 0がプロキシサーバまたはウェブキャッシュサーバとして実行されるような実施例では、N A S記憶装置からのデータをサーバ6 0 0に送ることができ、これが前記データをそのファイルキャッシュに格納して、そのデータに関する将来のリクエストに対するより早い応答を提供する。この実施例では、サーバ6 0 0上のファイルシステムが、N A S記憶装置6 4 2上のファイルデータを要求するようにI N I Cを指示し、これがファイルデータの最初の6 4 K Bを含む多数の約1 . 5 K Bパケットを送ることによって応答する。前記ファイルデータを含むパケットはI N I C 6 2 2が受信し、I N I C受信シーケンサにより分類されかつN A S・C C Bと整合され、最初のパケットからのセッションレイヤヘッダがホストスタックにより処理され、これがファイルシステムからI N I C 6 2 2ファイルキャッシュのアドレスのスキッタギャザリストを入手して、前記パケットからのデータを格納する。前記スキッタギャザリストはホストスタックによりI N I C 6 2 2に送られ、かつN A S・C C Bと共に格納され、I N I C 6 2 2が、N A S・C C Bに対応するあらゆる蓄積パケット及びその後続くパケットからのデータを、スキッタギャザリストに従ってファイルストリームとしてI N I C 6 2 2ファイルキャッシュ内へD M Aし始める。次に、ホストファイルシステムがI N I C 6 2 2に指示して、クライアントC C Bに基づいてヘッダを作成し、かつ前記ファイルストリームから読み出されたデータの packets に、該データをクライアント6 0 2に送るためにヘッダをプリペンドする。前記ファイルの残りの3 6 K Bは同様の手法

10

20

30

40

50

により送られ、かつ別のファイルストリームとして I N I C 6 2 2 ファイルキャッシュにキャッシュすることができる。ファイルストリームが I N I C 6 2 2 ファイルキャッシュに保持されることによって、クライアント 6 0 6 のようなクライアントへの前記ファイルへの後から来るリクエストが、より早く処理することができる。

#### 【 0 0 6 0 】

クライアント 6 0 2 により要求されたファイルがキャッシュに存在せず、その代わりにファイルブロックとしてサーバ接続記憶装置 6 3 4 上に格納されていた場合には、サーバ 6 2 2 ファイルシステムがホスト S C S I ドライバを指示して、サーバ接続記憶装置 6 3 4 から 1 0 0 K B のデータをサーバ 6 0 0 ファイルキャッシュ内にフェッチさせる（前記ファイルシステムは I N I C 6 2 2 ファイルキャッシュ上のデータをキャッシュしたくないと仮定する）。次に、ホスト S C S I ドライバが S C S I チャネル 6 3 8 上で前記データへの S C S I リクエストをサーバ直接接続型記憶装置 6 3 4 に送る。サーバ直接接続型記憶装置 6 3 4 上のコントローラが、要求されたブロックをそのディスクドライブから読み出しかつ該ブロックを S C S I チャネル 6 3 9 上で前記 S C S I ドライバに送ることによって前記コマンドに回答し、これが前記ファイルシステムの命令下でキャッシュマネージャと対話して、前記ブロックをファイルストリームとしてサーバ 6 0 0 ファイルキャッシュに格納する。次にファイルシステムリダイレクタが S M B を指示して、前記ファイルストリームのスキッタギャザリストを I N I C 6 2 2 に送り、これを I N I C 6 2 2 が用いて、サーバ 6 0 0 ファイルストリームからデータパケットを読み出す。I N I C 6 2 2 は、サーバ C C B に基づいて作成されたヘッダで前記データパケットをプリペンドしかつその結果のフレームをネットワーク 6 0 4 上に送る。

#### 【 0 0 6 1 】

この応答が到着したときクライアント 6 0 2 上で I N I C 6 0 6 が動作し、I N I C 6 0 6 は受け取った第 1 フレームからこの接続が高速バス 6 2 0 プロセッシング（T C P / I P、N e t B i o s、C C B をマッチング）を受け取っていることを認識し、S M B 6 1 6 がこの第 1 フレームを用いてメッセージのためのバッファステージを獲得する。バッファの割り当ては、あらゆる N e t B i o s / S M B ヘッダを含む前記フレームの最初の 1 9 2 バイトを A T C P 高速バス 6 2 0 を介してクライアント N e t B i o s 6 1 4 に直接送り、適当なヘッダを N e t B i o s / S M B に与えることによって提供される。N e t B i o s / S M B はこれらのヘッダを分析し、リクエスト I D との整合によって、これが元の読み出し接続への応答であることを認識し、かつその中にデータが配置されるクライアントファイルキャッシュのバッファの 6 4 K リストを A T C P コマンドドライバに与える。この段階ではただ 1 つのフレームが到着しているが、このプロセッシングが行なわれている間により多くのフレームが到着し得る。前記クライアントバッファリストが A T C P コマンドドライバ 6 2 8 に与えられるとすぐに、これはその転送情報を I N I C 6 0 6 に送り、I N I C 6 0 6 は D M A によりそれらのバッファ内に蓄積されている全てのフレームデータの送信を開始する。

#### 【 0 0 6 2 】

クライアント 6 0 2 が S M B ファイルをサーバ 6 0 0 に書き込みたい場合には、書き込みリクエストが、I N I C 6 2 2 により保持されている C C B とマッチングするであろうネットワーク 6 0 4 上を送信される。ファイル書き込みの最初のパケットからのセッションレイヤヘッダがサーバ S M B 6 3 2 により処理されて、サーバ 6 0 0 または I N I C 6 2 2 ファイルキャッシュ内にバッファが割り当てられ、それらバッファに関するアドレスのスキッタギャザリストが、高速バスプロセッシングが適当であると仮定した場合、I N I C 6 2 2 へ戻される。S M B ファイルデータを含むパケットが I N I C 6 2 2 により受け取られ、I N I C 受信シーケンサにより分類されかつキューの中に配置される。I N I C 6 2 2 プロセッサは、パケットがサーバ C C B に対応していることを認識しかつ前記パケットからのデータを I N I C 6 2 2 またはサーバ 6 0 0 ファイルキャッシュバッファ内にスキッタギャザリストに従って D M A して、ファイルストリームを形成する。

#### 【 0 0 6 3 】

次に、前記ファイルシステムが、前記ファイルストリームのサーバ記憶装置 634、ネットワーク記憶装置 640 または NAS 記憶装置 642 への送信を編成する。ファイルストリームをサーバ記憶装置 634 に送るために、ファイルシステムはサーバ 600 内の SCSI ドライバに命令して、前記ファイルストリームをファイルブロックとして記憶装置 634 に送る。ファイルストリームをネットワーク記憶装置 640 に送るために、ファイルシステムは INIC に命令して、SAN・CCB に基づいて SCSI または類似のヘッダを作成し、それらのヘッダを前記スキッタギャザリストに従って前記ファイルストリームから読み出したパケットにプリPENDし、その結果のフレームをネットワーク 644 上で記憶装置 640 に送る。例えば分散ファイルキャッシュまたはプロキシサーバの実行に有用で有り得る NAS 記憶装置 642 に前記ファイルストリームを送るために、ファイルシステムリダイレクタは適当な Net Bios / SMB ヘッダをプリPENDし、かつ INIC に命令して、NAS・CCB に基づいて IP / TCP ヘッダを作成し、かつそれらのヘッダを、スキッタギャザリストに従って前記ファイルストリームから読み出されたパケットにプリPENDし、その結果のフレームをネットワーク 644 上で記憶装置 642 に送る。

#### 【0064】

図 15 は図 14 に示すものと類似のものを示しているが、図 15 は、特にリアルタイムトランスポートプロトコル (RTP) 及びリアルタイムトランスポートコントロールプロトコル (RTCP) が関係し得る音声及び映像通信のコンテキストにおいて、本発明のユーザデータグラムプロトコル (UDP) に焦点を合わせている。サーバ 600 は、図 14 に示すプロトコルレイヤに加えて、UDP レイヤ 654、AUDP レイヤ 655、RTP / RTCP レイヤ 656 及びアプリケーションレイヤ 657 を有する。同様に、クライアント 602 は、図 14 に示すプロトコルレイヤに加えて UDP レイヤ 660、AUDP レイヤ 661、RTP / RTCP レイヤ 662 及びアプリケーションレイヤ 663 を有する。図 15 に RTP / RTCP レイヤ 656、662 は図示されているが、セッション初期設定プロトコル (SIP) または媒体ゲートウェイコントロールプロトコル (MGCP) のような他のセッション及びアプリケーションプロトコルを UDP の上に用いることができる。

#### 【0065】

オーディオ / ビデオインタフェース (AVI) 666 または類似の周辺装置がサーバ 600 の I / O バス 639 に接続され、かつ AVI 666 にはスピーカ 668、マイク 670、ディスプレイ 672 及びビデオカメラで有り得るカメラ 674 が接続されている。簡単にするために 1 つのインタフェースとして図示されているが、これに変えて AVI 666 は、サウンドカードまたはビデオカードのようなそれぞれ I / O バス 639 に接続される別個のインタフェースを有することができる。別の実施例では、AVI インタフェースが I / O バスではなくホストバスに接続され、かつホスト CPU のメモリコントローラまたは I / O コントローラと一体にすることができる。AVI 666 は、音声または映像データのための一時的記憶装置として機能するメモリを有することができる。また、クライアント 602 は I / O バス 675 を有し、AVI 677 は I / O バス 675 に接続されている。AVI 677 に接続されているのは、スピーカ 678、マイク 680、ディスプレイ 682、及びビデオカメラで有り得るカメラ 684 である。簡単にするために 1 つのインタフェースとして図示されているが、これに代えて AVI 666 は、サウンドカード及びビデオカードのようなそれぞれ I / O バス 675 に接続された別個のインタフェースを有することができる。

#### 【0066】

TCP と異なり、UDP は従属性の接続を提供しない。その代わりに、UDP パケットはベストエフォートベースで送られ、かつ失われたまたは壊れたパケットは UDP より上位のレイヤがそのようなサービスを提供しない限り、再送されない。UDP は、接続を確立する必要なく IP を介してアプリケーションがデータを送る方法を提供する。しかしながら、UDP またはネットワークファイルシステム (NFF)、TCP、RTCP、SIP

10

20

30

40

50



またはM G C Pのような別のプロトコルによるリクエストに回答してソケットが指定され、そのソケットがU D Pにより送られたメッセージを受け入れることができる受信装置上のポートを割り当てる。ソケットは、ソース及びデスティネーションI Pアドレス並びにソース及びデスティネーションU D Pポートを示すU D Pにより用いられるアプリケーションプログラミングインタフェース( A P I )である。

【 0 0 6 7 】

U D Pを介して従来どおりに通信を送るために、セッションレイヤまたはアプリケーションレイヤのヘッダを含むことができ、かつ通常約8 K BまでのN F Sデフォルトサイズのデータグラムである6 4 K Bまでのデータグラムを、ホストコンピュータによりU D PヘッダでプリペンドしてI Pレイヤに渡す。U D Pヘッダは、ソース及びデスティネーションポート並びに任意のチェックサムを含む。イーサネット伝送の場合、U D Pダイアグラムは、必要に応じて、I Pレイヤにより約1 . 5 K Bのフラグメントに分割され、I Pヘッダでプリペンドされかつイーサネットヘッダ及びトレーラでのカプセル化のためにM A Cレイヤに与えられ、かつネットワーク上に送られる。I Pヘッダは、U D Pデータグラムに固有のI P識別フィールド( I P ・ I D )と、そのダイアグラムのより多くのフラグメントが続いているかどうかを示すフラグフィールドと、前記フラグメントを正しい順序で再アセンブリするために、前記ダイアグラムのどの1 . 5 K BフラグメントがI Pヘッダに接続されたかを示すフラグメントオフセットフィールドを有する。

【 0 0 6 8 】

クライアント6 0 2からサーバ6 0 0への高速パスU D Pデータ転送の場合、例えばクライアントアプリケーション6 6 3またはオーディオ/ビデオインタフェース6 7 7からのデータのファイルストリームは、ファイルシステムの命令下でI N I C 6 0 6のメモリ上に格納することができる。アプリケーション6 6 3は、例えば約8 K Bを有するように前記ファイルストリームを構成することができ、それはアプリケーション6 6 3の命令下でI N I C 6 0 6により獲得される上位レイヤヘッダを含むことができ、I N I C 6 0 6により受け取られる各ファイルストリームは、指定されたソケットに従ってI N I C 6 0 6によりU D Pヘッダでプリペンドされ、U D Pデータグラムを作成する。U D Pデータグラムは、I N I C 6 0 6により6つの1 . 5 K Bメッセージフラグメントに分割され、それはそれぞれI P及びM A Cレイヤヘッダでプリペンドされて、ネットワーク6 0 4上に伝送されるI Pパケットを作成する。

【 0 0 6 9 】

I N I C 6 2 2は、I N I C 6 0 6から送られたイーサネットフレームをネットワーク6 0 4から受け取り、ヘッダをチェックサムしかつ処理して、関連するプロトコルを決定し、かつU D P及び上位レイヤヘッダをA U D Pレイヤ6 5 5に送って、そのU D Pデータグラムのパケットから前記データのデスティネーションアドレスのリストを入手する。U D P及び上位レイヤヘッダはU D Pデータグラムからの最大6つまでのパケットの1つに含まれ、かつそのパケットは通常そのデータグラムからの他のパケットより先に受け取られる。前記パケットが間違った順序で到着する場合には、I N I C 6 2 2が再アセンブリバッファ内でU D Pデータグラムからの前記最大6つのパケットを列に並ばせる。U D Pデータグラムに対応するパケットは、そのI P ・ I Dに基づいて識別され、フラグメントオフセットに基づいて連結される。U D P及び上位レイヤヘッダを含むパケットが処理されてデスティネーションアドレスが獲得された後、並べられた前記データをそれらのアドレスに書き込むことができる。U D Pデータグラムからの全てのパケットが来ない可能性を説明するために、I N I C 6 2 2は、受け取ったデータの落とし込み( dropping )をトリガするタイマを用いることができる。

【 0 0 7 0 】

リアルタイム音声または映像通信の場合、電気通信接続がサーバ6 0 0とクライアント6 0 2との間で最初に設定される。国際電気通信連合( I T U ) H . 3 2 3規格による通信では、電気通信接続の設定が、R T Pを用いたデータの流れるためにソース及びデスティネーションU D Pポートを指定するT C Pダイアログで実行される。電気通信接続設定の

10

20

30

40

50

ための別のTCPダイアログが、RTPを用いたデータフローをモニタするためにソース及びデスティネーションUDPポートを指定する。SIPによって電気通信接続を初期設定するための別のメカニズムが提供される。UDPソケットが指定された後、本発明によるサーバ600からクライアント602への音声または映像データの伝送を開始することができる。

#### 【0071】

例えば、オーディオ/ビデオインタフェース666は、マイク670及びカメラ674からの音声及びイメージを、INIC622に利用可能なオーディオ/ビデオデータに変換する。サーバ600ファイルシステムの指示のもとで、INIC622は、RTPヘッダを含むAVデータの8KBファイルストリームを獲得し、かつそのデータをINICファイルキャッシュに格納することができる。電気通信接続に従って、INIC622はUDPヘッダを各ファイルストリームにプリペンドすることができ、次にこれが1.5KBのフラグメントに断片化され、そのそれぞれがIP及びイーサネットヘッダでプリペンドされて、ネットワーク604上に伝送される。これに変えて、INIC622はINICメモリ上に格納されたAVデータから1.5KBのUDPデータグラムを作成することができ、断片化を回避し、かつUDP及びIP及びMACレイヤヘッダを同時に前記電気通信接続に対応するINIC622上に保持されたテンプレートから作成できるようにする。また、AVデータで8KBより大きいUDPデータグラムを作成することが可能であり、それが追加の分割を作成するが、データグラムごとにより大きなブロックのデータを転送することができる。RTPヘッダは、AVデータがオーディオ/ビデオインタフェース666によりパケット化される相対時間を表示するタイムスタンプを有する。

#### 【0072】

ホストCPUによるAVデータの従来のプロトコルプロセッシングと対照的にINIC622は、UDPヘッダ作成、IP断片化及びIPヘッダ作成のタスクをオフロードすることができ、ホストCPUの複数のプロセッシングサイクル、複数の割り込みを不要にし、かつホストバストラフィックを大幅に少なくすることができる。INIC622は、複数のヘッダを1度にプリペンドすることによって、より効率的にヘッダ作成タスクを実行することができる。更に、AVデータは、最初にI/Oバス639上をプロトコルプロセッシングスタックへ送られ、かつ次にRTP/UDP/IPパケットとしてI/Oバス639上を送り返されるよりはむしろ、I/Oバス639上でINIC622により直接アクセスすることができる。またこれらの効率上の利点は、従来のプロトコルプロセッシングと比較して、AVデータの送信における遅れを大幅に小さくすることができる。

#### 【0073】

リアルタイム音声及び映像通信システムの対象は、通信における遅れやジッタが該システムを介して通信する人々に実質的に知覚できないことである。ジッタは、パケットを受信する遅れの変化によって生じ、これが、音や景色が記録されたテンポと比較してそれらが表示されるテンポに変動を生じさせることになる。上述したパケット化された音声及び/または映像を伝送する際の遅れを減少させる利点は、受信装置でのバッファリングを介してジッタもまた減少させるのに用いることができるが、これは遅れの減少によってジッタを円滑化するための受信装置における時間が増加するからである。AVデータを通信する際の遅れ及びジッタを更に減らすために、INIC622によりUDPデータグラムフラグメントにプリペンドされたIPまたはMACレイヤヘッダは、そのサービス分野での高品質のサービス(QoS)表示を含んでいる。この表示は、ネットワーク604により用いられて、例えば優先的にバンド幅を割り当て、かつそれらのフレームについてキューに入れることによって、高QoSフレームの伝送を早めることができる。

#### 【0074】

AVデータを含むフレームがINIC606により受信されると、QoSの表示がINIC606受信論理によってフラグを立てられ、かつパケットがINIC606の高い優先順位の受信キューにバッファされる。INIC606のハードウェア論理によってパケットのヘッダを分類し、かつ検証することによって、従来のIP/UDPヘッダのプロセッシ

10

20

30

40

50

ングに比較して、A Vデータを受け取る際の遅れが減少する。各データグラムはUDPヘッダはUDPレイヤに送ることができ、これが前記ヘッダを処理して、関連するソケットをアプリケーションに表示する。次に、前記アプリケーションはINIC 606に指示して、そのUDPヘッダに対応するデータをオーディオ/ビデオインタフェース677のデスティネーションに送る。これに代えて、UDPデータグラムに対応する全てのフラグメントが前記受信キューに連結された後、前記データグラムをUDPヘッダ無しでファイルストリームとして前記ソケットに対応するオーディオ/ビデオインタフェース677のデスティネーションに書き込まれる。任意によって、UDPデータグラムは最初にINICメモリの別の部分に格納することができる。A Vデータ及びRTPヘッダは、次にオーディオ/ビデオインタフェース677に送られ、そこで復号されてスピーカ678及びディスプレイ682で再生される。

10

#### 【0075】

A Vデータフロー全体及び受信したIP/UDPパケットの大部分のプロトコルプロセッシングはINIC 606により取扱うことができ、クライアント602のCPUの複数のプロフェッシングサイクルを省略してホストバストラフィック及び割り込みを大幅に減少することに注意する。更に、INIC 606の専用ハードウェア論理によって、汎用CPUよりもより効率的にヘッダを分類することができる。また、A Vデータは、INIC 602により、最初にI/Oバス675上でプロトコルプロセッシングスタックに送られかつ次にI/Oバス675上でオーディオ/ビデオインタフェース677に送られるのではなく、I/Oバス675上をオーディオ/ビデオインタフェース677に直接提供することができる。

20

#### 【0076】

図16は、ネットワークインタフェース、ストレージコントローラ及びプロトコルプロセッサの機能を1つのASICチップ703に結合させたINIC 22の図を示している。本実施例におけるINIC 22は、サーバ及びネットワークストレージアプリケーションのための高速プロトコルプロセッシング用に設計された全2重4チャンネル10/100メガビット毎秒(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が提供される。また、INIC 22は、パーソナルコンピュータ、ワークステーション、ルータまたはTCP/IP、PPCP/IPまたはSPX/IPXプロトコルが使用されている他のホストに接続することができる。このようなINICの説明は、本願明細書の冒頭部分に挙げた関連特許出願に記載されている。

30

#### 【0077】

INIC 22はネットワークコネクタによって4つのネットワークライン702、704、706、708に接続され、これらはツイストペア同軸ケーブルまたは光ファイバのような、それぞれの接続によって、米国カリフォルニア州94538、フリーモント、レイサイドパークウェイ47200のSEEQ Technology Incorporatedから市販されているモデル80220/80221イーサネットメディアインタフェースアダプタのような物理的レイヤチェック712、714、716、718を介して媒体独立インタフェース(MII)を提供する多数の異なるコンジットに沿ってデータを運ぶことができる。前記ラインは好ましくは802.3コンプライアントでありかつINICに関連して4つの完全なイーサネットノード、INIC支援の10ベースT、10ベースT2、100ベースTX、100ベースFX及び100ベースT4に加えて将来のインタフェース標準を構成する。物理的レイヤの識別及び初期設定が、ホストドライバの初期設定ルーチンを介して実行される。ネットワークライン702、704、706、708とINIC 22との間の接続は、MACサブレイヤの基本機能を実行し、基本的にINICがネットワークライン702、704、706、708にアクセスする時に制御する論理回路を含むMAC装置MAC-A 722、MAC-B 724、MAC-C 726、MAC-D 728によって制御される。MAC装置722、724、726、728は、無差別モード、マルチキャストモードまたはユニキャストモードで動作することができ、INICがネットワークモニタとして機能し、ブロードキャストパケット及びマルチキャスト

40

50

トパケットを受信し、かつ各ノードについて複数のMACアドレスを実行できるようにする。また、MAC装置722、724、726、728によって、単一のネットワーク管理プロトコル(SNMP)について使用可能な統計的情報が得られる。

【0078】

MAC装置722、724、726、728はそれぞれ、送受信シーケンサXMT&RCV-A732、XMT&RCV-B734、XMT&RCV-C736、XMT&RCV-D738に接続される。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時、幾つかのプロトコルプロセッシングのステップをオンザフライで実行することができる。MAC装置と組み合わせて、送受信シーケンサ732、734、736、738は、ハードウェアにおけるデータリンク、ネットワーク、トランスポート、セッション及び適当な場合にはプレゼンテーション及びアプリケーションレイヤのプロトコルに関するパケットステータスを集めることができ、従来のシーケンシャル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プロセッシングのための時間を大幅に減少させることができる。送受信シーケンサ732、734、736、738はSRAM及びDMAコントローラ740に接続され、それらはDMAコントローラ742及びSRAM744を有し、これがスタティックランダムアクセスメモリ(SRAM)バッファ748を制御する。SRAM及びDMAコントローラ740は外部メモリコントロール750と対話して、外部メモリバス752を介して、ICチップ700の近くに配置されたダイナミックランダムアクセスメモリ(DRAM)バッファ755との間でフレームを送りかつ受信する。DRAMバッファ755は、4MB、8MB、16MBまたは32MBとして構成することができ、かつ任意によりチップ上に配置することができる。SRAM及びDMAコントローラ740は、本実施例ではINIC22とPCIインタフェースバス757との間のインタフェースを管理するPCIバスインタフェース装置(BIU)756であるI/Oブリッジに接続されている。64ビットの多重化BIU756によって、スレイブ及びマスタ機能双方についてPCIバス757との直接のインタフェースが得られる。INIC22は64ビットまたは32ビットのPCI環境で動作することができ、そのいず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する。

【0079】

マイクロプロセッサ780がSRAM及びDMAコントローラ740とPCI・BIU756に接続されているマイクロプロセッサ780の命令及びレジスタファイルはSRAMの書き込み可能なオンチップコントロールストア(WCS)及びリードオンリメモリ(ROM)を有するオンチップコントロールストア780に存在する。マイクロプロセッサ780は、入力するフレームを処理し、ホストコマンドを処理し、ネットワークトラフィックを指示しかつPCIバストラフィックを指示することができるプログラマブルステートマシンを提供する。3つのプロセッサが、クロックサイクルごとに1つの命令を発しかつ完了する3段階パイプラインアーキテクチャの共有ハードウェアを用いて実行される。受信プロセッサ782は、基本的に通信を受け取るために用いられるのに対し、送信プロセッサ784は、基本的に全2重通信を容易にする目的で通信を送るために用いられ、ユーティリティプロセッサ786は、PCIレジスタのアクセスを監督しかつ制御することを含む様々な機能を提供する。

【0080】

プロセッサ782、784、786のための命令がオンチップコントロールストア781に存在するので、前記3つのプロセッサの機能は容易に再定義することができ、それによってマイクロプロセッサ780は所定の環境に適應することができる。例えば、受信機能に必要なプロセッシングの量が送信またはユーティリティ機能に必要なそれを上回ることがある。そのような場合、幾つかの受信機能が送信プロセッサ784及び/またはユーティリティプロセッサ786によって実行することができる。それに代えて、追加レベルのパイプライン方式を作って3つではなく4つまたはそれ以上の仮想プロセッサを設けることができ、その追加レベルに受信機能が指示される。

【0081】

10

20

30

40

50

本実施例の INIC 22 は、DRAM 755 のテーブルに保持される最大 256 の CCB をサポートすることができる。また、しかしながら、順次検索を省略するために SRAM 748 には順番がハッシュ状態の CCB インデックスが存在する。一旦ハッシュが生成されると、CCB が SRAM にキャッシュされ、本実施例では 60 までの CCB が SRAM にキャッシュされる。SRAM にキャッシュされた 60 の CCB の割り当ては、以下に説明するように、最小使用頻度のレジスタによって操作される。これらキャッシュの配置は、送信プロセッサ 784 と受信プロセッサ 786 との間で共有され、それによってより大きなロードを有す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、前記シーケンサ間で共有されるべき 8 つのヘッダバッファ及び 8 つのコマンドバッファが存在する。所定のヘッダまたはコマンドバッファは、特定の CCB バッファに統計的にリンクされないが、そのリンクは各フレームごとにダイナミックだからである。

10

#### 【0082】

図 17 はパイプラインマイクロプロセッサ 780 の概観を示しており、受信、送信及びユーティリティプロセッサへの命令がクロックインクリメント I、II 及び III によって交互に変わるそれぞれにパイプラインの段階に対応する 3 つのフェーズで実行される。各フェーズは異なる機能について責任を有し、前記 3 つのプロセッサのそれぞれが各クロックインクリメントの際に異なるフェーズを占有する。各プロセッサは通常コントロールストア 781 からの異なる命令ストリームに基づいて動作し、それぞれに前記フェーズのそれぞれ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。

#### 【0083】

20

一般にパイプラインマイクロプロセッサの第 1 インストラクションフェーズ 800 は、命令を完了しかつ結果をデスティネーションオペランドに格納し、次の命令をフェッチし、かつその次の命令をインストラクションレジスタに格納する。第 1 のレジスタセット 790 によって、命令レジスタを含む多数のレジスタが提供され、かつ前記第 1 レジスタセットのためのコントロール 792 セットによって、第 1 レジスタセット 790 への記憶のためのコントロールが提供される。幾つかのアイテムがコントロール 790 による変更なしに第 1 フェーズを通過し、かつその代わりに単に第 1 レジスタセット 790 または RAM ファイルレジスタ 833 内にコピーされる。第 2 インストラクションフェーズ 860 はインストラクションデコード及びオペランドマルチプレクサ 798 を有し、これは一般に第 1 レジスタセット 490 のインストラクションレジスタ内に格納された命令をデコードし、かつ生成されている全てのオペランドを集め、これは次に第 2 のレジスタセット 796 のデコードレジスタに格納される。第 3 インストラクションフェーズ 900 で使用される第 1 レジスタセット 790、第 2 レジスタセット 796 及び第 3 レジスタセット 801 は、図 18A ~ C により詳細に示すように、多数の同じレジスタを有する。インストラクションデコード及びオペランドマルチプレクサ 798 は、第 1 フェーズ 806 及び第 2 フェーズ 860 の双方で動作する RAM ファイルレジスタ 833 の 2 つのアドレスおよびデータポートから読み出すことができる。プロセッサ 780 の第 3 フェーズ 900 は、演算論理装置 (ALU) 902 を有し、これは一般に第 2 レジスタセットからのオペランドについて全ての ALU 動作を実行し、結果を第 3 レジスタセット 801 に含まれるリザルトレジスタに格納する。スタックエクスチェンジ 808 がレジスタスタックを順序付けすることができ、かつキューマネージャ 803 がプロセッサ 780 のキューを編成し、その結果が第 3 レジスタセットに格納される。

30

40

#### 【0084】

インストラクションは、循環パイプライン 805 に示されるように第 1 フェーズに続き、第 3 フェーズが続く。全ての所定のフェーズにおける組み合わせ遅れを最小にするために、様々な機能が命令実行の 3 つのフェーズに渡って分散されていることに注意する。本実施例では 66 MHz の周波数で、各クロックインクリメントは 15 ナノ秒で完了し、3 つの前記プロセッサのそれぞれについて 1 つの命令を完了するのに合計 45 ナノ秒を要する。インストラクションフェーズが循環する様子が図 18A ~ C により詳細に示されており、各フェーズが異なる図で示されている。

50

## 【 0 0 8 5 】

より詳細に言えば、図 1 8 A は、第 1 フェーズ 8 0 0 の幾つかの特定のハードウェア機能を示しており、一般に第 1 レジスタセット 7 9 0 及び関連するコントロール 7 9 2 を有する。第 1 レジスタセット 7 9 2 のコントロールには、アドレス及び書込みデータを S R A M アドレス及びデータレジスタ 8 2 0 内にロードするための論理コントロールである、S R A M コントロール 8 2 0 を有する。従って、第 3 フェーズ 9 0 0 からの A L U 9 0 2 の出力は、S R A M コントロール 8 0 2 によって S R A M アドレス及びデータレジスタ 8 2 0 のアドレスレジスタまたはデータレジスタ内に配置される。同様にロードコントロール 8 0 4 がコンテキストレジスタ 8 2 2 にファイルするファイルのためのコンテキストを書き込むためのコントロールを提供し、かつ別のロードコントロール 8 0 6 が、様々な雑データをフリップフロップレジスタ 8 2 5 に格納するためのコントロールを提供する。キャリアビットが設定されているかどうかというような A L U 状態コードが、第 1 フェーズ 8 0 0 において実行される動作無しで、A L U 状態コードレジスタ 8 2 8 内にクロックされる。フラグレコード 8 0 8 がロックを設定するような、フラグレジスタ 8 3 0 内に格納される様々な機能を実行する。

10

## 【 0 0 8 6 】

R A M ファイルレジスタ 8 3 3 は、アドレス及びデータのための 1 つの書込みポートとアドレス及びデータのための 2 つの読み出しポートとを有し、それによって 2 以上のレジスタが同時に読み出すことができる。上述したように、R A M ファイルレジスタ 3 3 は基本的に第 1 及び第 2 フェーズにまたがるが、それは第 1 フェーズ 8 0 0 において書き込まれかつ第 2 フェーズ 8 6 0 において読み出されるからである。コントロールストアインストラクション 8 1 0 によって、この図面には記載されていないコントロールストア 7 8 1 からの新しいデータによる前記プロセッサの再プログラミングが可能になり、前記インストラクションはインストラクションレジスタ 8 3 0 のみ格納される。このためのアドレスがフェッチコントロールレジスタ 8 1 1 において生成され、それによりどのアドレスにフェッチするかが決定され、アドレスがフェッチアドレスレジスタ 8 3 8 に格納される。ロードコントロール 8 3 1 がプログラムカウンタ 8 4 0 の命令を出し、これはコントロールストアのためのフェッチアドレスに良く似たように動作する。3 つのレジスタからなる後入れ先出し方式のスタック 8 4 4 が、このフェーズで他の動作を受けることなく、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 8 4 8 のためのロードコントロール 8 1 7 が任意により含まれ、それによって起こり得るエラーの修正が可能になる。

20

30

## 【 0 0 8 7 】

図 1 8 B は、R A M ファイルレジスタ 8 3 3 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 8 6 0 を示している。スクラッチ S R A M 8 6 5 が、最初の 2 つのフェーズを通過して第 3 で増分されるレジスタを有する、前記第 1 レジスタセットの S R A M アドレス及びデータレジスタ 8 2 0 から書き込まれる。スクラッチ S R A M 8 6 5 は、スタック 8 4 4、デバッグアドレス 8 4 8、及び上述した S R A M アドレス及びデータレジスタを除いて、前記第 1 レジスタセットからの大部分のレジスタがそうであるように、インストラクションデコード及びオペランドマルチプレクサ 7 9 8 によって読み出される。インストラクションレコード及びオペランドマルチプレクサ 7 9 8 は、レジスタセット 7 9 0 の各レジスタ及び S R A M 8 6 5 を見て、インストラクションをデコードし、かつ次のフェーズにおける動作のためのオペランドを集め、特に以下の A L U 9 0 2 に提供するオペランドを決定する。インストラクションデコード及びオペランドマルチプレクサ 7 9 8 の成果は、A L U オペランド 8 7 9、8 8 2、A L U 状態コードレジスタ 8 8 0、及び本実施例では 3 2 のキューを制御することができるキューチャネル及びコマンド 8 8 7 レジスタを含む第 2 レジスタセット 7 9 6 の多数のレジスタに格納される。レジスタセット 7 9 6 のレジスタの幾つかは、デコード 7 9 8 により実質的にデコードすることなく、上述したインストラクションレジスタ 8 3 5 からほぼ直接にロードされ、プログラムコントロール 8 9 0、リテラルフィールド 8 8 9、テストセレクト 8 8 4 及びフラグ 8 8 5 が含まれる。第 1 フェーズ 8 0 0 のファイルコンテキスト 8 2 2 のような他のレ

40

50

ジスタは、常に第2フェーズ860のファイルコンテキスト877に格納されるが、マルチプレクサ870により集められたオペランドとして取扱うこともできる。スタックレジスタ844は、単にスタックレジスタ894にコピーされる。プログラムカウンタ840はこのフェーズで増分され868、かつレジスタ892に格納される。また、任意のデバッグアドレス848が増分され、870、かつロードコントロール875が各フェーズでのエラーを制御できるようにする目的で、この時点でパイプライン805から供給することができ、その結果がデバッグアドレス898に格納される。

#### 【0088】

図18Cは、ALU及びキュー動作を含む第3のマイクロプロセッサフェーズ900を示している。ALU902は加算器、優先順位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ALU出力918、ALU状態コード920及びデステーションオペランドリザルト922に格納される。ファイルコンテキストレジスタ925、フラグセレクトレジスタ926及びリテラルフィールドレジスタ930は、先のフェーズ860から単にコピーされるだけである。テストマルチプレクサが設けられて、条件付ジャンプが結果としてジャンプになるかどうか決定され、その結果がテストリザルトレジスタ924に格納される。テストマルチプレクサ904は、そうではなくフェッチコントロール811のような類似の決定と共に第1フェーズ800で実行することができる。スタックエクスチェンジ808が、プログラムカウンタをスタック879からフェッチし、またはそのスタック上にプログラムカウンタを置くことによってスタックを上下にシフトさせ、その結果がプログラムコントロール934、プログラムカウンタ938及びスタック940レジスタに格納される。SRAMアドレスは、任意によりこのフェーズ900で増分させることができる。このフェーズでもエラーを制御できるようにするために、別のデバッグアドレス942の別のロードコントロール910をこの時点でパイプライン805から強制することができる。この図面では一体に表されているQRAM&QARU906が、キューチャンネル及びコマンドレジスタ887から読出し、SRAMに格納し、かつキューを再編成し、データのキューを管理するのに必要なようにデータ及びポインタを追加または取り外し、結果をテストマルチプレクサ904及びキューフラグ及びキューアドレスレジスタ928に送る。このようにQRAM&QARU906は、従来はCPU上のソフトウェアにより順に実行されていたタスク、前記3つのプロセッサのためのキューを管理する義務を引き受け、代わりにキューマネージャ906によって加速されかつ実質的にパラレルなハードウェアキューイングが行なわれる。

#### 【0089】

図19は、キューマネージャ906によって管理される32のハードウェアキューの内の2つを示しており、前記キューのそれぞれはSRAMヘッド、SRAM末尾及びTRAM本体内で情報を列に並べる能力をも有し、各キューの拡張及び個々の構成を可能にしている。このように、FIFO1000はSRAM記憶装置1005、1007、1009、1011を有し、それぞれが8バイトで合計32バイトを有するが、これら装置の数及び能力は、他の実施例において変えることができる。同様に、FIFO1002はSRAM記憶装置1013、1015、1017、1019を有する。SRAM装置1005、1007はFIFO1000のヘッドであり、かつ装置1009、1011はそのFIFOの末尾であるのに対し、装置113、115はFIFO1002のヘッドでありかつ装置117、119はそのFIFOの末尾である。FIFO1000の情報は矢印1022で示すように、ヘッドの装置1005または1007内に書き込むことができ、かつ矢印1025で示すように、末尾の装置1011または1009から読み出される。しかしながら、特定の項目はヘッドの装置1005または1007への書込み及び読出し両方を行なうことができ、または末尾の装置1009または1011への書込み及び読出し双方を行なうことができ、データの移動及び待ち時間を最小にする。同様に、FIFO1002の情報は一般に、矢印1033で示すように、ヘッドの装置1013または1015に書き込まれ、かつ矢印1039で示すように、末尾の装置1017または1019から読み出されるが、そうではなく、書き込まれた同じヘッドまたは末尾の装置から読み出すことが

できる。

#### 【 0 0 9 0 】

S R A M ・ F I F O 1 0 0 0、1 0 0 2 は双方共 D R A M 7 5 5 に接続されており、それら F I F O が実質的に無限に拡張して S R A M ヘッド及び末尾部分が一杯であるような状況进行处理できるようにしている。例えば、Q 0 とラベル付けされた 3 2 のキューの最初のものは、矢印 1 0 2 7 で示すように、F I F O 7 0 0 のヘッド部分または末尾部分で待ち行列に入れられる代わりに、キューマネージャの指示下で動作する D M A 装置によって、D R A M 7 5 5 のエントリを待ち行列に入れることができる。D R A M 7 5 5 に格納されたエントリは、矢印 1 0 3 0 で示すように、S R A M 装置 1 0 0 9 に戻り、その F I F O の長さ及びフォールスルー時間を延長させる。S R A M から D R A M への転換は、D R A M の方がより遅かつ D M A の動作が別の待ち時間を生じさせることから、一般に S R A M が一杯である場合のために取っておかれる。このように Q 0 は、キューマネージャ 8 0 3 によって F I F O 1 0 0 0 及び D R A M 7 5 5 の双方に格納されたエントリで構成することができる。同様に、例えば Q 2 7 に対応する F I F O 1 0 0 0 に行きの情報は、矢印 1 0 3 5 で示すように、D M A によって D R A M 7 5 5 内に動かすことができる。遅いとはいえコスト効果のよい D R A M 8 0 3 で待ち行列に入れる能力は、初期設定の際にユーザ定義可能なものであり、前記キューを所望のサイズに変更できるようにしている。D R A M 7 5 5 にキューされた情報は、矢印 1 0 3 7 で示すように、S R A M 装置 1 0 1 7 に戻される。

#### 【 0 0 9 1 】

3 2 個のハードウェアキューのそれぞれのステータスが、図 2 0 に示すように、4 つの 3 2 ビットのレジスタからなる組み 1 0 4 0 に保持されかつそれからアクセスされ、各レジスタの特定のビットが特定のキューに対応していると好都合である。前記レジスタには、Q - O u t \_ \_ R e a d y 1 0 4 5、Q - I n \_ \_ R e a d y 1 0 5 0、Q - E m p t y 1 0 5 5 及び Q - F u l l 1 0 6 0 のラベルが付されている。特定のビットが Q - O u t \_ \_ R e a d y レジスタ 1 0 5 0 で設定されると、そのビットに対応するキューは、読み出される準備のできた情報を有するのに対し、同じビットの Q - I n \_ \_ R e a d y 1 0 5 2 レジスタにおける設定は、前記キューが書き込まれる準備ができていることを意味している。同様に、特定のビットの Q - E m p t y レジスタ 1 0 5 5 の肯定的設定がそのビットに対応するキューが空であることを意味するのに対し、特定のビットの Q - F u l l レジスタ 1 0 6 0 における肯定的設定は、そのビットに対応するキューが一杯であることを意味する。このように Q - O u t \_ \_ R e a d y 1 0 4 5 は、ビット 2 7 ・ 1 0 5 2、2 8 ・ 1 0 5 4、2 9 ・ 1 0 5 6 及び 3 0 ・ 1 0 5 8 を含むビット 0 ・ 1 0 4 6 から 3 1 ・ 1 0 4 8 を有する。Q - I n \_ \_ R e a d y 1 0 5 0 は、ビット 2 7 ・ 1 0 6 6、2 8 ・ 1 0 6 8、2 9 ・ 1 0 7 0 及び 3 0 ・ 1 0 7 2 を含めて、ビット 0 ・ 1 0 6 2 ~ 3 1 ・ 1 0 6 4 を有する。Q - E m p t y 1 0 5 5 は、ビット 2 7 ・ 1 0 7 8、2 8 ・ 1 0 8 0、2 9 ・ 1 0 8 2 及び 3 0 ・ 1 0 8 4 を含めて、ビット 0 ・ 1 0 7 4 ~ 3 1 ・ 1 0 7 6 を有し、Q - F u l l 1 0 7 0 は、ビット 2 7 ・ 1 0 9 0、2 8 ・ 1 0 9 2、2 9 ・ 1 0 9 4 及び 3 0 ・ 1 0 9 6 を含めて、ビット 0 ・ 1 0 8 6 ~ 3 1 ・ 1 0 8 8 を有する。

#### 【 0 0 9 2 】

F I F O 1 0 0 0 に対応する Q 0 はフリーバッファキューであり、全ての使用可能なバッファのアドレスのリストを保持している。このキューは、前記マイクロプロセッサまたは他のデバイスがフリーバッファのアドレスを必要とする時にアドレス指定され、従って一般に分かる程度の D R A M 7 5 5 を有する。従って、フリーバッファのアドレスを必要とするデバイスは、Q - 0 を調べてそのアドレスを得ることになる。F I F O 1 0 0 2 に対応する Q - 2 7 は、受信バッファ記述子のキューである。受け取ったフレームを受信シーケンサにより処理した後、シーケンサは Q - 2 7 に前記フレームの記述子を格納すると思われる。そのような記述子のロケーションが S R A M においてすぐに使用可能な場合には、Q - I n \_ \_ R e a d y 1 0 5 0 のビット 2 7 ・ 1 0 6 6 が設定されることになる。そうでない場合には、前記シーケンサはキューマネージャが S R A M から D R A M へ D M A の



移動を開始させるのを待たなければならず、それにより空間を解放して受信記述子を格納する。

【 0 0 9 3 】

S R A Mと前記プロセッサ間、送信シーケンサと受信シーケンサ間、及びS R A MとD R A M間におけるキューエントリの動きを管理するキューマネージャの動作が、図 2 1 により詳細に示されている。前記キューを利用するリクエストは、プロセッサリクエスト 1 1 0 2、送信シーケンサリクエスト 1 1 0 4、及び受信シーケンサリクエスト 1 1 0 6 を含む。前記キューへの他のリクエストは、D R A M対 S R A Mリクエスト 1 0 8 6 及び S R A M対 D R A Mリクエスト 1 1 1 0 であり、これは D R A Mと前記キューの S R A Mヘッド部分または末尾部分との間でデータを行ったり来たり動かす際にキューマネージャのために動作する。これら様々なリクエストのどれが次のサイクルでキューマネージャを使用することになるかを決定することは、優先順位論理アービタ 1 1 1 5 によって処理される。高周波数動作を可能にするために、キューマネージャはパイプライン化され、レジスタ A 1 1 1 8 及びレジスタ B 1 1 2 0 によって一時的ストレージが提供されるのに対し、ステータスレジスタ 1 1 1 2 は、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、送受信シーケンサが要求する D M A について偶数サイクルを、プロセッサのリクエストのために奇数サイクルを確保する。デュアルコアの Q R A M 1 1 2 5 が前記キューのそれぞれに関する変数、前記キューの S R A M 状態に対応するヘッドライトポイント、ヘッドリードポイント、テールライトポイント及びテールリードポイント、並びに前記キューの D R A M 状態及び前記キューのサイズに対応するボディライトポイント及びボディリードポイントを含む各キューの変数を格納する。

【 0 0 9 4 】

アービタ 1 1 1 5 が実行すべき次の動作を選択した後 Q R A M 8 2 5 の変数は Q A L U 1 1 2 8 により選択された動作に従って設置されかつ変更され、S R A M 読出しリクエスト 1 1 3 0 または S R A M 書込みリクエスト 1 1 4 0 が生成される。前記変数が更新されかつ更新されたステータスが Q R A M 1 1 2 5 に加えてステータスレジスタ 1 1 2 2 に格納される。前記ステータスはまたアービタ 1 1 1 5 に供給されて、先に要求された動作が完了したという信号を送り、リクエストの重複を抑制する。ステータスレジスタ 1 1 2 2 は、4 つのキューレジスタ Q - O u t \_ R e a d y 1 0 4 5、Q - I n \_ R e a d y 1 0 5 0、Q - E m p t y 1 0 5 5 及び Q - F u l l 1 0 6 0 を更新して、アクセスされた前記キューの新しいステータスを反映する。同様に、D M A を介してそのキューについて S R A M ヘッド部分及び末尾部分にアクセスしかつそれからアクセスされた S R A M アドレス 1 1 3 3、ボディ書込みリクエスト 1 1 3 5 及びボディ読出しリクエスト 1 1 3 8 が更新される。これに代えて、Q 書込みデータ 1 1 4 4 で示すように、マルチプラクサ 1 1 4 6 により選択されかつ S R A M 書込みリクエスト 1 1 4 0 にパイプライン化されたキューに様々な処理を書き込みたい場合がある。S R A M コントローラは、アクセスされた前記キューの末尾部分に書き込みまたはヘッド部分を読み出し、かつ応答を返すことによって、読出し及び書込みリクエストを処理する。このようにして、様々なキューが用いられかつそれらのステータスが更新される。

【 0 0 9 5 】

図 2 2 A ~ D は、どのコンテキストまたは C C B を I N I C キャッシュメモリ内に保持するかを選択するために使用される最小使用頻度レジスタ 1 2 0 0 を示している本実施例の I N I C は、所定の時間で最大 6 0 の C C B を S R A M にキャッシュすることができ、従って新しい C C B がキャッシュされると、古いものは多くの場合捨てなければならず、捨てられる C C B は通常このレジスタ 1 2 0 0 によって、その時点で最も使用されていない C C B が選択される。本実施例では、最大 2 5 6 の C C B のハッシュテーブルが S R A M に保持されるのに対し、最大 2 5 6 の満杯の C C B が D R A M に保持される。最小使用頻度レジスタ 1 2 0 0 は、R 0 ~ R 1 5 とラベル付けされた 1 6 の 4 ビットブロックを有し、そのそれぞれが S R A M キャッシュ装置に対応している。初期設定時、前記ブロックには 0 ~ 1 5 の番号が付され、番号 0 が前記最小使用頻度 ( L R U ) キャッシュ装置を表す

ブロックに任意に格納され、かつ番号15が最大使用頻度(MRU)キャッシュ装置を表すブロックに格納される。図22Aは、LRUブロックR0が番号9を保持し、かつMRUブロックR15が番号6を保持する任意の時点におけるレジスタ1200を示している。現在SRAMに保持されている別のCCBをキャッシュしようとする場合、LRUブロックR0が読み出されこれは図22Aにおいて番号9を保持し、かつ新しいCCBが番号9に対応するSRAMキャッシュ装置に格納される。番号9に対応する新しいCCBはこの時点で最大使用頻度CCBであるから、番号9が、図22Bに示すように、MRUブロックに格納される。他の番号は、番号1をLRUブロックに残して、全て1レジスタブロックずつ左へシフトされる。番号9に対応してSRAM装置に以前キャッシュされていたCCBは、速度は遅いがよりコスト効率的なDRAMに移動される。

10

#### 【0096】

図22Cは、使用された次のCCBがすでにSRAMにキャッシュされた結果を示している。この実施例では、番号10に対応するSRAM装置にCCBがキャッシュされており、従ってそのCCBを用いた後、番号10がMRUブロックに格納される。以前により多く使用されていた番号10以外の番号のみ(レジスタブロックR9~R15)が、番号1をLRUブロックに残して左側へシフトされる。このようにして、INICはSRAMキャッシュ内に最もアクティブなCCBを保持する。

#### 【0097】

ある場合には、使用されているCCBは、制限されたキャッシュメモリ内に保持することが望ましくないものである。例えば、閉じつつあることが分かっているコンテキストについてCCBをキャッシュしないことが好ましく、従って他のキャッシュされたCCBをSRAMにより長く留めることができる。この場合には、デキャッシュ可能なCCBを保持するキャッシュ装置を表す番号が、MRUブロックR15ではなくむしろLRUブロックR0に格納され、従ってデキャッシュ可能なCCBは、LRUブロックR0に保持される番号に対応するSRAM装置内にキャッシュされた新しいCCBの使用の直後に置き換えられることになる。図22Dは、番号8(図22CにおいてブロックR9にあったもの)が使用されかつその後閉じることになっているCCBに対応する場合を示している。この場合、番号8はブロックR9から取り外され、かつLRUブロックR0内に格納される。先にブロックR9の左側に格納されていた全ての数(R1~R8)は、ここで1ブロック右側へシフトされる。

20

30

#### 【0098】

図23は、最小使用頻度レジスタ1200を動作するのに用いられる論理装置のいくつかを示している。16個の3または4入力マルチプレクサ1210のアレイ、その内のマルチプレクサMUX0、MUX7、MUX8、MUX9及びMUX15のみが簡単にするために示されているが、最小使用頻度レジスタ1200の対応する16個のブロックに接続された出力を有する。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。各レジスタブロックの値は、ブロック番号をシフトする際に使用するために、その対応するマルチプレクサの入力及び両側に隣接するマルチプレクサの入力に接続されている。例えば、R8に格納される番号はMUX7、MUX8及びMUX9の入力に供給される。MUX0及びMUX15はそれぞれ隣接するブロックを1つだけ有し、それらマルチプレクサの特別の入力がそれぞれLRU及びMRUブロックの選択のために使用される。MUX15は4入力マルチプレクサとして図示されており、入力1215がR0に格納された番号を供給する。

40

#### 【0099】

16個のコンパレータ1220からアレイはそれぞれ、最少使用頻度レジスタ1200の対応するブロックに格納された値を受け取る。また、各コンパレータは、ライン1235に沿ってプロセッサ470から信号を受け取り、それによってプロセッサ470により送られたものと整合する番号を有するレジスタブロックが、論理回路1230に真を出力するのに対し、他の15個のコンパレータが偽を出力する。論理回路1230が、マルチプレクサへの入力を選択しかつ従ってレジスタブロック番号のシフトを制御するために、

50

前記マルチプレクサのされぞれに繋がる一対の選択ラインを制御する。このように、選択ライン1239がMUX0を制御し、選択ライン1244がMUX7を制御し、選択ライン1249がMUX8を制御し、選択ライン1254がMUX9を制御し、かつ選択ライン1259がMUX15を制御する。

#### 【0100】

CCBを使用しようとする時、プロセッサ470は、CCBが現在16個のキャッシュ装置の1つに保持されているCCBと整合するかどうかを調べる。整合が認められると、前記プロセッサは、そのキャッシュ装置に対応するブロック番号、例えば番号12と共に信号をライン1235に沿って送る。コンパレータ1220は、そのライン1235からの前記信号をブロック番号と比較し、かつコンパレータC8が前記信号に整合するブロックR8について真を出力するのに対し、他の全てのコンパレータは偽を出力する。論理回路1230は、プロセッサ470からの制御下で、選択ライン1259を用いて、MUX15についてライン1235からの入力を選択し、番号12をMRUブロックR15に格納する。また、論理回路1230は、MUX15は別にして、MUX8及びそれより上位のマルチプレクサに関する各対の選択線に信号を送り、MUX8及びそれより上位の各マルチプレクサへの入力として、1ブロック右側のレジスタブロック(R9~R15)に格納されていた値を選択することによって、その出力を1ブロック左側へシフトする。MUX8の左側にあるマルチプレクサの出力は、一定に選択される。

10

#### 【0101】

他方、プロセッサ470が16個のキャッシュ装置の中で前記CCBとの整合を見つけない場合には、前記プロセッサはライン1266に沿ってLRUブロックR0から読出しを行い、LRUブロックに対応するキャッシュを識別し、かつそのキャッシュに格納されていたデータをDRAMに書き込む。R0に格納されていた番号、本実施例では番号3、がMRUブロックR15に格納するためにMUX15への入力1215として選択ライン1259により選択される。他の15のマルチプレクサは、そのそれぞれのレジスタブロックに、各レジスタブロックのすぐ右側に格納されていた番号を出力する。

20

#### 【0102】

前記プロセッサがCCBを使用後にキャッシュから取り外したいような場合には、MRUブロックR15よりもむしろLRUブロックR0が、そのCCBを保持するキャッシュ装置に対応する番号を置くために選択される。SRAMから取り外すためにLRUブロックR0に置かれるべきCCBに対応する数(例えば、ブロックR9に保持される番号1)は、プロセッサ470によりライン1235に沿って送られ、これはコンパレータC9によって整合される。前記プロセッサは論理回路1230に指示して、ライン1239でMUX0への入力1235を選択することによって、番号1をR0に入力させる。MUX9への選択ライン1254が入力として、レジスタブロックR8に保持されている番号を選択し、それによってR8からの番号がR9に格納される。R0とR9間の他のレジスタブロックにより保持されている番号が同様に右側へシフトされ、R9の右側のレジスタブロックの番号が一定のまま残される。これによって、多くのサイクルについて閉じたCCBを保持することから不十分なキャッシュメモリが解放され、それらの識別番号がレジスタブロックを介してMRUからLRUブロックへ移動する。

30

40

#### 【0103】

図24は、INICの詳細を追加して示しており、この記載では単一のネットワークシステムに焦点を置いている。INIC22は、PHYチップ712、ASICチップ700及びDRAM755を有する。PHYチップ712はINICカード22をネットワークコネクタ2101を介してネットワークライン2105に接続している。INIC22は、カードエッジコネクタ2107及びPCIバス757を介してホストのCPU(例えば、図1のホスト20のCPU30)に接続されている。ASICチップ700は、媒体アクセスコントロール(MAC)装置722、シーケンサブロック732、SRAMコントロール744、SRAM748、DRAMコントロール742、キューマネージャ803、プロセッサ780及びPCIバスインタフェース装置756を有する。シーケンサブ

50

ック732は、送信シーケンサ2104、受信シーケンサ2105、及びコンフィグレーションレジスタ2106を有する。MACデスティネーションアドレスはコンフィグレーションレジスタ2106に格納される。プロセッサ780により実行されるプログラムコードの一部分がROM(図示せず)に保有され、かつ一部分が書込み可能なコントロールストアSRAM(図示せず)に送られている。前記プログラムは、初期設定時にホスト20から書込み可能コントロールストアSRAMにダウンロードすることができる。

#### 【0104】

図25は、図24の受信シーケンサ2105のより詳細な線図である。受信シーケンサ2105は、データ同期バッファ2200、パケット同期シーケンサ2201、データセンブリレジスタ2202、プロトコルアナライザ2203、パケットプロセッシングシーケンサ2204、キューマネージャインタフェース2205、及びダイレクトメモリアクセス(DMA)コントロールブロック2206を有する。パケット同期シーケンサ2201及びデータ同期バッファ2200はMAC722のネットワーク同期式クロックを利用し、受信シーケンサ2105の残りの部分は固定周波数のクロックを利用する。破線2221は、クロックのドメイン境界を示している。

#### 【0105】

図24及び図25の受信シーケンサ2105の動作について、ネットワークライン702からのTCP/IPパケットのINIC22上での受信に関連して説明する。初期設定時には、プロセッサ780がDRAM755をバッファに区分する。受信シーケンサ2105がDRAM755のバッファを用いて、到来するネットワークパケットデータを前記パケットのステータス情報と共に格納する。プロセッサ780は各バッファについて32ビットのバッファ記述子を作成する。バッファ記述子は、その関連するバッファのDRAMにおけるサイズ及びロケーションを示す。プロセッサ780は、これらのバッファ記述子をキューマネージャ803に書き込むことによって、「フリーバッファキュー」2108上にバッファ記述子を置く。キューマネージャ803は、「フリーバッファキュー」2108を含む複数のキューを保持する。この実施例では、様々なキューのヘッド部分及び末尾部分がSRAM748に置かれ、前記キューの中間部分がDRAM755に置かれる。

#### 【0106】

ライン2229は、リクエストライン及びアドレスラインに関連するリクエストメカニズムを有する。同様にライン2230は、リクエストライン及びアドレスラインに係するリクエストメカニズムを有する。キューマネージャ803はライン2229、2230を用いて、DRAMからSRAMへまたはSRAMからDRAMへキュー情報を転送するリクエストを出す。

#### 【0107】

受信シーケンサのキューマネージャインタフェース2205は、パケットプロセッシングシーケンサ2204による使用のために、常にフリーバッファ記述子2207を保持しようとする。ビット2208は、フリーバッファ記述子2207がパケットプロセッシングシーケンサ2204により使用可能であることを表示するレディビットである。キューマネージャインタフェース2205がフリーバッファ記述子を持っていない(ビット2208が設定されていない)場合には、キューマネージャインタフェース2205はリクエストライン2205を介してキューマネージャ803から1を要求する。(リクエストライン2209は実際、動作がキューへの書込みである場合、リクエスト、キューID、読出し/書込み信号及びデータを通信するバスである。)

応答の際、キューマネージャ803は「フリーバッファキュー」2108の末尾部分からフリーバッファ記述子を検索し、次に応答ライン2210上の応答信号を介してキューマネージャインタフェース2205に知らせる。キューマネージャインタフェース2205が応答信号を受け取ると、キューマネージャインタフェース2205はフリーバッファ記述子2207をロードし、かつレディビット2208を設定する。前記フリーバッファ記述子がSRAM748のフリーバッファキューの末尾部分にあったことから、キューマネージャインタフェース2205は、実際はSRAMコントロールブロック744の読出し

10

20

30

40

50

データバス 2 2 2 8 からフリーバッファ記述子 2 2 0 7 を受け取る。パケットプロセッシングシーケンサ 2 2 0 4 は、リクエストライン 2 2 1 1 を介してフリーバッファ記述子 2 2 0 7 を要求する。キューマネージャインタフェース 2 2 0 5 がフリーバッファ記述子 2 2 0 7 を検索し、かつフリーバッファ記述子 2 2 0 7 がパケットプロセッシングシーケンサによる利用が可能である場合、キューマネージャインタフェース 2 2 0 5 はグラントライン 2 2 1 2 を介してパケットプロセッシングシーケンサ 2 2 0 4 に知らせる。この処理によって、フリーバッファ記述子はパケットプロセッシングシーケンサ 2 2 0 4 による使用が可能になり、かつ受信シーケンサ 2 1 0 5 が到来するパケットを処理できる状態となる。

#### 【 0 1 0 8 】

次に、TCP/IP パケットがネットワークコネクタ 2 1 0 1 及び物理的レイヤインタフェース (PHY) 7 1 2 を介してネットワークライン 2 1 0 5 から受け取られる。PHY 7 1 2 は、前記パケットを MAC 7 2 2 に媒体独立インタフェース (MII) パラレルバス 2 1 0 9 を介して供給する。MAC 7 2 2 は前記パケットの処理を開始し、かつパケットの開始部分が受信されていることを示す「パケットの開始」信号をライン 2 2 1 3 上にアサートする。データのバイトが MAC に受信されかつ MAC 出力 2 2 1 5 において利用可能になると、MAC 7 2 2 は「データ有効」信号をライン 2 2 1 4 上にアサートする。前記「データ有効」信号を受け取ると、パケット同期シーケンサ 2 2 0 1 がロード信号ライン 2 2 2 2 を介してデータ同期バッファ 2 2 0 0 を指示し、データライン 2 2 1 5 から受け取ったバイトをロードする。データ同期バッファ 2 2 0 0 は 4 バイトの深さを有する。次に、パケット同期シーケンサ 2 2 0 1 がデータ同期バッファの書き込みポインタを増分させる。このデータ同期バッファの書き込みポインタは、ライン 2 2 1 6 を介してパケットプロセッシングシーケンサ 2 2 0 4 に利用可能となる。データライン 2 2 1 5 から連続するデータのバイトは、このようにしてデータ同期バッファ 2 2 0 0 内に刻時される。

#### 【 0 1 0 9 】

ライン 2 2 1 9 上で利用可能なデータ同期バッファ読出しポインタがパケットプロセッシングシーケンサ 2 2 0 4 によって保持される。パケットプロセッシングシーケンサ 2 2 0 4 は、ライン 2 2 1 6 上のデータ同期バッファ書き込みポインタをライン 2 2 1 9 上のデータ同期バッファ読出しポインタと比較することによって、データ同期バッファ 2 2 0 0 においてデータが使用可能であることを決定する。

#### 【 0 1 1 0 】

データアセンブリレジスタ 2 2 0 2 は 1 6 バイト長のシフトレジスタ 2 2 1 7 を有する。このレジスタ 2 2 1 7 は、一度に 1 バイトが順にロードされ、かつパラレルにアンロードされる。データがレジスタ 2 2 1 7 にロードされると、書込みポインタが増分される。この書込みポインタが、ライン 2 2 1 8 を介してパケットプロセッシングシーケンサ 2 2 0 4 に利用可能となる。同様に、データがレジスタ 2 2 1 7 からアンロードされると、パケットプロセッシングシーケンサ 2 2 0 4 により保持された読出しポインタが増分される。この読出しポインタは、ライン 2 2 2 0 を介してデータアセンブリレジスタ 2 2 0 2 に利用可能である。従って、パケットプロセッシングシーケンサ 2 2 0 4 は、ライン 2 2 1 8 上の書込みポインタをライン 2 2 2 0 上の読出しポインタと比較することによって、レジスタ 2 2 1 7 に利用可能な余裕があるかどうかを決定することができる。

#### 【 0 1 1 1 】

レジスタ 2 2 1 7 に利用可能な余裕があるとパケットプロセッシングシーケンサ 2 2 0 4 が決定した場合、パケットプロセッシングシーケンサ 2 2 0 4 はデータアセンブリレジスタ 2 2 0 2 に命令して、データ同期バッファ 2 2 0 0 からデータのバイトをロードする。データアセンブリレジスタ 2 2 0 2 は、ライン 2 2 1 8 上のデータアセンブリレジスタ書込みポインタを増分し、かつパケットプロセッシングシーケンサ 2 2 0 4 がライン 2 2 1 9 上のデータ同期バッファ読出しポインタを増分する。レジスタ 2 2 1 7 にシフトされたデータは、チェックサムを検証しかつ「ステータス」情報 2 2 2 3 を生成するプロトコルアナライザ 2 2 0 3 によってレジスタ出力で調べられる。

10

20

30

40

50

## 【0112】

DMAコントロールブロック2206は、64バイト受信FIFO2110を介して情報をレジスタ2217からバッファ2114へ動かす責任を有する。DMAコントロールブロック2206は、64バイトのSRAM748を用いて受信FIFO2110を2つの32バイトピンポンバッファとして実行する。DMAコントロールブロック2206は、書込みポインタ及び読出しポインタを用いて前記受信FIFOを実行する。転送されるべきデータがレジスタ2217において利用可能でありかつ空間がFIFO2110において利用可能であると、DMAコントロールブロック2206は、ライン2225を介してSRAMコントローラ744にSRAM書込みリクエストをアサートする。次にSRAMコントローラ744がデータをレジスタ2217からFIFO2110へ移動させ、かつライン2225を介して応答信号をDMAコントロールブロック2206へアサートする。ここで、DMAコントロールブロック2206が受信FIFO書込みポインタを増分させ、かつデータアセンブリレジスタ読出しポインタを増分させる。

10

## 【0113】

32バイトのデータが受信FIFO2110に預けられていると、DMAコントロールブロック2206はライン2226を介してDRAMコントローラ742にDRAM書込みリクエストを示す。この書込みリクエストは、DRAMリクエストアドレスに関して「バッファロードカウンタ」でOR演算したフリーバッファ記述子2207と、SRAM読出しアドレスに関する受信FIFO読出しポインタとから構成される。受信FIFO読出しポインタを用いて、DRAMコントローラ742はSRAM744に読出しリクエストをアサートする。SRAMコントローラ744は、SRAM748の受信FIFO2110から表示されたデータを戻しかつ応答信号をアサートすることによって、DRAMコントローラ742に応答する。DRAMコントローラ742は、前記データをDRAM書込みデータレジスタに格納し、DRAMリクエストアドレスをDRAMアドレスレジスタに格納し、かつ応答DMAコントロールブロック2206にアサートする。次に、DMAコントロールブロック2206は、受信FIFO読出しポインタをデクリメントする。次に、DRAMコントローラ742は、前記データをDRAM書込みデータレジスタからバッファ2114へ移動させる。このようにして、連続的な32バイトのチャンクからなるデータがSRAM748に格納されるので、DRAMコントロールブロック2206は、これら32バイトのチャンクからなるデータを一度にSRAM748からDRAM755のバッファ2214に移動させる。32バイトチャンクのデータをこのようにしてDRAM755に転送することによって、DRAMの比較的効率的なバーストモードを用いて、データをDRAMに書き込むことができる。

20

30

## 【0114】

パケットデータは、全てのパケットデータが受け取られるまで、ネットワークライン2105からバッファ2114へ流れ続ける。次に、MAC722は、ライン2227上に「フレーム終了」（即ちパケットの終了）信号をアサートすることによって、かつ最後のパケットステータス（MACパケットステータス）をパケット同期シーケンサ2204に示すことによって、到来するパケットが完了したことを表す。次にパケットプロセッシングシーケンサ2204が、最後にバッファ2114に転送するためにステータス2223（「プロトコルアナライザステータス」とも称する）及びMACパケットステータスをレジスタ2217へ移動させる。前記パケットの全データがバッファ2214内に置かれた後、ステータス2223及びMACパケットステータスは、図23に示すように関連するデータにプリPENDされて格納されるように、バッファ2214に転送される。

40

## 【0115】

全データ及びステータスがバッファ2114に転送された後、パケットプロセッシングシーケンサ2204が、フリーバッファ記述子2207、バッファロードカウンタ、MAC・ID及びステータスビット（「アテンションビット」とも称する）を連結することによって、サマリ2224（「受信パケット記述子」とも称する）を作成する。前記アテンションビットが1である場合、前記パケットは「高速パス対象」でなく、前記アテンション

50

ビットが0である場合には、前記パケットは「高速バス対象」となる。アテンションビットの値は、さもなければプロセッサ780が前記パケットが「高速バス対象」であるかどうかを決定するためにしなければならない多大な量の処理の結果を表す。例えば、アテンションビットが0であることは、前記パケットがTCPプロトコル及びIPプロトコルの両方を用いていることを示している。この多大な量の処理を事前にハードウェアで実行し、かつ次に結果をアテンションビットにエンコードすることによって、パケットが本当の「高速バスパケット」であるかどうかに関してプロセッサ780が次に行なう決定が加速される。

【0116】

次に、パケットプロセッシングシーケンサ2204がサマリ2224に関連するレディビット（図示せず）を設定し、かつサマリ2224をキューマネージャインタフェース2205に提示する。次にキューマネージャインタフェース2205が、「サマリキュー」2112（「受信記述子キュー」とも称する）のヘッド部分に書込みを要求する。キューマネージャ803はこのリクエストを受け取り、サマリ2224をサマリキュー2212のヘッドに書き込み、かつライン2210を介してキューマネージャインタフェースに応答信号をアサートする。キューマネージャインタフェース2205が前記応答を受け取ると、キューマネージャインタフェース2205は、前記サマリに関連するレディビットをクリアすることによって、サマリ2224がサマリキュー2212にあることをパケットプロセッシングシーケンサ2204に知らせる。また、パケットプロセッシングシーケンサ2204は、MACパケットステータスとMAC・IDを連結することによって、前記パケットに関する追加のステータス情報（「ベクタ」とも称する）を生成する。パケットプロセッシングシーケンサ2204は、このベクタに関連するレディビット（図示せず）を設定し、かつこのベクタをキューマネージャインタフェース2205に提示する。キューマネージャインタフェース2205及びキューマネージャ803は、上述したようにサマリ2224がサマリキュー2112のヘッドに書き込まれたのと同様に、このベクタを協同して「ベクタキュー」2113のヘッドに書き込む。前記パケットのベクタがベクタキュー2113に書き込まれると、キューマネージャインタフェース2205は、前記ベクタに関連するレディビットをリセットする。

【0117】

サマリ2224（バッファ2114を指すバッファ記述子を含む）が一旦サマリキュー2112に置かれ、かつパケットデータがバッファ2114に置かれると、プロセッサ780はサマリ2224をサマリキュー2112から検索して、前記「アテンションビット」を調べることができる。

【0118】

サマリ2224からのアテンションビットが数字の1である場合、プロセッサ780は、前記パケットが「高速バス対象」でないと決定し、プロセッサ780はパケットヘッダを調べる必要がない。バッファ2114からのステータス2223（最初の16バイト）のみがSRAMにDMA転送され、従ってプロセッサ780がこれを調べることができる。ステータス2223は、前記パケットがホストに転送されることになっていない種類のパケットであること（例えば、ホストが受け取るようにレジスタされていないマルチパストフレーム）を示す場合には、前記パケットは捨てられる（即ち、ホストに送られない）。ステータス2223が、前記パケットがホストに転送されるべきでないタイプのパケットであることを示していない場合には、パケット全体（ヘッダ及びデータ）が、「低速バス」移動及びホスト20のプロトコルスタックによるネットワークレイヤプロセッシングのためにホスト20のバッファに送られる。

【0119】

前記アテンションビットが0である場合、プロセッサ780は、前記パケットが「高速バス対象」であることを決定する。プロセッサ780が前記パケットを「高速バス対象」と決定した場合、プロセッサ780は前記サマリからのバッファ記述子を用いて、バッファ2114からの最初の約96バイトの情報をDRAM755からSRAM748の部分内

10

20

30

40

50

にDMA転送し、それによってプロセッサ780がこれを調べることができる。この最初の約96バイトは、IPヘッダのIPソースアドレス、IPヘッダのIPデスティネーションアドレス、TCPヘッダのTCPソースアドレス及びTCPヘッダのTCPデスティネーションアドレスと共に、ステータス2223を含む。IPヘッダのIPソースアドレス、IPヘッダのIPデスティネーションアドレス、TCPヘッダのTCPソースアドレス及びTCPヘッダのTCPデスティネーションアドレスは一緒になって、前記パケットが関連する1つの接続コンテキスト(TCB)を独自に定義する。プロセッサ780は、これらのTCP及びIPヘッダのアドレスを調べ、前記パケットの接続コンテキストを決定する。次に、プロセッサ780は、INIC22の制御下にある接続コンテキストのリストを調査し、前記パケットがINIC22の制御下で接続コンテキスト(TCB)に關連しているかどうかを決定する。

10

#### 【0120】

前記接続コンテキストが前記リストにない場合、前記「高速パス対象」パケットは「高速パスパケット」でないと決定される。そのような場合、パケット全体(ヘッダ及びデータ)は、ホスト20のプロトコルスタックによる「低速パス」処理のためにホスト20のバッファに転送される。

#### 【0121】

他方、前記接続コンテキストが前記リストにある場合、ソフトウェアステートマシン2231、2232を含むプロセッサ780により実行されるソフトウェアが、多くの例外状態の中の1つを調べて、前記パケットが「高速パスパケット」であるか「高速パスパケット」でないかを決定する。これらの例外状態に含まれるのは、1)IP断片化が検出されたこと、2)IPオプションが検出されたこと、3)予期しないTCPフラグ(緊急ビットセット、リセットビットセット、SYNビットセットまたはFINビットセット)が検出されたこと、4)TCPヘッダのACKフィールドがTCPウィンドウの前にあり、またはTCPACKフィールドがTCPウィンドウの後にあり、またはTCPヘッダのACKフィールドがTCPウィンドウをシュリンクさせていること、5)TCPヘッダのACKフィールドが複製のACKでありかつACKフィールドが前記複製ACKのカウントを越えている(前記複製ACKのカウントがユーザが設定可能な値であること)、及び6)TCPヘッダのシーケンス番号が順番になっていない(パケットが順番どおりに受け取られていない)ことである。プロセッサ780により実行されるソフトウェアがこれら例外状態の1つを検出した場合、プロセッサ780は前記「高速パス対象」が「高速パスパケット」でないと決定する。このような場合、前記パケットの接続コンテキストは「フラッシュ」され(接続コンテキストはホストに戻される)、それによって前記接続コンテキストはもはやINIC22の制御下で接続コンテキストのリストに存在しないことになる。パケット全体(ヘッダ及びデータ)は、ホスト20のプロトコルスタックによる「低速パス」トランスポートレイヤ及びネットワークレイヤの処理のためにホスト20のバッファへ転送される。

20

30

#### 【0122】

他方、プロセッサ780がそのような例外状態を発見しない場合、前記「高速パス対象」のパケットは本当の「高速パスパケット」であると決定される。次に、受信ステートマシン2232がTCPを介して前記パケットの処理を行なう。次にバッファ2114の前記パケットのデータ部分が別のDMAコントローラ(図21には図示せず)によって、バッファ2114からホストストレージ35のホスト割り当て型ファイルキャッシュに転送される。ある実施例では、ホスト20は「高速パスパケット」のTCP及びIPヘッダの分析を行なわない。「高速パスパケット」のTCP及びIPヘッダの全ての分析はINICカード20上で行なわれる。

40

#### 【0123】

図26は、「高速パスパケット」(64kバイトセッションレイヤメッセージ2300のパケット)のデータのINIC22からホスト20への転送を示す線図である。破線2301の左側の線図の部分はINIC22を表し、破線2301の右側の線図の部分はホス

50



ト 20 を表す。64 k バイトのセッションレイヤメッセージ 2300 は約 45 個のパケットを有し、図 24 ではその内の 4 つ (2302、2303、2304、2305) にラベルが付されている。第 1 パケット 2302 は、トランスポート及びネットワークレイヤヘッダ (例えば、TCP 及び IP ヘッダ) を含む部分 2306 と、セッションレイヤヘッダを含む部分 2307 と、データを含む部分 2308 とを有する。最初のステップでは、部分 2307、部分 2308 からの最初の数バイトのデータ、及びパケット 2300 の接続コンテキスト識別子が INIC 22 からホスト 20 の 256 バイトのバッファ 2309 に転送される。第 2 ステップでは、この情報をホスト 20 が調査し、かつ INIC 22 に前記データのデスティネーション (例えばストレージ 35 におけるファイルキャッシュ 2311 のロケーション) を戻す。また、ホスト 20 は、バッファ 2309 からのデータの最初の数バイトをファイルキャッシュ 2311 の第 1 部分 2312 の開始部分にコピーする。第 3 ステップでは、INIC 22 が部分 2308 からの前記データの残りの部分をホスト 20 へ、前記データの残りの部分がファイルキャッシュ 2311 の第 1 部分 2312 の残りの部分に格納されるように転送する。ネットワーク、トランスポートまたはセッションレイヤヘッダはファイルキャッシュ 2311 の第 1 部分 2312 に格納されない。次に、第 2 パケット 2303 のデータ部分 2313 がホスト 20 に、第 2 パケット 2303 のデータ部分 2313 がファイルキャッシュ 2311 の第 2 部分 2314 に格納されるように転送する。第 2 パケット 2303 のトランスポートレイヤ及びネットワークレイヤヘッダ部分 2315 はホスト 20 に転送されない。第 1 パケット 2302 のデータ部分と第 2 パケット 2303 のデータ部分との間でファイルキャッシュ 2311 に格納されるネットワーク、トランスポートまたはセッションレイヤヘッダはない。同様に、第 2 セッションレイヤメッセージの次のパケット 2304 のデータ部分 2316 はファイルキャッシュ 2311 に、第 2 パケット 2303 のデータ部分とファイルキャッシュ 2311 の第 3 パケット 2304 のデータ部分との間でネットワーク、トランスポートまたはセッションレイヤヘッダが存在しないように転送される。このようにして、前記セッションレイヤメッセージのパケットのデータ部分のみがファイルキャッシュ 2311 に置かれる。セッションレイヤメッセージ 2300 からのデータはブロックとしてファイルキャッシュ 2311 内に存在し、それによってこのブロックはネットワーク、トランスポートまたはセッションレイヤヘッダを含まない。

#### 【0124】

より短い単一パケットのセッションレイヤメッセージの場合、前記セッションレイヤメッセージの部分 2307 及び 2308 は、上述した長いセッションレイヤメッセージの場合のように、接続コンテキスト識別子 2310 と共にホスト 20 の 256 バイトバッファ 2309 に転送される。単一パケットセッションレイヤメッセージの場合、しかしながら、この時点で転送が完了する。ホスト 20 はデスティネーションを INIC 22 に戻さず、かつ INIC 22 はそのようなデスティネーションに後のデータを転送しない。

#### 【0125】

全体として、上述したデータ通信を処理するためのデバイス及びシステムによって大きな接続ベースのメッセージを処理するために必要な時間及びホストリソースが大幅に低減される。プロトコル処理速度及び効率は従来のプロトコルソフトウェアを実行する汎用 CPU と比較して、特別に設計されたプロトコル処理ハードウェアを有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) により著しく加速され、かつホスト CPU への割り込みが大幅に減少する。これらの利益は、ネットワークストレージの用途の場合に大きく、そのような場合、ファイル転送からのデータはホストによりファイル転送の制御が維持されて、ホストメモリバス及びホストメモリ I/O バス双方を回避することができる。

#### 【図面の簡単な説明】

【図 1】 INIC に接続された記憶装置のための I/O コントローラ及びファイルキャッシュを有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) により複数のネットワークにより接続されたホストコンピュータを有するネットワークストレ

ジシステムを示す図である。

【図2】 本発明による複数のネットワーク間でデータを転送する際のI N I C及びホストコンピュータの機能を示す図である。

【図3】 図1のシステムによりネットワークからメッセージパケットを受信する際に必要なステップの順序を示すフロー図である。

【図4】 図1のシステムによりネットワークからのリクエストにตอบสนองしてメッセージパケットをネットワークに送る際に必要なステップの順序を示すフロー図である。

【図5】 ホストコンピュータにより管理される複数のI N I Cにより複数のネットワーク及び複数の記憶装置に接続されたホストコンピュータを有するネットワークストレージシステムを示す図である。

10

【図6】 I / Oコントローラ無しでインテリジェントネットワークインタフェースカード(I N I C)により複数のL A N及び複数のS A Nに接続されたホストコンピュータを有するネットワークストレージシステムを示す図である。

【図7】 ネットワークラインと記憶装置間に接続されたイーサネットS C S Iアダプタを有する、図6のS A Nの1つを示す図である。

【図8】 図6のイーサネット - S C S Iアダプタの1つを示す図である。

【図9】 ホストコンピュータにより管理される複数のI N I Cにより複数のL A N及び複数のS A Nに接続されたホストコンピュータを有するネットワークストレージシステムを示す図である。

【図10】 パケットコントロールシーケンサ及びフライバイシーケンサを有する、図1に示すI N I Cの実施例のためのハードウェア論理を示す図である。

20

【図11】 I N I Cにより受け取られるヘッダバイトを分析するための図10のフライバイシーケンサを示す図である。

【図12】 低速パスでパケットを処理することに加え、高速パスのために通信制御ブロックを作成しかつ制御するための図1の専用ホストプロトコルスタックを示す図である。

【図13】 N e t B i o s 通信用に構成されたマイクロソフトのT C P / I Pスタック及びAlacrytechのコマンドドライバを示す図である。

【図14】 ネットワーク記憶装置を有するサーバとクライアント間でのN e t B i o s 通信エクスチェンジを示す図である。

【図15】 図14のサーバとクライアント間で音声または映像データを転送することができるユーザデータグラムプロトコル(U D P)エクスチェンジを示す図である。

30

【図16】 図1のI N I Cに含まれるハードウェア機能を示す図である。

【図17】 それぞれにプロセッサを有する3つのフェーズを含む、図16のI N I Cに含まれる3つ組のパイプライン化マイクロプロセッサを示す図である。

【図18A】 図17のパイプライン化マイクロプロセッサの第1フェーズを示す図である。

【図18B】 図17のパイプライン化マイクロプロセッサの第2フェーズを示す図である。

【図18C】 図17のパイプライン化マイクロプロセッサの第3フェーズを示す図である。

40

【図19】 図17のマイクロプロセッサと対話しかつS R A M及びD R A Mを有する複数のキュー記憶装置を示す図である。

【図20】 図19のキュー記憶装置のためのステータスレジスタセットを示す図である。

【図21】 図19及び図20のキュー記憶装置及びステータスレジスタと対話するキューマネージャを示す図である。

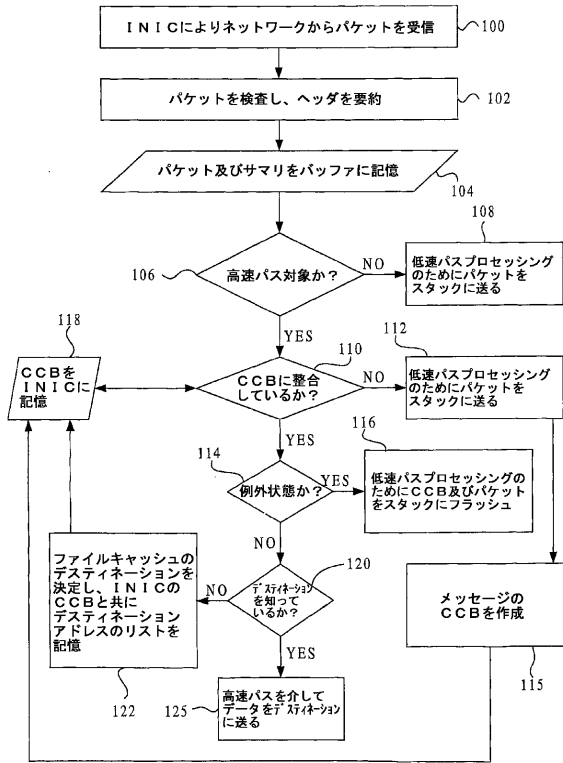
【図22】 A図~D図は、キャッシュメモリを割り当てるために用いられる最小使用頻度レジスタの様々な段階をそれぞれ示す図である。

【図23】 図22A~Dの最小使用頻度レジスタを動作するために用いられるデバイスを示す図である。

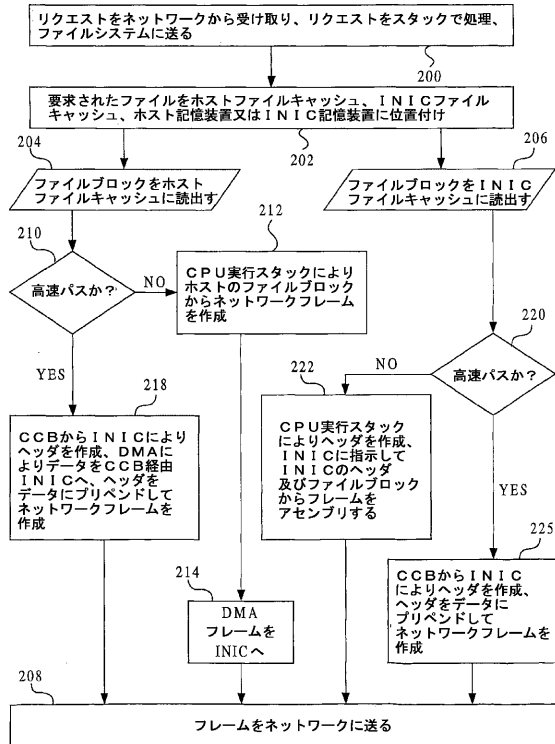
50



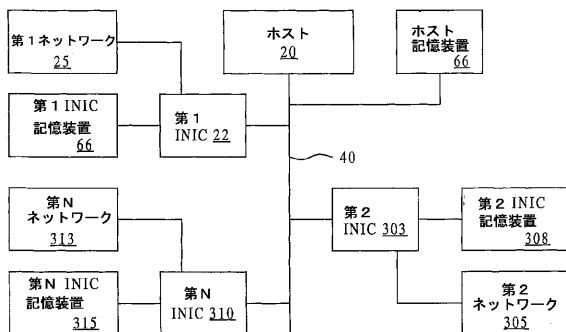
【図 3】



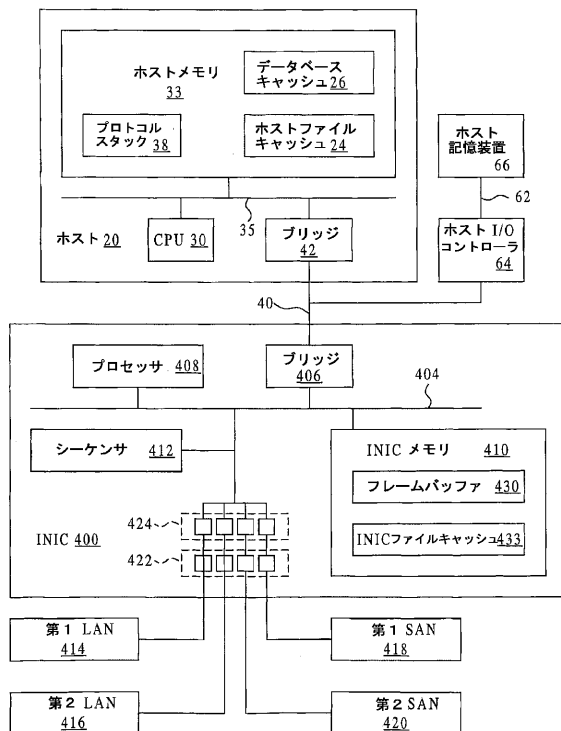
【図 4】



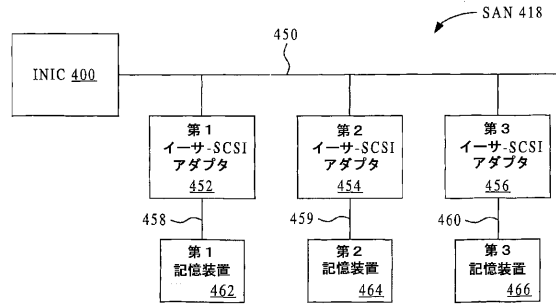
【図 5】



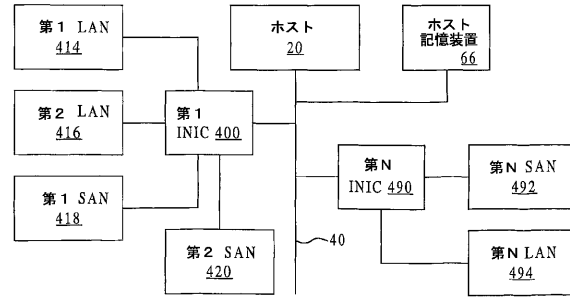
【図 6】



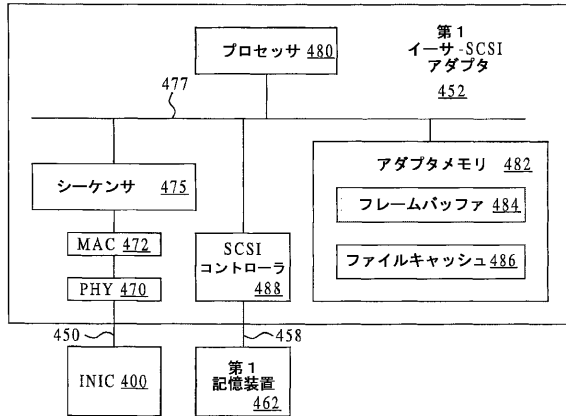
【図 7】



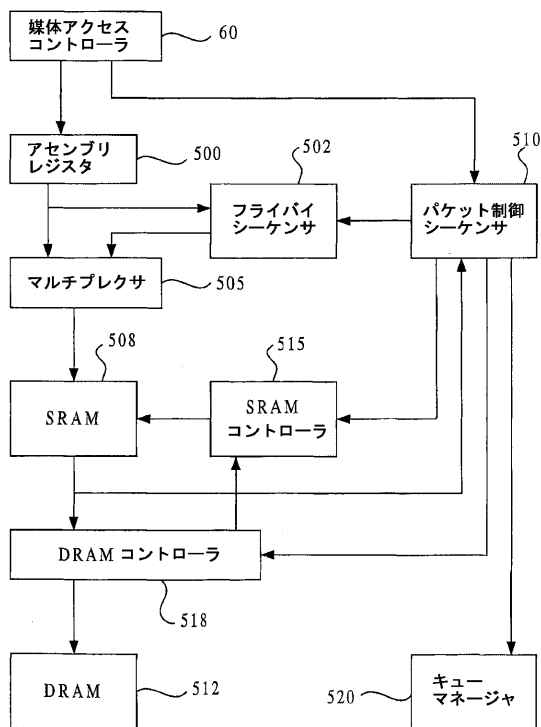
【図 9】



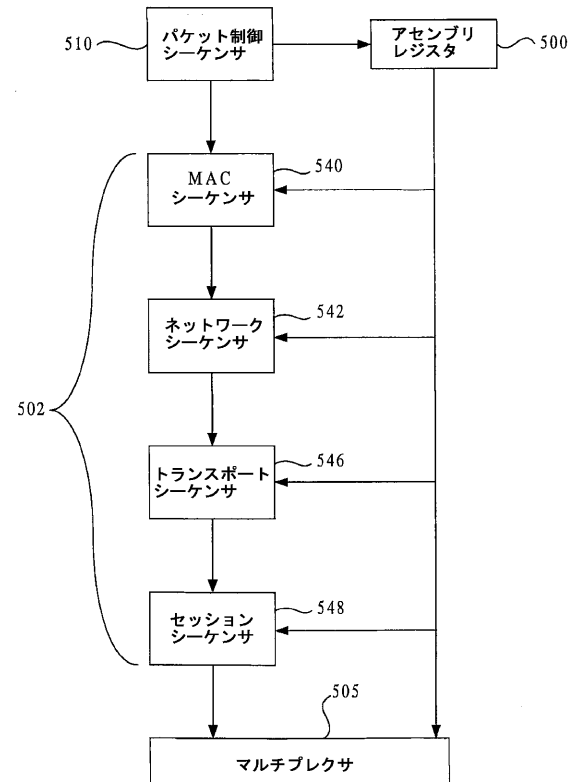
【図 8】



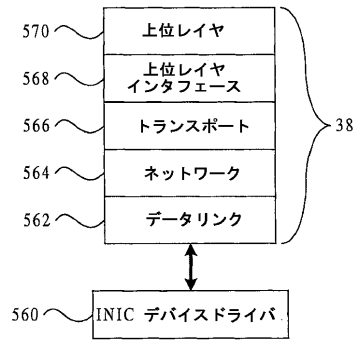
【図 10】



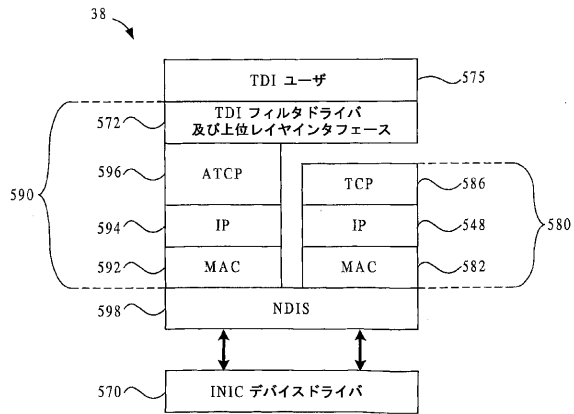
【図 11】



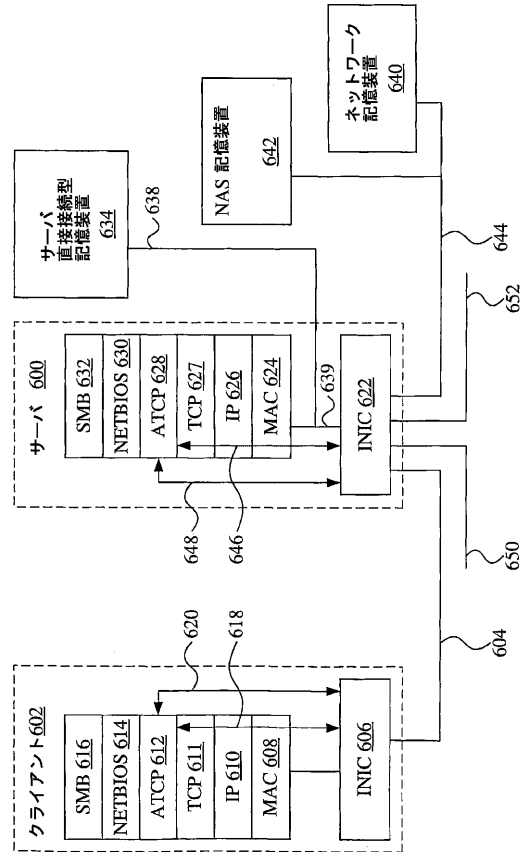
【図 12】



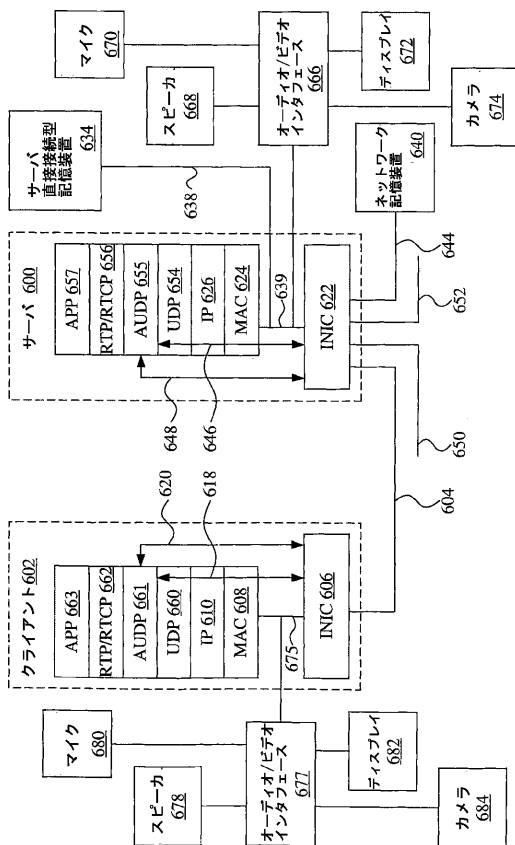
【図 13】



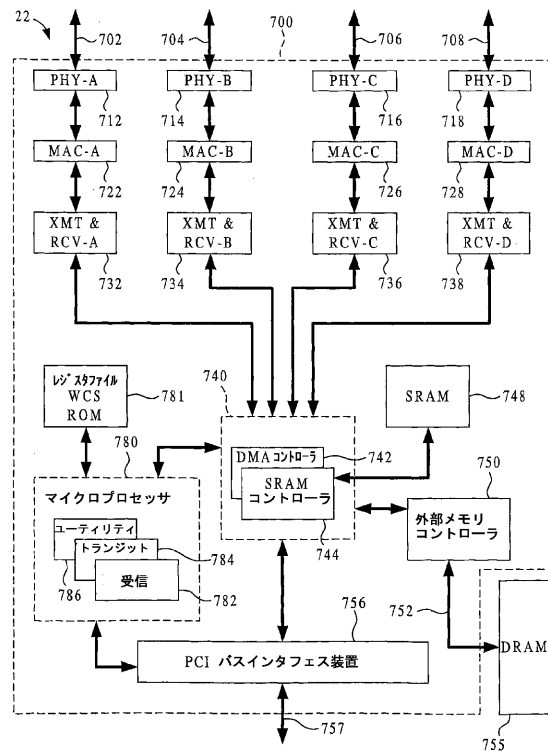
【図 14】



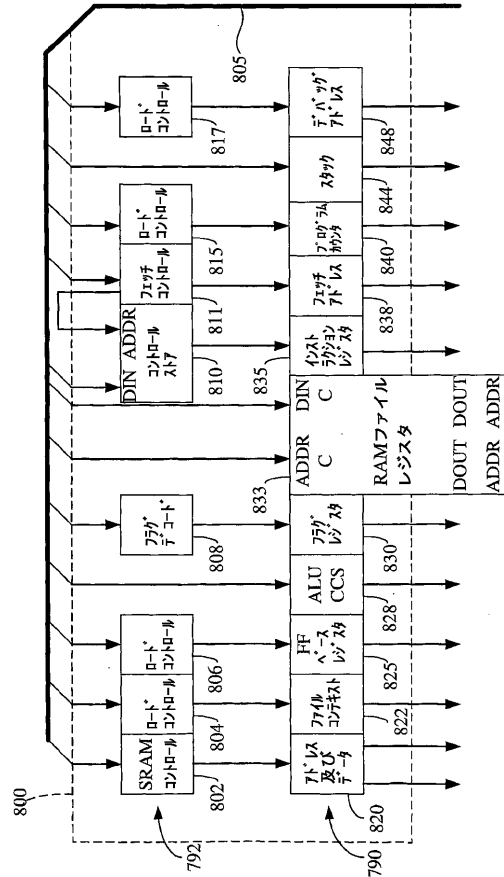
【図 15】



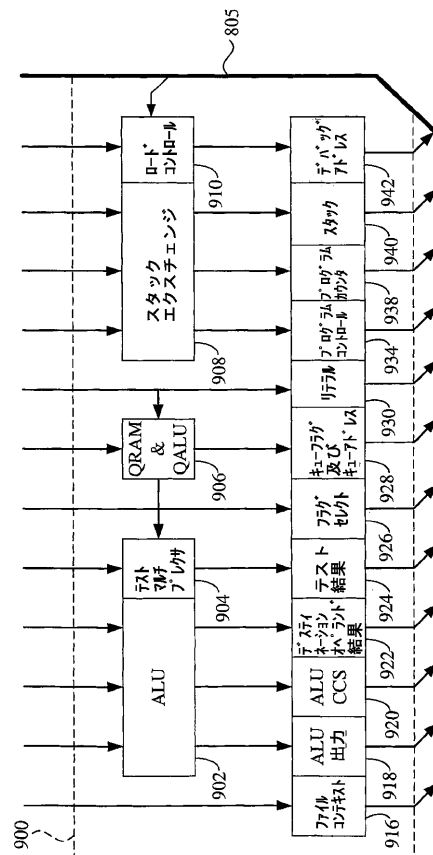
【図 16】



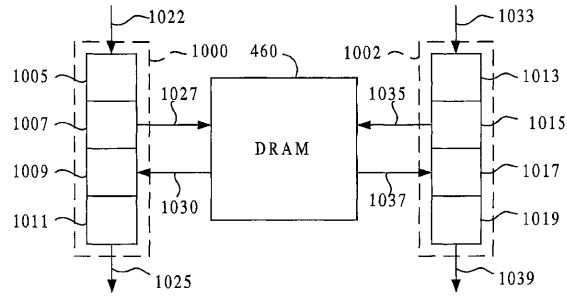
【 図 1 8 - A 】



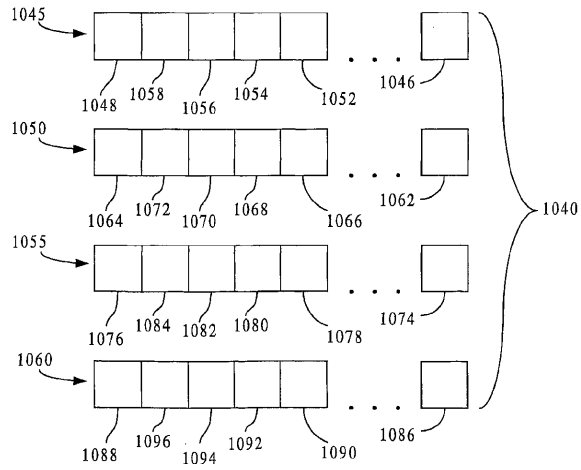
【 図 1 8 - C 】



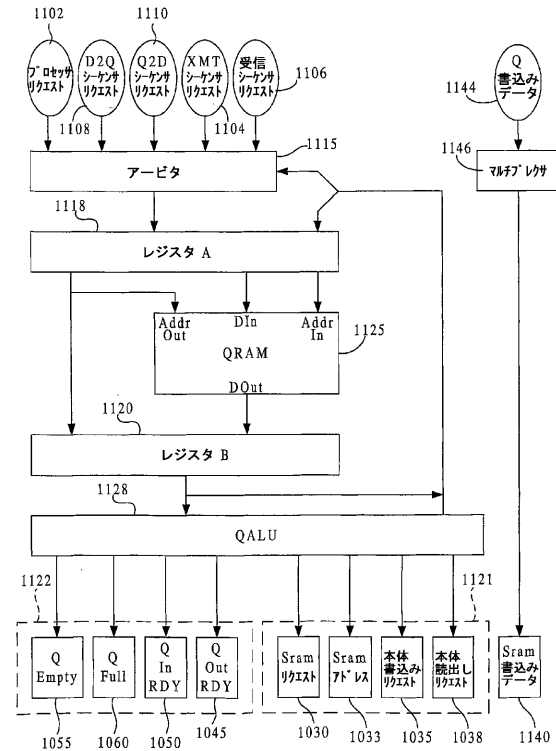
【図 19】



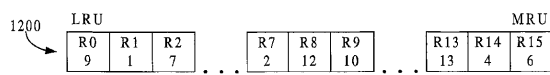
【図 20】



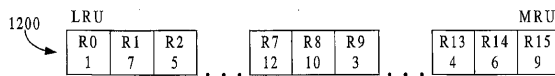
【図 21】



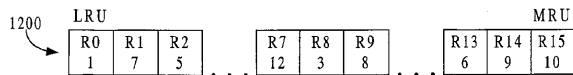
【図 22 - A】



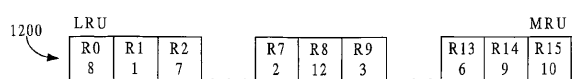
【図 22 - B】



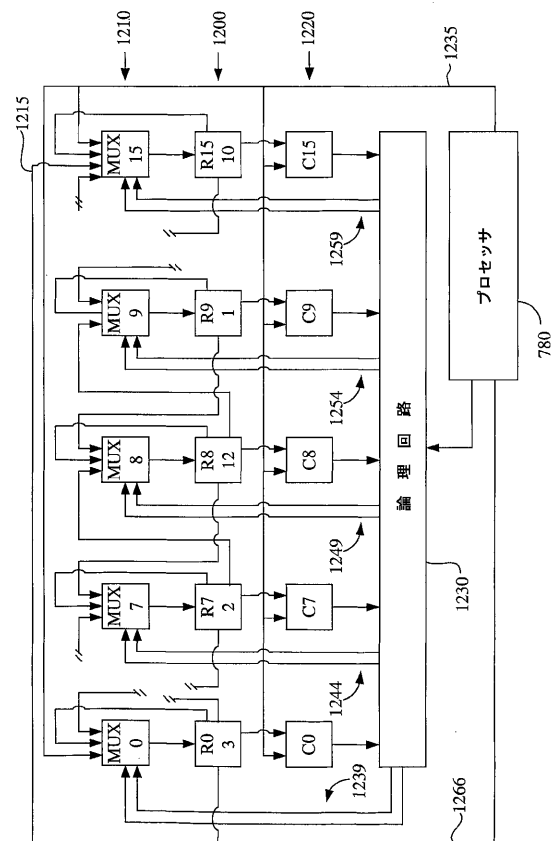
【図 22 - C】



【図 22 - D】

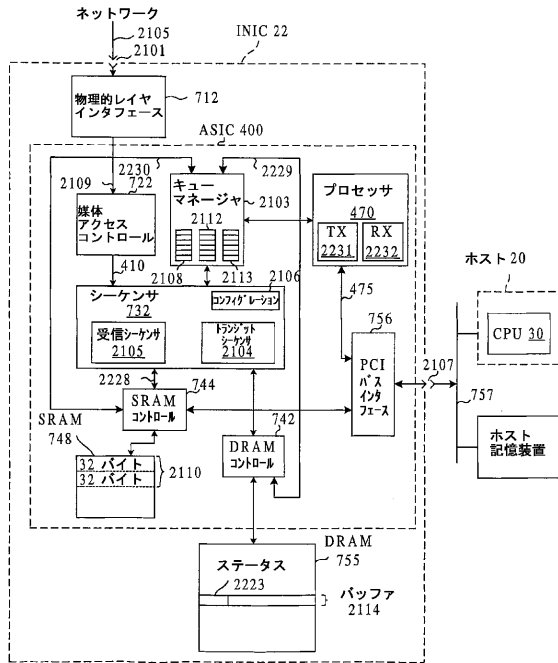


【図 23】

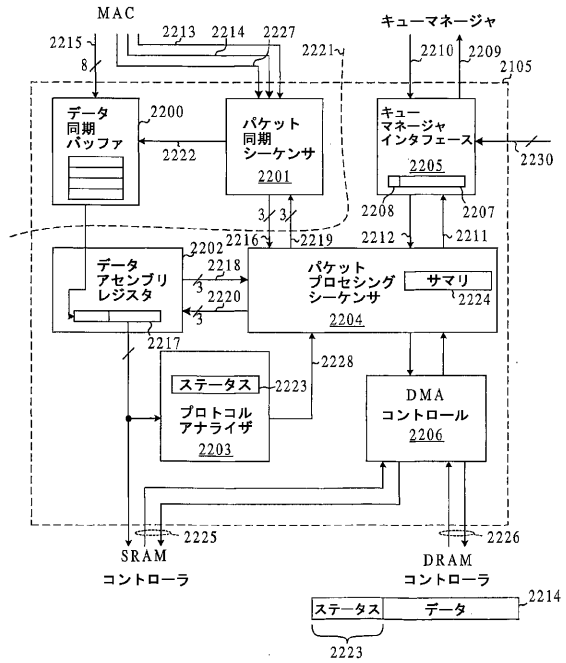




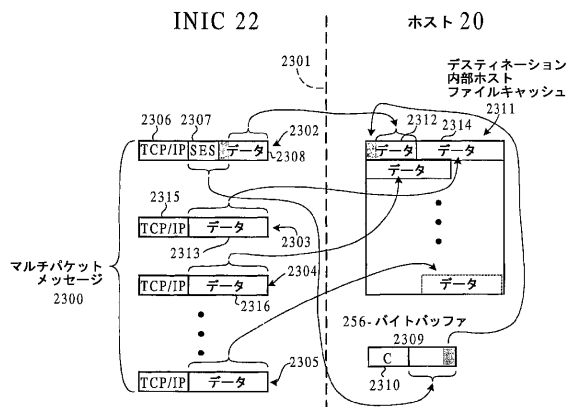
【図 24】



【図 25】



【図 26】



---

フロントページの続き

(31)優先権主張番号 09/802,551

(32)優先日 平成13年3月9日(2001.3.9)

(33)優先権主張国 米国(US)

(72)発明者 ブーシェ, ローレンス, ビー

アメリカ合衆国, カリフォルニア州・95070, サラトガ, モンタルボ・ハイツ・ドライブ・2  
0605

審査官 横山 佳弘

(56)参考文献 国際公開第2000/013091(WO, A1)

特開平11-045203(JP, A)

特開平09-128314(JP, A)

国際公開第00/004453(WO, A1)

特開平11-184639(JP, A)

特開2000-101537(JP, A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/38

G06F 12/00

G06F 13/10

G06F 13/00