

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4406604号  
(P4406604)

(45) 発行日 平成22年2月3日(2010.2.3)

(24) 登録日 平成21年11月13日(2009.11.13)

(51) Int.Cl. F I  
G O 6 F 13/10 (2006.01) G O 6 F 13/10 3 4 O B

請求項の数 10 (全 57 頁)

(21) 出願番号	特願2004-511951 (P2004-511951)	(73) 特許権者	504457223
(86) (22) 出願日	平成15年6月10日 (2003.6.10)		パンドヤ アシシュ エイ
(65) 公表番号	特表2006-516054 (P2006-516054A)		アメリカ合衆国 カリフォルニア州 95
(43) 公表日	平成18年6月15日 (2006.6.15)		762 エル ドラド ヒルズ ラファイ
(86) 国際出願番号	PCT/US2003/018386		エット ドライヴ 4318
(87) 国際公開番号	W02003/104943	(74) 代理人	100082005
(87) 国際公開日	平成15年12月18日 (2003.12.18)		弁理士 熊倉 禎男
審査請求日	平成18年6月7日 (2006.6.7)	(74) 代理人	100067013
(31) 優先権主張番号	60/388, 407		弁理士 大塚 文昭
(32) 優先日	平成14年6月11日 (2002.6.11)	(74) 代理人	100074228
(33) 優先権主張国	米国 (US)		弁理士 今城 俊夫
		(74) 代理人	100086771
			弁理士 西島 孝喜

最終頁に続く

(54) 【発明の名称】 TCP/IP、RDMA、及びIPストレージアプリケーションのための高性能IPプロセッサ

(57) 【特許請求の範囲】

【請求項1】

トランスミッションコントロールプロトコル(TCP)又はストリームコントロールトランスミッションプロトコル(SCTP)又はユーザデータグラムプロトコル(UDP)オペレーション、又はその任意の組合せをインターネットプロトコル(IP)上で提供するハードウェアプロセッサであって、前記オペレーションは、イニシエータから又はイニシエータへ、及びターゲットへ又はターゲットからネットワーク上でデータを転送するリモートダイレクトメモリアクセス(RDMA)機能を含み、スモールコンピュータシステムインタフェース(SCSI)コマンド層及びIPストレージドライバを有するホストプロセッサによって要求され、前記ハードウェアプロセッサは、

- a . RDMA 機構、
- b . プロセッサ内の作動のためのホストプロセッサのコマンド層からのコマンドをスケジュールに入れるためのコマンドスケジューラ、
- c . 前記コマンドスケジューラを経由する既存セッションのための前記ホストプロセッサからのスケジュールされたコマンドを待ち行列に入れるための第1のコマンド待ち行列、
- d . 前記コマンドスケジューラを経由する現在存在しないセッションのための前記ホストプロセッサからのスケジュールされたコマンドを待ち行列に入れるための第2のコマンド待ち行列、
- e . 前記コマンドが移送されたセッションの状態を記録するための、また、該コマンド

の R D M A を使用するものに対する R D M A の進行を記録するためのデータベース、

f . コマンド実行の状況処理のために前記 S C S I 層に通信するためのプロセッサと前記 S C S I 層の間の通信パス、及び

g . コマンドを解釈し、前記ハードウェアプロセッサと前記ターゲット及び前記イニシエータの少なくとも1つとの間でデータを転送するための R D M A オペレーションを実行するために協働するように連結した少なくとも1つの送信 / 受信エンジン及び少なくとも1つのコマンドエンジン、

を含むことを特徴とするプロセッサ。

【請求項2】

少なくとも1つのデータプロセッサを複数のインターネットプロトコル ( I P ) プロセッサにインタフェースで接続するための複数の I P プロセッサに連結した少なくとも1つのデータプロセッサを含むマルチプロセッサシステムであって、

前記複数の I P プロセッサの各々は、リモートダイレクトメモリアクセス ( R D M A ) プロトコルと、トランスミッションコントロールプロトコル ( T C P )、ストリームコントロールトランスミッションプロトコル ( S C T P )、及びユーザデータグラムプロトコル ( U D P ) のうちの1つのプロトコルとを実行するよう構成され、データを I P ネットワークのセット上に転送し、前記 I P プロセッサの1つは前記少なくとも1つのプロセッサからのコマンドを受信するよう構成され、 前記複数の I P プロセッサの1つは、

a . R D M A を介してデータを転送するための R D M A 機構、

b . 前記データを含む複数の I P パケットを処理するための少なくとも1つのパケットプロセッサ、

c . 前記データの転送に使用されるセッションに関する T C P / I P セッション情報を保存するためのセッションメモリ、

d . 前記データを保存するメモリに対するメモリアクセスを制御するための少なくとも1つのメモリコントローラ、

e . 前記 I P ネットワークの少なくとも1つに連結し、リモートホストからの要求を受信するための少なくとも1つのメディアインタフェース、及び

f . 前記少なくとも1つのデータプロセッサに連結するためのホストインタフェース、又は前記受信した要求に応答してファブリックを介してストレージシステムにアクセスするためのファブリックインタフェース、

g . 前記コマンドの状態を記録し、前記コマンドのターミネーションに基づいて前記コマンドの状態を更新するよう構成されたコマンド状態テーブル、

を含む、

ことを特徴とするシステム。

【請求項3】

前記 I P プロセッサの1つは、

a . I P ストレージセッション情報を保存するための I P ストレージセッションメモリ

b . I P パケットを分類するための分類プロセッサ、及び、前記少なくとも1つのパケットプロセッサによって使用される分類タグと前記 I P パケットとの対応付け、

c . 前記データのフローを制御するためのフローコントローラ、

d . ポリシーを適用するためのポリシープロセッサ、

e . 認証、暗号化、及び復号化のうちの1又は複数のセキュリティ作動を実行するためのセキュリティプロセッサ、

f . パケットを保存するためのパケットメモリ、

g . プレーン処理を制御するためのコントローラ、

h . パケットスケジューラ、

i . ピアプロセッサにインタフェース接続するための共用プロセッサインタフェース、 又は

j . これらの任意の組合せ、

10

20

30

40

50

のうちの少なくとも1つを含むことを特徴とする請求項2に記載のマルチプロセッサシステム。

【請求項4】

トランスミッションコントロールプロトコル(TCP)又はインターネットプロトコル(IP)ネットワーク上の他の接続指向プロトコルを使用しリモートダイレクトメモリアクセス(RDMA)を介してデータ転送が可能なハードウェアプロセッサを有するストレージエリアネットワークであって、

前記ハードウェアプロセッサはプログラム可能なように構成され、ポリシー処理エンジンを含み、前記ポリシー処理エンジンは、パケット基準毎、フロー基準毎、及びコマンド基準毎の1又は複数に基づくポリシーを適用するものであって、

前記ハードウェアプロセッサは、

ホストプロセッサから受信したデータ転送に対応するコマンドを実行し、前記受信したコマンドに基づいて前記ホストプロセッサのアドレスを決定し、前記コマンドの完了に際して前記決定されたアドレスにおいて前記ホストプロセッサのみに割り込みを行うようプログラム可能である

ことを特徴とするストレージエリアネットワーク。

【請求項5】

前記少なくとも1つの送信/受信エンジンは、分離された送信エンジン及び分離された受信エンジンとして実装されることを特徴とする、請求項1に記載のハードウェアプロセッサ。

【請求項6】

前記少なくとも1つの送信/受信エンジン及び少なくとも1つのコマンドエンジンは、送信、受信、及びコマンドエンジン機能を与える少なくとも1つの複合エンジンとして実装されることを特徴とする、請求項1に記載のハードウェアプロセッサ。

【請求項7】

前記ハードウェアプロセッサは、イニシエータからターゲットは、あるいはターゲットからイニシエータへのデータ転送についてのiSCSIコマンドをiSCSIセッションのいて受信し、前記セッションの各々の状態を前記データベースにセッションエントリとして記録することを特徴とする、請求項1に記載のハードウェアプロセッサ。

【請求項8】

前記複数のIPプロセッサの2以上のものが互いに結合されることを特徴とする、請求項2に記載のマルチプロセッサシステム。

【請求項9】

前記ハードウェアプロセッサは深いパケット分類エンジンを更に含むことを特徴とする、請求項4に記載のストレージエリアネットワーク。

【請求項10】

前記ハードウェアプロセッサは、サーバ、ホストバスアダプタ、スイッチ、スイッチラインカード、ゲートウェイ、ゲートウェイのラインカード、ストレージエリアネットワークアプライアンス、アプライアンスのラインカード、ストレージシステム、及びストレージシステムのラインカード、のうちの1又は複数において配置されることを特徴とする、請求項4に記載のストレージエリアネットワーク。

【発明の詳細な説明】

【技術分野】

【0001】

関連出願

2002年6月11日出願の米国特許仮出願出願番号第60/388,407号に対して優先権を請求する。2003年6月10日出願の「RDMAを使用する高性能IPプロセッサ」という名称の米国特許出願(出願番号未定)、2003年6月10日出願の「RDMAを使用するTCP/IPプロセッサ及びエンジン」という名称の米国特許出願(出願番号未定)、2003年6月10日出願の「RDMAを使用するIPストレージプロセ

10

20

30

40

50

ッサ及びそのためのエンジン」という名称の米国特許出願（出願番号未定）、2003年6月10日出願の「高性能IPプロセッサ用メモリシステム」という名称の米国特許出願（出願番号未定）、2003年6月10日出願の「インターネットプロトコル及びRDMAを使用するデータ処理システム」という名称の米国特許出願（出願番号未定）、2003年6月10日出願の「高性能IPプロセッサ」という名称の米国特許出願（出願番号未定）、及び2003年6月10日出願の「インターネットプロトコルを使用するデータ処理システム」という名称の米国特許出願（出願番号未定）は、上記の仮出願に関連する。

本発明は、一般的に、ストレージネットワークング半導体に関し、より詳細には、インターネットプロトコル（IP）ベースのストレージネットワークを作り出すために使用される高性能ネットワークストレージプロセッサに関する。

10

#### 【背景技術】

##### 【0002】

インターネットプロトコル（IP）は、ローカルエリアネットワーク（LAN）、メトロエリアネットワーク（MAN）、ワイドエリアネットワーク（WAN）などの様々なネットワークに亘って配備される最も普及しているネットワークングプロトコルである。ストレージエリアネットワーク（SAN）は、主にファイバチャネル（FC）技術に基づいている。IPベースのストレージネットワークを作り出す必要性が存在する。

##### 【0003】

データストリームを移送するように設計されたIP上でブロックストレージトラフィックを移送する時に、データストリームは、IP上で実行されるように階層化された伝送制御プロトコル（TCP）を使用して移送される。TCP/IPは、オペレーティングシステム内のソフトウェアで実行される確実な接続/セッション志向プロトコルである。TCP/IPソフトウェアスタックは、将来配備されることになる高回線速度を扱うには非常に低速である。現在、1Gbpsネットワーク接続を使用してTCP/IPスタックを実行する1GHzプロセッサベースサーバは、50～70%又はそれ以上のプロセッササイクルを使用し、プロセッサがサーバ上で実行されるアプリケーションに割り当てるのに利用可能な最小限のサイクルを残す。このオーバーヘッドは、高性能IPネットワークに対する場合と同じく、ストレージデータをTCP/IP上で移送する場合は許容できないことである。従って、新しいハードウェアソリューションは、ストレージ及びデータトラフィックを運び、FCベースソリューションと競合するようにTCP/IPスタックを加速するであろう。TCPプロトコルに加えて、SCTP及びUDPプロトコルのような他のプロトコル、及びデータストリームを移送するのに適切な他のプロトコルを使用することもできる。

20

30

##### 【0004】

【特許文献1】米国特許仮出願出願番号第60/388,407号

【特許文献2】2003年6月10日出願「RDMAを使用する高性能IPプロセッサ」米国特許出願（出願番号未定）

【特許文献3】2003年6月10日出願「RDMAを使用するTCP/IPプロセッサ及びエンジン」米国特許出願（出願番号未定）

【特許文献4】2003年6月10日出願「RDMAを使用するIPストレージプロセッサ及びそのためのエンジン」米国特許出願（出願番号未定）

40

【特許文献5】2003年6月10日出願「高性能IPプロセッサ用メモリシステム」米国特許出願（出願番号未定）

【特許文献6】2003年6月10日出願「インターネットプロトコル及びRDMAを使用するデータ処理システム」米国特許出願（出願番号未定）

【特許文献7】2003年6月10日出願「高性能IPプロセッサ」米国特許出願（出願番号未定）

【特許文献8】2003年6月10日出願「インターネットプロトコルを使用するデータ処理システム」米国特許出願（出願番号未定）

#### 【発明の開示】

50

## 【 0 0 0 5 】

本発明人は、ホストプロセッサからTCP/IPプロトコルスタックのオーバーヘッドを大幅に低減し、IPに基づく高回線速度のストレージ及びデータ移送ソリューションを可能にする高性能ハードウェアプロセッサについて説明する。

従来的に、TCP/IPネットワークスタックは、ソフトウェアスタックとしてオペレーティングシステムカーネルの内部で実行されている。ソフトウェアTCP/IPスタックの実行は、上述のように、1 G b p s ネットワークに情報提供する時に、1 G H z プロセッサで使用可能な処理サイクルの50%よりも多くを費やす。オーバーヘッドは、チェックサム計算、メモリバッファコピー、パケット到着時のプロセッサ割り込み、セッションの確立、セッションの取外し、及び他の確実な移送サービスを含むソフトウェアTCP/IPスタックの様々な局面から生じる。ソフトウェアスタックのオーバーヘッドは、回線速度がより高い時は法外に大きくなる。より性能が低いホストプロセッサを使用する無線ネットワークのような回線速度がより低いネットワークでも同様の問題が発生する。ハードウェアの実施によってホストプロセッサからオーバーヘッドを取り除くことができる。

10

## 【 0 0 0 6 】

オペレーティングシステムから提供されたソフトウェアTCP/IPネットワークスタックは、ホストプロセッササイクルの大部分を使い果たす。TCP/IPは、不安定なデータリンク上で実行することができる確実な移送手段である。従って、ネットワークパケットが欠落するか又はエラーを有する場合、TCPは、パケットの再送信を行う。パケットのエラーは、パケット内に有するチェックサムを使用して検出される。TCPパケットの受け取り側は、受け取ったパケットのチェックサムを実行し、それを受け取ったチェックサムと比較する。これは、パケット内の各受信バイトに関わる各パケット上で実行される高価な計算集約的作業である。発信元と宛先の間のパケットは、順序が狂って到着する場合があります。TCP層は、それをより上の層に提示する前にデータストリームの順序付けを行う。IPパケットはまた、リンク層の最大移送ユニット(MTU)に基づいてフラグメント化される場合があります。従って、受け取り側は、パケットを非フラグメント化することを期待される。これらの機能は、順序の狂ったパケット、フラグメント化したパケット、又は肯定応答されなかったパケットの例えばネットワークカード上のメモリへの一時的保存をもたらす。回線速度が1 G b p s を超えて増加すると、これらによって生じるメモリサイズオーバーヘッド及びメモリ速度の妨害により、ネットワークカードにかなりのコストが追加され、同じく莫大な性能的オーバーヘッドを引き起こす。多くのプロセッサのリソースを費やす別の機能は、ネットワークカードバッファ、カーネルバッファ、及びアプリケーションバッファへのノからのデータのコピーである。

20

30

## 【 0 0 0 7 】

マイクロプロセッサは、深いパイプライン化とスーパースカラーアーキテクチャを使用して、益々高性能と高速度を実現している。小さなパケットの到着時にこれらのプロセッサに割り込むと、コンテキスト切替オーバーヘッド、パイプラインフラッシュ、及びパイプラインの補充のために厳しい性能劣化を引き起こすことになる。従って、プロセッサの割り込みは、最も基本的な割り込みだけにして最小限にすべきである。ブロックストレージトラフィックがTCP/IPネットワーク上で移送される時、これらの性能の問題が重大になり、ストレージトラフィックの処理能力と待ち時間に大幅な影響を与える。従って、ストレージトラフィックの移送の処理全体におけるプロセッサの介入は、IPベースストレージソリューションが、ハードウェア実行を目的として特化されたファイバチャンネルのような他の特化されたネットワークアーキテクチャと同等な性能と待ち時間を有するために最小にする必要がある。i S C S I、F C I P、i F C P、及び他のもの(N F S、C I F S、D A F S、H T T P、X M L、X M L 派生物(音声X M L、E B X M L、マイクロソフトS O A P、及びその他など)、S G M L、及びH T M L フォーマットのような)のような出現してきたIPベースストレージ基準は、TCP/IPセグメント内のストレージ及びデータトラフィックをカプセル化する。しかし、通常は、TCPセグメント

40

50

とTCPパケットによってカプセル化されたプロトコルデータユニットとの間には、アラインメント関係が存在しない。これは、今日のネットワークでは非常に頻繁なイベントであるパケットの順序が狂って到着した時に問題となる。ストレージ及びデータブロックは、ストリーム内の中間パケットが到着するまで、使用のための抽出を順序が狂ったパケットから行うことができず、この中間パケットは、ネットワークアダプタにこれらのパケットをメモリに保存させ、中間パケットの到着時にそれらを検索させ、それらを順序化させることになる。これは、要求されるメモリストレージのサイズから、及びメモリサブシステムがサポートすることを期待される性能からも、特に回線速度が1 Gbpsを超える時に高価になる可能性がある。各TCPセグメントがプロトコルデータユニット及びそのシーケンスを一意的に識別することができる場合、このオーバーヘッドを排除することができる。これは、パケットをホストシステム内のそれらの末端メモリ位置に直接移送することを可能にすることができる。ホストプロセッサの介入はまた、ストレージサブシステムに転送される場合があるか又はクラスター化した環境又は他のクライアントサーバ環境で他のプロセッサと共有される大きなデータブロックの転送では最小限にされるべきである。プロセッサは、影響を最小限にするためにストレージコマンド境界上でのみ割り込まれるべきである。

#### 【0008】

本明細書に示すIPプロセッサは、革新的な構造的特徴と設計により、以上概説した様々な問題の影響を除去するか又は著しく低減する。説明したプロセッサのアーキテクチャは、ストレージ及びデータペイロードを運ぶTCPトラフィックを終了させる機能を提供し、それによってホストプロセッサ上のTCP/IPネットワークスタックのオーバーヘッドを除去するか又は大幅に低減し、パケットを入力から出力まで最小の待ち時間で通過させるパケットストリームアーキテクチャをもたらす。高回線速度のストレージ又はデータトラフィックがIP上で搬送されることを可能にするために、従来はホストカーネル又はドライバソフトウェアによって維持されていた様々な接続(セッション)に対する伝送制御ブロック情報を維持することが必要になる。本特許で使用されるように、用語「IPセッション」は、IP上で実行されるセッション志向プロトコルに対するセッションを意味する。例としては、TCP/IP及びSCTP/IPなどがある。各パケットに対するセッション情報へのアクセスは、大幅な処理オーバーヘッドを付加する。説明したアーキテクチャは、このオーバーヘッドを大幅に低減する高性能メモリサブシステムを作り出す。このプロセッサのアーキテクチャは、主としてコマンド又はデータ転送終了境界でのホストプロセッサに対する割り込みを最小限にする知的フロー制御の機能を提供する。

#### 【0009】

今日、セキュリティを有するTCP/IPプロセッサは提供されていない。

説明したプロセッサのアーキテクチャはまた、一体化されたセキュリティ機能を提供する。ストレージトラフィックが、ネットワーク上でサーバからSAN又は他のストレージシステム内のストレージレイまで搬送される時、それは、直接取付けられたストレージシステムが処理する必要のない様々なセキュリティの脆弱性に露出される。このプロセッサは、ストリーム内でストレージトラフィックの暗号化と暗号解読を可能にするので、高回線速度を可能にし、同時にストレージデータトラフィックの機密性を提供する。

#### 【0010】

ネットワークトラフィックの分類は、深いパケットの検査及び処理のための僅かなサイクルを残してパケットプロセッサ上で利用可能な処理サイクルの半分までを消費する別のタスクである。IPベースのストレージトラフィックは、プロトコルの性質上、高速で待ち時間が少ない深いパケット処理を要求する。説明したIPプロセッサは、プログラマブル分類エンジンを提供することにより、分類のオーバーヘッドを大幅に低減する。

#### 【0011】

ストレージ容量とストレージネットワークにおける大きな成長は、IT部門の主なコストアイテムとしてストレージ領域管理を作り出している。管理コストを抑えるためにポリシーベースのストレージ管理が要求される。説明したプログラマブル分類エンジンは、パ

10

20

30

40

50

ケット、トランザクション、フロー、及びコマンド境界に対して強制することができるストレージポリシーの配備を可能にする。これは、ストレージ領域の管理コストを大幅に改善することになる。

#### 【0012】

プログラマブルIPプロセッサアーキテクチャはまた、顧客特定アプリケーションの配備を可能にするのに十分なヘッドルームを提供する。これらのアプリケーションは、例えば、ネットワーク管理、ストレージファイアウォール又は他のセキュリティ機能、帯域幅管理、サービス品質、仮想化、性能モニタリング、ゾーニング、及びLUNマスキングなどの複数のカテゴリに属するであろう。

#### 【発明を実施するための最良の形態】

#### 【0013】

本発明人は、ソフトウェアTCP/IPスタックの厳しい性能の影響からホストプロセッサを解放するために、ハードウェアにTCP/IPスタックを実施する新しい高性能で待ち時間が短い方法を提供する。このハードウェアTCP/IPスタックは、次に、付加的な処理要素とインタフェースで接続され、高性能で待ち時間が短いIPベースのストレージアプリケーションを使用可能にする。

これは、TCP/IPの終了、高性能で待ち時間が短いIPストレージ機能、遠隔DMA(RDMA)機能、セキュリティ機能、及びプログラマブル分類及びポリシー処理機能などの恩典をもたらすために、様々な形式で実行することができる。以下は、これを実行することができるいくつかの実施形態である。

#### 【0014】

##### サーバ

説明するアーキテクチャは、TCP/IPソフトウェアの1つ又は複数のホストサーバプロセッサ及び性能オーバーヘッドを緩和するハードウェアベースのTCP/IP機能を提供する高性能サーバ環境に組み込むことができる。IPプロセッサは、ハードウェアTCP/IPとの高性能ネットワーキングインタフェースを提供する、サーバチップセットに対する対照プロセッサとすることができる。サーバは、ブレードサーバ、アプライアンスサーバ、ファイルサーバ、シンサーバ、クラスター化サーバ、データベースサーバ、ゲームサーバ、グリッド計算サーバ、VOIPサーバ、無線ゲートウェイサーバ、セキュリティサーバ、ネットワーク取付ストレージサーバ、又は従来のサーバのような様々な形状因子にすることができる。本実施形態は、サーバマザーボード上に高性能ネットワークインタフェースの作成を可能にすると考えられる。

#### 【0015】

##### サーバチップセットに対する対照プロセッサ

サーバ環境はまた、高性能TCP/IP及び/又はRDMA機能に加えて、説明したプロセッサの高性能IPストレージ処理性能に影響を及ぼすことがある。そのような実施形態では、プロセッサは、サーバプロセッサからのTCP/IPオフロードに加えて高性能ネットワークストレージI/O性能を提供する、サーバチップセットに対する対照プロセッサとすることができる。この実施形態は、マザーボード上に高性能IPベースのネットワークストレージI/Oの作成を可能にすると考えられる。換言すれば、それは、マザーボード上の「IP SAN」を使用可能にするであろう。

#### 【0016】

##### ストレージシステムチップセット

プロセッサはまた、ストレージシステム内のチップセットの対照として使用することができる。これは、ストレージネットワーキング環境でストレージデータサーバの機能を実行するストレージアレイ(又は、他の適切なストレージシステム又はサブシステム)コントローラとすることができる。IPベースSANにおいてネットワーキングするために、プロセッサは、ストレージアレイコントローラに対してIPネットワークストレージ機能を提供すると考えられる。その構成は、ストレージアレイにアクセスして別のストレージ中枢機能を提供するために、システム内に付加的な機能を持ったサーバ環境の構成と同等に

10

20

30

40

50

することができる。

【0017】

サーバ/ストレージホストアダプタカード

IPプロセッサはまた、高速度TCP/IPネットワークングを提供するサーバホストアダプタカードに組み込むことができる。同じアダプタカードはまた、IPベースストレージネットワークに対して、高速度ネットワークストレージ機能を提供することができる。アダプタカードは、従来のサーバで使用することができ、同様にブレードサーバ構成においてブレードとして使用することもできる。プロセッサはまた、IPベースストレージネットワークング機能を提供するストレージアレイ（又は、別のストレージシステム又はサブシステム）のフロントエンド内のアダプタで使用することができる。

10

【0018】

プロセッサチップセット構成要素

TCP/IPプロセッサは、TCP/IPオフロード機能を提供するプロセッサチップセットの内部に組み込むことができる。そのような構成は、ハイエンドサーバ、ワークステーション、又は高速ネットワークとインタフェースで接続する高性能パーソナルコンピュータに使用することができる。そのような実施形態はまた、IPストレージ又はRDMA機能、又は本発明の組合せを含むことができ、チップセットに組み込まれたIPベースストレージネットワークング、及び/又はRDMA機能を有するTCP/IPを提供する。説明したアーキテクチャの複数の機能の使用は、機能要件、開発時間とコスト、及びシリコンダイのコストなどのトレードオフとして、上記又は他の実施形態の他の機能の使用と独立して行うことができる。

20

【0019】

ストレージ又はSANシステム又はサブシステム切替回線カード

IPプロセッサはまた、高性能で待ち時間が少ないSAN切替システム（又は、別のシステム又はサブシステム）回線カードを作り出すために使用することができる。プロセッサは、回線カードへの/からのIPベースストレージトラフィックを終了させたり開始させたりするメインプロセッサとして使用することができる。このプロセッサは、ホストのように作動することができる切替システムのファブリックコントローラを使用して作動し、切替システム内に存在する転送情報ベースによって判断された適切な切替回線カードに対するIP宛先に基づいて、終了したストレージトラフィックを移送する。そのような切替システムは、純粋にIPベースネットワークングをサポートすることができ、又は、マルチプロトコルサポートをサポートすることもでき、ファイバチャンネルのような他のデータセンターのSANファブリックに加えて、IPベースSANとインタフェースで接続することを可能にする。ゲートウェイコントローラシステムの内には、LAN又はWANからIPストレージトラフィックを終了させ、ストレージトラフィックをSANに搬送するために新しいセッションを開始させる非常に類似の構成が存在することができ、上述のSANは、IPベースSANか、又は、より可能性があるのは、ファイバチャンネルのようなデータセンター内部の他のファブリックから構成されたSANである。プロセッサはまた、SANゲートウェイコントローラに組み込むことができる。

30

【0020】

ストレージアプライアンス

ストレージネットワークの管理コストは急速に増加している。ネットワーク及びストレージ容量の大幅な成長を管理する能力は、ストレージ領域の管理機能をもたらす特別なアプライアンスを作り出すことを要求する。高性能IPベースSANに対して説明した管理アプライアンスは、その機能を達成することができるように、TCP/IPパケット内部で移送されたストレージトラフィック上で高性能IPプロセッサを実行する。これらのシステムは、ポリシーベースの管理及び実施機能を提供するために、深いパケットの検査を実行してIPトラフィック内のストレージペイロードを抽出するための高性能プロセッサを要求する。セキュリティを有するプログラブル分類及びポリシーエンジンは、説明した高速TCP/IP及びIPストレージエンジンと共に、本特許に説明するこれらのアプ

40

50



ライアンスと他の実施形態が、高回線速度かつ短い待ち時間で深いパケット検査及び分類を実行し、パケット毎のベースで必要なポリシーを適用することを可能にする。更に、これらの機能は、仮想化、ポリシーベース管理、セキュリティ実施、アクセス制御、侵入検出、帯域幅管理、トラフィックシェーピング、サービス品質、迷惑メール防止、ウイルス検出、暗号化、暗号解読、LUNマスキング、ゾーニング、リンク集約、及びストレージ領域ネットワークトラフィックに対してインバンドである類似のもののような、その機能を実行することができるストレージ管理アプライアンスの作成を可能にする。同様のポリシーベースの管理、及びセキュリティ作動又は機能性はまた、本特許に説明する他の実施形態の内部でサポートすることができる。

#### 【 0 0 2 1 】

##### クラスター化された環境

システム性能と、クラスター化データベースなどのようなアプリケーションに対する拡張容易性とを増強するために、クラスター化環境ではサーバシステムが使用される。高性能クラスターサーバ上で作動しているアプリケーションは、プロセス間通信に対して高速でデータを共有する能力を要求する。従来のソフトウェアTCP/IPネットワーク上のクラスタープロセッサ間でのこのプロセス間通信トラフィックの移送は、厳しい性能のオーバーヘッドを受ける。従って、ファイバチャンネルのような専用ファブリックがそのような構成では使用されてきた。しかし、通信プロセスメモリ間で直接のメモリアクセスを可能にするTCP/IPベースのファブリックは、ファイバチャンネルのような専用ファブリックに変換することなく、任意のTCP/IPネットワーク上で作動するアプリケーションによって使用される。説明したIPプロセッサは、その高性能TCP/IP処理機能とRDMA機能を使用してクラスターサーバ環境に組み込むことが可能であり、高性能で短い待ち時間の直接メモリという恩典をメモリデータ転送にもたらす。この実施形態はまた、グローバルなクラスター化を作り出すために使用することができ、同様にグリッドコンピュータ及びグリッドネットワーク内のデータ転送を可能にするために使用することができる。

#### 【 0 0 2 2 】

##### 付加的な実施形態

プロセッサアーキテクチャは、部分的にソフトウェアで、また部分的にハードウェアで実施することができる。性能の必要性和コストの推測は、本発明のシステムアーキテクチャ全体のハードウェア及びソフトウェア分割に対するトレードオフを押し進める可能性がある。このアーキテクチャをハードウェア及びソフトウェア分割と共にチップセットの組み合わせとして、又は分割とは独立して実行することも可能である。例えば、セキュリティプロセッサと分類エンジンは、別々のチップ上にあり、同様の機能を提供することができる。これは、開発及び製造コストを含めて、シリコンコストの廉価なIPプロセッサをもたらすことができるが、いくつかの例では、システムも含めた要素を増加させ、設置面積とソリューション全体のコストを増大させる可能性がある。セキュリティ及び分類エンジンも別々のチップとすることができるであろう。本明細書で使用されるように、チップセットは、複数チップのチップセット、又は単一チップだけを含むチップセットを意味することができる、それはアプリケーションに依存する。

#### 【 0 0 2 3 】

ストレージフローコントローラと待ち行列は、ホスト上のソフトウェア内に維持することができる、又はチップセット内の別のチップの一部になることもできる。従って、高性能プロセッサの将来の出現に伴って要求される高性能IPベースストレージとTCP/IPオフロードアプリケーションとを達成するために、このアーキテクチャを分割する複数の方法が実現可能である。iSCSIと関連してストレージエンジンが説明されたが、制御プロセッサと共に、TCP/IP及びストレージエンジンのプログラム可能性、分類子のプログラム可能性、及びストレージフローコントローラを使用して、iFCP、FCIP、及びその他のような他のIPストレージプロトコルを適切なファームウェアで実施することができる。iSCSI作動はまた、IPストレージ作動とすることができる。高性能

10

20

30

40

50

IPプロセッサコアは、より低い回線速度の複数の入力出力ポートと連結し、全処理量を適合させてマルチポートIPプロセッサ実施形態を同様に作り出すことができる。

【0024】

ストレージエンジンを使用しないでメインプロセッサからオフロードする高性能TCP/IPに対してこのアーキテクチャを使用することが実現可能である。これは、データ及び遠距離通信アプリケーションの次世代の高性能ネットワークに対してシリコンとシステムのソリューションをもたらすことができる。TCP/IPエンジンは、アプリケーション特定の packets 加速器を使用して増強し、コアのアーキテクチャを利用してこのプロセッサの新しい味を引き出すことができる。本発明の他の機能に加えて、ストレージエンジンをファイアウォールエンジン、又はルート検索エンジン、又は遠距離通信/ネットワーク加速エンジンのような別のアプリケーション特定の加速器と取り替え、このプロセッサアーキテクチャのターゲットを遠距離通信/ネットワークング及び他のアプリケーションにすることができる。

10

【0025】

詳細な説明

ストレージのコストと需要は、過去数年間に亘って速いペースで増加してきた。これは、近い将来に亘って同じ割合で成長することが見込まれる。eビジネスの出現に伴って、サーバ又はシステムの機能停止時間に関係なく、任意の時間の任意の場所でのデータ使用可能度が重要になっている。これは、大きな必要性を駆り立て、サーバ接続ストレージをネットワークに移動させ、ストレージの統合、データの利用可能性、及びデータ管理の容易性をもたらしている。ストレージエリアネットワーク(SAN)は、現在では大部分がファイバチャンネル技術に基づいており、TCP/IP技術と比べてハードウェア志向のスタックを使用した短い待ち時間と高性能というような様々な恩典をもたらしている。

20

【0026】

一部のシステムは、データストリームを移送するように設計されたIP上でブロックストレージトラフィックを移送する。データストリームは、IPの上で実行されるように階層化された伝送制御プロトコル(TCP)を使用して移送される。TCPは、オペレーティングシステム内のソフトウェアで実施される確実な接続志向プロトコルである。TCP/IPソフトウェアスタックは、将来配置されるであろう高回線速度を処理するには低速である。新しいハードウェアのソリューションは、ストレージとネットワークトラフィックを搬送するためにTCP/IPスタックを加速し、FCベースのソリューションに対して競争力があることになる。

30

【0027】

高性能サーバ、ワークステーション、ストレージコントローラ、及びネットワークトラフィックで普及しているストレージプロトコルは、約20年前から存在するSCSIプロトコルである。SCSIアーキテクチャは、階層化プロトコルアーキテクチャとして構築されている。図1は、イニシエータ(ブロック101)とターゲットサブシステム(ブロック102)内の様々なSCSIアーキテクチャ階層を示す。本特許で使用されるように、用語「イニシエータ」及び「ターゲット」は、データ処理装置又はそれらを含むサブシステム又はシステムを意味する。用語「イニシエータ」及び「ターゲット」はまた、クライアント又はサーバ又はピアを意味することができる。同様に、用語「ピア」は、対等の処理装置又はそれらのサブシステム又はシステムを意味することができる。「遠隔ピア」は、世界中に亘って又は部屋中に亘って配置されたピアとすることができる。

40

【0028】

図1のイニシエータとターゲットサブシステムは、クライアントサーバ要求と応答処理を提供するために使用されるSCSIアプリケーションプロトコル層(ブロック103)を使用して互いに対話する。それは、ディスクアレイやテープドライブなどのような多くの形態をとることができるイニシエータ及びターゲットの大容量ストレージ間で装置サービス要求と応答を同様に提供する。従来、ターゲットとイニシエータは、SCSIプロトコル(ブロック104)を搬送するSCSIバスアーキテクチャを使用して相互接続され

50

た。SCSIプロトコル層は、SCSIアプリケーションプロトコルを使用して、クライアントとサーバの対話を互いに可能にする移送層である。移送層は、上部層のプロトコルとアプリケーションが移送プロトコルを独立させておくことができるように、上部層に対して同じ意味を呈するべきである。

【0029】

図2は、ベースとなる移送層に加えてSCSIアプリケーション層を示す。IETF基準の追跡プロトコル、すなわち、iSCSI(IP上のSCSI)は、IPベースのストレージ移送プロトコルを提供するための試みである。FCIP(IPにカプセル化されたFC)、iFCIP(IP上のFC)などを含む別の同様の試みがある。移送機構としてのTCP/IPに加えて、これらのプロトコルの多くは、図2に示す方法と同様の方法で階層化される。図2に示すように、iSCSIプロトコルサービス層(ブロック204)は、SCSIアプリケーション層(ブロック203)に対して層状インタフェースを形成する。iSCSIは、基準によって定められるように、iSCSIプロトコルデータユニット(PDU)としてSCSIコマンド及びデータを搬送する。これらのプロトコルデータユニットは、次に、TCP/IP(ブロック205)などを使用してネットワーク上で移送することができる。この基準は、「iSCSI PDU」を運ぶ基本的な移送を実施するための手段を特定しない。図2は、「iSCSI PDU」に対して移送を提供するTCP/IP上で階層化されたiSCSIを示す。

【0030】

iSCSIのようなIPベースのストレージプロトコルは、ソフトウェアベースのTCP/IPスタックに加えて、ソフトウェア内で階層化することができる。しかし、そのような実施は、それに加えてソフトウェアTCP/IPと階層化されたストレージプロトコルから発生する厳しい性能ペナルティに苦しむことになる。それらの実施は、ホストプロセッサの性能に厳しく影響し、1Gbpsを超える回線速度での他の全てのタスクに対してプロセッサを使用不能にするであろう。従って、ハードウェア内でTCP/IPスタックを実施し、ホストプロセッサを解放し、その上にストレージプロトコルを構築することができるであろう。iSCSIと同様にストレージプロトコルは、ホストプロセッサ上で実施されるソフトウェア内に構築することができるか、又は本特許に説明するようにハードウェアの実施によって加速することができる。ソフトウェアiSCSIスタックは、受信TCPセグメントからPDUを抽出してそれらに作用することができるように、ホストプロセッサに対して多くの割り込みを呈することになる。そのような実施は、ソフトウェアベースTCPスタックが実施したであろう同様の理由で、厳しい性能ペナルティに苦しむことになる。説明したプロセッサは、ホストプロセッサに対する性能ペナルティと得られる待ち時間の影響とを排除又は大幅に低減する、TCP/IPベースネットワーク上でストレージプロトコルを移送するための高性能で待ち時間の短いアーキテクチャを提供する。

【0031】

図3は、OSIネットワークスタックとして参照されたファイバチャンネルに対するTCP/IPスタックの比較を示す。TCP/IPスタック(ブロック303)には、本特許の「発明の開示」の節で上述したように、ホスト上のソフトウェアの実施に起因する性能上の問題がある。それに比べて、ファイバチャンネル(ブロック304)のような又は他の専用ネットワークプロトコルは、ハードウェア内で実施されるように設計される。ハードウェアでの実施により、ネットワークソリューションがIPベースソリューションよりも高性能になる。しかし、IPのコピキタス性とITユーザ及び開発者の観点からのIPの習熟性とは、広範囲の開発に対してIPを更に適切なものにする。これは、TCP/IPに起因する性能ペナルティが、他の競争力のある専用プロトコルに相当するまで減少した場合に達成することができる。図4は、図示の他の専用プロトコルと競合することができるようになるために、TCP/IP(ブロック403)に使用されるハードウェア及びソフトウェア内で階層化するプロトコルのレベルを示す。

【0032】

図5は、ハードウェアベースのTCP/IPと、本特許のストレージプロトコルの実行を使用するホストオペレーティングシステムスタックとを示す。プロトコルは、それがホストオペレーティングシステムスタック(ブロック513)内にオペレーティングシステム層がその上で不変のまま導入されようように実施される。これは、SCSIアプリケーションプロトコルを一切変えないで作動させる。ドライバ層(ブロック515)と、IPベースのストレージインタフェースのための下のスタック(ブロック501)とは、非ネットワークSCSIインタフェース(ブロック506及び503)又はファイバチャンネルインタフェース(ブロック502)と同等のインタフェースを表すことになる。

【0033】

図6は、ソフトウェアTCP/IPスタックに包含されたデータ転送を示す。TCP/IPスタックのそのような実施は、データ転送のメモリコピーからの非常に大きな性能ペナルティを有する。図は、クライアントとサーバネットワークスタックの間のデータ転送を表す。クライアントからユーザまで又はその反対に移送する必要があるユーザレベルのアプリケーションバッファ(ブロック601)は、図示の様々なレベルのデータ転送を通過する。発信元上のユーザアプリケーションバッファは、OSカーネルのスペースバッファ(ブロック602)にコピーされる。このデータは、次に、ネットワークドライバのバッファ(ブロック603)にコピーされ、そこからネットワークインタフェースカード(NIC)又はホストバスアダプタ(HBA)にDMA転送される。バッファコピー作動は、ホストプロセッサを伴い、大事なプロセッササイクルを使い果たす。更に、転送中のデータは、ホストからの付加的な計算サイクルを使い果たし、ホスト上のチェックサム計算を通過して進行する。ホスト上のシステムメモリへの又はからのデータの複数回の移動も、同じくメモリ帯域幅の妨害を作り出す。NIC/HBAに移送されたデータは、次にネットワーク上に送られ(ブロック609)、宛先システムに到着する。宛先システムでは、データパケットは、ホストと反対方向であるが同様のバッファコピーとチェックサム作動に従って、ソフトウェアネットワークスタックを通過してトラバースする。TCP/IPスタックのそのような実施は、ブロックストレージデータの移送に対して、及び発信元と宛先の間で大量のデータが移送されるような場合のクラスター化アプリケーションに対しては非効率である。

【0034】

図7は、本特許に説明するアーキテクチャの遠隔の直接メモリアクセス(RDMA)機能を可能にするという特徴を有するイニシエータ及びターゲットでのネットワークスタックを示す。以下は、RDMA能力又はRDMA機構又はRDMA機能と呼ぶことができる。そのようなシステムでは、イニシエータ又はターゲット上で実行されるアプリケーションは、メモリ領域(ブロック702)を登録し、この領域は、実質的にホストの介入なしでそのピアにNIC/HBAから直接利用可能にされる。これらのアプリケーションはまた、それらのピアにRDMA(ブロック708)に関して利用可能であるメモリ領域について知らせるのである。通信の両方のピアがRDMA機構を使用する準備ができた状態で、RDMA領域からのデータの移送は、ピア内のNIC/HBAハードウェアがRDMA機能を実行すれば、実質的なホストの介入がなくとも発信元から宛先まで基本的にゼロコピーのオーバーヘッドで起こり得る。発信元、すなわちイニシエータは、その要求のそのピアに特定RDMA使用可能バッファを読み取るか又は書き込むように通知し、次に、宛先、すなわちターゲットにそのRDMAバッファへ/からデータをプッシュ又はプルさせる。イニシエータ及びターゲットNIC/HBAは、次に、本特許に説明するTCP/IPハードウェア実施、RDMA703、TCP/IPオフロード704、RDMA708、及びオフロード709を使用して、ホストプロセッサの実質的な記入なしで互いの間でデータを移送し、それによってプロセッサのオーバーヘッドを大幅に低減するのである。この機構は、ホストプロセッサ上のTCP/IP処理のオーバーヘッドを大幅に低減し、図6に示すデータ転送のための複数バッファのコピーの必要性を排除すると考えられる。RDMA実行可能システムは、従って、高速又は低速の如何に関わらず、システムがそのピアの機能の妨害を作成せずにデータ転送を実行することを可能にする。IPソリュ

10

20

30

40

50

ーションの上でストレージ内のこのプロセッサに実施されたRDM A機能は、通常はデータ転送の開始時と終了時を除いて、ホストの介入を排除する。これは、ターゲット及びイニシエータの両方のシステムでホストプロセッサを解放し、各パケット到着又は転送時に介入されることなく有用なタスクを実行する。RDM Aの実行はまた、システムを安全にし、認可されていないアクセスを阻止することを可能にする。これは、領域のIDと共にアクセス制御キーを有するHBA/NICを使用してエクスポートされたメモリを登録することによって達成される。HBA/NICは、遠隔のホストからRDM Aバッファへのメモリ領域要求のアドレス変換を実行し、セキュリティキー検証のような安全作動を実行し、次にデータを移送させる。この処理は、HBA/NICに常駐するか、又は、例えばマザーボード上のホストプロセッサに対するコンパニオンプロセッサとして、本発明のプロセッサのホストプロセッサから切り離されて実行される。この機能はまた、クライアントサーバアプリケーションと同様に、厳しいクラスター化アプリケーションに対する大量のデータ転送に使用することができる。発信元すなわちイニシエータと宛先すなわちターゲットの間で大量のデータを移送するリアルタイムのメディアアプリケーションは、これによって恩典を得ることができる。

10

## 【0035】

図8は、ソフトウェアで実行されたホストファイルシステムとSCSIスタックとを示す。上述の通り、IPベースのストレージスタック(ブロック805、806、807、808、及び809)は、SCSI層(ブロック803及び804)に対して、SCSI移送層(ブロック811)又はファイバチャンネル移送(ブロック810)によって提供されたような一貫したインタフェースを呈示すべきである。この図は、性能に影響されやすいブロックデータを移送するように設定されていないIPの様々な問題によって課せられるもの以外に、システムレベルからIPベースのストレージの実行に課せられる高レベルの要求を示す。

20

## 【0036】

図9は、図8に示すiSCSIスタックをより詳細に示す。iSCSIスタック(ブロック805から809)は、SCSIコマンド統合層(ブロック803及び804)に対して、この層とその上の他の層の挙動が変わらないように、OSで定められたドライバインタフェースレベルの機能性を提供すべきである。図9は、IPストレージ機能を提供するために実施されると考えられる一組の機能を示す。iSCSIの機能性を提供する機能は、関連する組の機能に分類されるが、どの当業者も認めるように、これらの多くの変形が存在する可能性がある。標準的(例えば、ターゲット及びイニシエータのログイン及びログアウト)機能(ブロック916)や接続の確立及び取り外し機能(ブロック905)を満足するように要求される一組の機能がある。図は、「OS SCSI」ソフトウェアスタックがiSCSI装置(ブロック916)を見つけ、オプション/パラメータ(ブロック903及び909)を設定して取得し、装置を始動して(ブロック913)装置を解放する(ブロック911)ことを可能にする機能を示す。上述の制御機能以外に、iSCSIの実施は、待ち行列(912及び917)を通じて大量のデータの移送機能を提供し、iSCSI基準によって規定されたPDUを移送する。iSCSIスタックはまた、直接データ転送/配置(DDT)又はRDM A機能又はその組合せを含むことができ(ブロック918)、これらは、イニシエータ及びターゲットシステムによって使用することができ、ストレージ及び他の大量のブロックデータ転送を含めて、実質的にゼロバッファのコピーとホスト介入のないデータ転送とを実行する。SCSIコマンド及びこれらに関連するブロックデータ転送は、説明したプロセッサ上で実行されるコマンド待ち行列(ブロック912及び917)として実施される。ホストは、主にコマンドの終了時に割り込まれる。終了したコマンドは、ホストの都合のよい時に作用するように、ホストに対する待ち行列に入れられる。図は、iSCSIプロトコル層と、同じくプロセッサから切り離されて実行される本明細書に説明したIPプロセッサシステム上のTCP/IPスタック(ブロック907及び908)上で階層化されたドライバ層とを示す。

30

40

## 【0037】

50

図10は、説明したIPプロセッサシステムで実施されるTCP/IPスタック機能性を示す。これらの機能は、このプロセッサの高性能TCP/IP実施を利用するために、上部層プロトコル機能に対するインタフェースを提供し、直接「OS TCP/IP」バイパス、RDMA又はネットワークソケット直接機能、又はその組合せから恩典を得ることができる他のアプリケーションと共にIPストレージトラフィックを搬送する。TCP/IPスタックは、上部層データ(ブロック1017及び1031)とコマンドPDUを送受信し、移送接続と取外し機能(ブロック1021)を確立し、データ転送機能やチェックサム機能(ブロック1019)、並びにエラー処理機能(ブロック1022)やセグメント化及び順序付け及びウィンドウ操作作動(ブロック1013)を送受信するための機能を提供する。チェックサム確認/作成のような特定の機能は、データ転送の全てのバイトに関連するのに対して、データパケットを移送して伝送制御ブロック又はセッションデータベースを更新するいくつかの機能は、データ転送の各パケットのために呼び出される。セッションDB(ブロック1025)は、TCP/IP状態情報と共に、アクティブセッション/接続に関する様々な情報を維持するために使用される。TCP層は、基準による要求に応じてIP機能を提供するIP層に加えて構築される。この層は、パズMTUによってパケットをフラグメント化/フラグメント解消する機能(ブロック1033)を提供し、例えば、ICMP(ブロック1029)のようなエラーを通信するために必要な他の機能とのインタフェースと同様に、ルート及び転送情報(ブロック1032)を提供する。IP層は、イーサネット層又は他のメディアアクセス層技術とインタフェースで接続し、ネットワーク上にTCP/IPパケットを移送する。この説明の様々な図には、より低い層がイーサネットとして示されているが、例えば、MAN/WAN上のSONETに亘ってパケットを移送するためのSONETのような別の技術とすることができる。イーサネットはまた、類似のアプリケーションで使用されるが、例えば、LAN及び専用ローカルSAN環境内で更に多く使用することができる。

#### 【0038】

図11は、iSCSIデータフローを示す。図は、データフローの受信及び送信パスを示す。ホストのSCSIコマンド層は、iSCSIDライバと協働し、両方ともブロック1101に示されているが、図26により詳細に見られるストレージフローコントローラ内のコマンドスケジューラ(ブロック1108)に対して処理されるようにコマンドをスケジュールに入れるであろう。コマンドスケジューラ1108は、図17により詳細に示すプロセッサ内で作動するために、新しいコマンドをスケジュールに入れる。既存の接続を用いてターゲット装置と定められた新しいコマンドは、その既存の接続に対する待ち行列に入れられる(ブロック1111)。ターゲット装置に対する接続が存在しない場合は、新しいコマンドは、割り当てられていないコマンド待ち行列に対してエンキューされる(ブロック1102)。図47とブロック905及び1006に示すようなセッション/接続の確立処理が次に呼び出され、ターゲットに接続される。接続が確立された状態で、待ち行列1102からの対応するコマンドは、図に示すように、コマンドスケジューラ1108により新しく作り出された接続コマンド待ち行列1111に対してエンキューされる。コマンドが実行ステージに達した状態で、コマンドが読取又は書込トランザクション次第で、受信1107又は送信1109パスがアクティブにされる。コマンドが移送される接続/セッションの状態は、次に説明するように、セッションデータベース内のコマンド実行の進行を記録するために使用される。データ転送に関連するバッファは、転送が終了する時までロックすることができる。イニシエータとターゲットの間でデータを送信するためにRDMA機構が使用される場合、適切な領域バッファ識別子、アクセス制御キー、及び関連RDMA状態データは、選択された実施次第でプロセッサ上のメモリに維持され、同じくオフチップメモリにも維持することができる。複数のTCPセグメントに亘ることがあるコマンドに関連するデータ転送が終了すると、コマンド実行のステータスは、次に適切な処理を行うホストのSCSI層に転送される。これは、アプリケーションや統計データの更新などに対するデータ転送のために使用されているバッファの解放を含むことができる。転送中は、「iSCSI PDU」は、伝送コマンドエンジン(ブロック

10

20

30

40

50

1110)と協働する伝送エンジン(ブロック1109)によって伝送され、このコマンドエンジンは、PDUを変換し、ストレージプロセッサに対するDMAを使用し、ホストメモリからのアプリケーションバッファを検索するなどの適切な作動を実行し、進行に伴って更新されたiSCSI接続データベース内のストレージコマンドフロー情報を維持する。本特許で使用されるように、用語「エンジン」は、エンジンの機能又は使用に適するデータプロセッサ又はデータプロセッサの一部とすることができる。同様に、受信エンジン(ブロック1107)は、受信コマンドを新しい要求、応答、エラー、又は他のコマンド、又は適切に作用されるべきデータPDUに翻訳する。コマンドエンジン(ブロック1106)と協働するこれらの受信エンジンは、読取データ又は受信データを、直接データ転送/配置、又はiSCSIセッションテーブル内のセッションのために維持されたRDMA制御情報を通じて、適切に割り当てられたアプリケーションバッファにルーティングする。コマンドが終了すると、それぞれのバッファに対する制御(ブロック1103及び1112)は、使用するアプリケーションのために解放される。受信及び送信エンジンは、図23のTCP/IPプロセッサ又は図24のIPプロセッサから見る時にグローバルメモリとして見るることができる、セッションデータベースエントリ1704に記録されたセッション情報と協働するこのIPプロセッサの図17のSANパケットプロセッサ1706(a)から1706(n)とすることができる。ブロック1704と図17の1708のセッションデータベースの一部に関連して以下により詳細に説明される、セッションデータベースによって提供された適切なストレージフローコンテキストを使用して、同じエンジンを異なるパケットとコマンドのために再使用することができる。明確にするために、IPネットワークアプリケーションプロセッサ、IPストレージプロセッサ、NEIPストレージネットワークアプリケーションプロセッサ、及びIPプロセッサという用語は、アプリケーション次第で同じエンティティとすることができる。IPネットワークアプリケーションプロセッサコア及びIPストレージネットワークアプリケーションプロセッサコアは、アプリケーション次第で同じエンティティとすることができる。

#### 【0039】

同様に、制御コマンドは、送信パスを使用することができ、一方、受信応答は、受信パスを使用するであろう。類似のエンジンは、ターゲットと同様にイニシエータ上に存在することができる。データフローの方向は、それがイニシエータかターゲットかによって異なる。しかし、ターゲットでは付加的な段階を伴うが、本質的にはイニシエータとターゲットの両方に類似のデータフローが存在する。例えば、書込コマンドデータを取得するのに必要なバッファを保証するために、ターゲットは、付加的な作動を実行する必要があるか又はイニシエータにデータが供給される前に読取データを準備する必要があるであろう。中間装置の場合は同様の例が存在するであろうが、そのようなスイッチ又はアプライアンスとすることができる装置では、あるレベルの仮想化又はフレーム濾過、又は、一方でセッションの終了を他方でセッションの開始を要求する場合がある類似の別の作動を実施することができる。この機能性は、このアーキテクチャによってサポートされるが、当業者には公知の範囲なので図には明確には示されていない。

#### 【0040】

図12から図15は、移送セッションに関する特定のプロトコル情報と、その情報をメモリ内のデータベースに保存することができる方法とを示す。

図12は、iSCSIプロトコル及び関連TCP/IP接続のために維持されたデータ構造を示す。各iSCSIセッション(ブロック1201)に属するデータは、基本的にイニシエータとターゲット接続のネクサスであり、適切な接続(ブロック1202)上で搬送される。従属コマンドは、同じ接続の待ち行列上でスケジュールに入れられ、コマンドの順序を維持する(ブロック1203)。しかし、関連のないコマンドは、別の移送接続に割り当てることができる。実施がセッション当たり1つの接続だけをサポートする場合、全てのコマンドを同じ接続の待ち行列に入れることが可能である。しかし、イニシエータとターゲットの間の回線中継をサポートするために、セッション当たり複数の接続が実現可能である。例えば、いくつかの実施形態では、イニシエータとターゲットは、互い

10

20

30

40

50

に通信が可能であり、複数接続を受け入れるための交渉を通じて決めることになる。別の実施形態では、イニシエータとターゲットは、1つのセッション又は接続だけを通じて通信する。図13と図14は、TCP/IP及びiSCSIセッションデータベース、又は、セッション及び接続毎の伝送制御ブロックを示す。これらのエントリは、別々のテーブルとして搬送することができ、又は、次に図23、24、26、及び29に関連して見られるように、複合テーブルとして一緒に搬送することもできるが、TCP/IPだけ、RDMAを有するTCP/IP、IPストレージだけ、TCP/IPを有するIPストレージ、及びRDMAを有するIPストレージなど、選択された実施と実施された機能性とに依存する。TCP/IP及びストレージフロー制御を実施する様々なエンジンは、これらのフィールドの全部又はいくつか、又は、示されていないより多くのフィールドを使用して、TCP/IPに亘ってブロックデータ転送を方向付けする。データ転送過程の間に複数の状態を通じて接続が進行する時、適切なフィールドが更新される。図15は、オンチップセッションキャッシュ(ブロック1501及び1502)、及びデータ転送の連続した進行のために必要な状態情報を保持するオフチップセッションメモリ(ブロック1503、1504、1505、1506、及び1507)から成るメモリサブシステムに転送制御エントリを保存する一方法を示す。

#### 【0041】

図16は、抽象化レベルの高いIPプロセッサアーキテクチャを示す。プロセッサは、モジュールの拡張可能IPネットワークアプリケーションプロセッサコア(ブロック1603)から成る。その機能性ブロックは、IPネットワーク上で高速保存とデータ転送を利用可能にする機能性を提供する。プロセッサコアは、インテリジェントフローコントローラプログラマブル分類エンジン、及びストレージ/ネットワークポリシーエンジンを含むことができる。それぞれは、独立したプロセッサと考えることができ、それらの任意の組合せは、単一プロセッサとして実行することができる。開示されたプロセッサはまた、セキュリティ処理ブロックを含み、ネットワークパケットに対して高回線速度の暗号化及び暗号解読機能を提供する。同様に、これは、単一プロセッサにすることができ、又は上述の他のものと組み合わせることができる。開示されたプロセッサは、メモリサブシステムを含み、このサブシステムは、オンチップセッションキャッシュ/メモリを管理するメモリコントローラインタフェースと、オフチップメモリへのアクセスを管理するメモリコントローラ(ブロック1602)とを含み、このオフチップメモリは、SRAM、DRAM、FLASH、ROM、EEPROM、DDR SDRAM、RDRAM、FCRAM、QDR SRAM、又は、静的又は動的ランダムアクセスメモリの他の派生物、又はその組合せとすることができる。IPプロセッサは、適切なシステムインタフェースを含み、それを目標のマーケットセグメントで使用することができるようにし、LAN、SAN、WAN、及びMANネットワーク、並びに同様のネットワーク、並びに適切なホストインタフェース(ブロック1606)に対して、適正なメディアインタフェース(ブロック1601)を提供する。メディアインタフェースブロックとホストインタフェースブロックとは、マルチポートの形態で存在することができ、そこでは、いくつかのポートは、開示されたプロセッサが使用されるネットワーク及びシステムで冗長性及びフェイルオーバー機能を供給することができる。プロセッサはまた、例えば、複数プロセッサシステムを作り出してメインプロセッサの機能を拡張するために、共用プロセッサインタフェース(ブロック1605)を含むことができる。ブロック1604のシステムコントローラインタフェースは、このプロセッサが汎用マイクロコントローラとインタフェースで接続することを可能にし、このマイクロプロコントローラは、開示されたプロセッサを使用することができるシステムに対して、システムコントローラとして作動することができる。プロセッサアーキテクチャはまた、システムコントローラ又はセッションマネージャとして作動することができる搭載された制御プレーンプロセッサをサポートする。システムコントローラインタフェースを依然として提供し、外部プロセッサの使用を可能にすることができる。このプロセッサのこのようなバージョンは、ダイのコストの理由から制御プロセッサを含むことができない。例えば、サーバアダプタ、又はストレージコントローラ、又は

10

20

30

40

50



スイッチ回線カード、又は他のネットワーキングシステムのような特定のシステム要件に的を絞って作り出すことができる様々な形式のコアアーキテクチャがある。主な違いは、本特許の前出の節で上述したようなことであろう。これらのプロセッサブロックは、i S C S I、F C I P、i F C I Pなどのような標準的プロトコルを使用し、高性能IPベースのストレージを達成するための機能と性能をもたらす。これらのブロックの詳細なアーキテクチャを以下に説明する。

#### 【0042】

図17は、IPプロセッサアーキテクチャをより詳細に示す。上述のアーキテクチャは、完全なTCP/IPの終了と深いパケットの検査とを通じてメディアアクセス制御装置(MAC)又は他の適切な装置からの着信IPパケットを処理するための機能を提供する。この図表は、図16のMAC層のブロック1601、又はブロック1602、1604、又は1605を示していない。メディアインタフェース(ブロック1601)内のプロセッサの入力待ち行列(ブロック1701)と出力待ち行列(ブロック1712)とに対するMAC層インタフェースブロックは、図16に示されている。MACの機能性は、ネットワークに依存する特定の形式を有する標準ベースとすることができるであろう。イーサネット及びSONET上パケットは、今日最も広範囲に使用されるインタフェースの例であり、それらは、同じシリコン上、又は各々を使用して作り出されたプロセッサの異なるバージョンに含めることができる。

#### 【0043】

図17のブロック図は、入力待ち行列及び出力待ち行列ブロック1701及び1712を2つの別のブロックとして示す。機能性は、組み合わせられたブロックを使用して提供することができる。入力待ち行列1701は、MACインタフェースブロックからの着信パケットを検索するために、論理、制御、及びストレージから成る。ブロック1701は、パケットスケジューラ1702及び分類エンジン1703と協働してインタフェースから到着するパケットを待ち行列に入れ、パケットの最初、パケットの最後、及び、フラグメント化したパケット又は安全なパケットのような他の属性を識別するために適切な標識を作り出す。パケットスケジューラ1702は、入力待ち行列コントローラからパケットを検索し、分類のためにそれらを分類エンジンに転送することができる。分類ブロック1703は、スケジューラに続くように示されているが、論理的な観点から、分類エンジンが入力待ち行列からパケットを受け取り、パケットを分類して分類タグをパケットに提供し、そのパケットが次にスケジューラによってプロセッサアレイ1706(a) . . . 1706(n)に対してスケジューラに入れられる。従って、分類エンジンは、通過分類エンジンとして作用することができ、その構成を最高回線速度で使用してパケットフローを維持する。分類エンジンは、プログラマブルエンジンであり、様々なカテゴリーのネットワークからパケットを受け取り、識別の結果を使用して、スケジューラと、使用すべき他のパケットプロセッサとのためにパケットにタグ付けする。ネットワークトラフィックの分類は、非常に計算集約的なタスクであり、パケットプロセッサで利用可能なプロセッササイクルの半分を占める。この一体化された分類エンジンは、第2層から第7層の検査を実行するためにプログラム可能である。分類されるフィールドは、適合するものがあれば、それらに対応する比較と処置のための期待値を使用してプログラムされる。図30に関連して次に見られるように、分類子は、分類ウォークの結果を収集し、分類結果を識別するパケットに対するタグとしてこれらを呈示することができる。これは、ツリー構造によく似ており、「ウォーク」として理解されている。分類されたパケットは、次に、処理のパイプラインの次の段階としてスケジューラ1702に供給される。

#### 【0044】

パケットスケジューラブロック1702は、状態コントローラと、開示されたプロセッサ上の適切な実行エンジンにパケットを割り当てるシーケンサとを含む。実行エンジンは、所望する実施次第で、TCP/IP、及び/又はストレージエンジン、並びにストレージフロー/RDMAコントローラ(ブロック1708)、又はホストバイパス、及び/又は他の適切なプロセッサを含むSANパケットプロセッサ(ブロック1706(a))から

10

20

30

40

50

1706(n))である。明確にするために、記号「/」は、本特許でハードウェア構成要素を指す時は、適宜「及び/又は」を意味することができる。例えば、構成要素「ストレージフロー/RDMAコントローラ」は、実行に対して適宜、ストレージフロー及びRDMACコントローラ、ストレージフローコントローラ、又はRDMACコントローラとすることができる。スケジューラはまた、パケットから同じ接続/セッション上のパケットへの状態の従属性が着信パケットの正しい処理のために重要である場合、プロセッサを通じてパケット順序化を維持する。パケットのリタイヤメントまで、プロセッサを使用してスケジュールに入れたパケットの進行を追跡するために、スケジューラは、様々なテーブルを維持する。スケジューラはまた、出ていくコマンド上のパケットプロセッサと、ホストプロセッサ又は切替ファブリックコントローラ又はインタフェースからのパケットとに対してスケジュールに入れる必要があるコマンドを受け取る。

10

## 【0045】

プログラマブルパケットプロセッサに加えて、TCP/IP及びストレージエンジンは、図17のSANパケットプロセッサ1706(a)から1706(n)として一緒にラベル付けされる。これらのパケットプロセッサは、独立のプログラマブル構成要素であるエンジンであり、特定の役割を果たす。代替的に、2つ又はそれ以上のそれらを望ましい実行次第で単一プロセッサとして実行することができる。図23のTCP/IPエンジンと図24のストレージエンジンは、この例で、図21のプログラマブルパケットプロセッサエンジンブロック2101に対する共有プロセッサとして構成される。このアーキテクチャは、コスト、製造、マーケットセグメントのような理由からストレージエンジンを置換/除去することにより、ストレージ以外のアプリケーションに比較的容易に適用することができる。純粋なネットワーキング環境では、ストレージエンジンは、専用のTCP/IPエンジンを有するパケットプロセッサを残して除去され、ネットワーキングトラフィックの代わりに適用されることがあり、TCP/IPソフトウェアスタックによる同様の処理のオーバーヘッドに直面する。代替として、1つ又はそれ以上のエンジンを望ましい実施のために落としてもよく、例えば、IPストレージ機能だけをサポートするプロセッサに対して、別のチップ内にあるTCP/IPエンジン及び/又はパケットエンジンを落とすことができる。従って、コアの拡張可能モジュラーアーキテクチャの複数の変形が実現可能である。従って、コアのアーキテクチャは、例えば、アプリケーション専用の加速度的ためにサポートを提供する高性能ネットワークセキュリティ及びポリシーエンジン、高性能ルーティングエンジン、高性能ネットワーク管理エンジン、ストリング検索を提供する深いパケット検査エンジン、XMLのエンジン、仮想化のためのエンジンのような他の専用エンジンでストレージエンジンを置換することにより、IPアプリケーション上のストレージに加えて、アプリケーション内で強化することができる。このIPプロセッサの処理機能は、ネットワークインタフェースの回線速度の要求に合致するように、チップ内のSANパケットプロセッサブロック1706(a)から1706(n)の数を増減することにより、スケールアップすることができる。拡張容易性からの根本的な制約は、要求されるシリコンのスペースとシリコン処理技術によって課せられる制約から生じる。根本的にこのアーキテクチャは、より多くのSANパケットプロセッサブロックを追加して処理機能を向上させることにより、かなりの高回線速度まで拡張することができる。同様の結果を達成する他の手段は、プロセッサの作動のクロック周波数を処理技術の限界内で実現可能な程度まで上げることである。

20

30

40

## 【0046】

図17はまた、IPセッションキャッシュ/メモリ及びメモリコントローラブロック1704を示す。このキャッシュは、内部メモリ又はローカルセッションのデータベースキャッシュと見ることができる。このブロックは、TCP/IPセッションデータベースと、同様に特定数のアクティブセッションに対するストレージセッションデータベースとをキャッシュして保存するために使用される。キャッシュされるセッションの数は、選択されたシリコンスペースと、製造するための経済的な可能性との直接の結果である。オンチップでないセッションは、ブロック1704の一部か又は別様の高性能メモリコントロー

50

ラブロックを使用して、外部メモリとして見られるオフチップメモリへノから保存及び検索される。このプロセッサの様々な処理要素は、高速内部バスを使用してこのコントローラを共有し、セッション情報を保存して検索する。メモリコントローラはまた、フラグメント化したパケットを一時的に保存するために、又はホストインタフェース又は外部行き待ち行列がバックアップされる時に使用することができる。コントローラはまた、統計データ情報、又は、開示プロセッサか又は開示又はホストプロセッサ上で作動しているアプリケーションによって収集することができる他の任意の情報を保存するために使用することができる。

#### 【 0 0 4 7 】

図 17 のプロセッサのブロック図はまた、ホスト入力待ち行列 (ブロック 1707) とホスト出力待ち行列 (ブロック 1709)、並びにストレージフロー / RDMA コントローラ (ブロック 1708) を示す。これらのブロックは、ホスト (「ピア」とも呼ばれる) メモリ又は切替ファブリックへ及びそこからデータを転送するために必要な機能を提供する。これらのブロックはまた、ホストベースのドライバがコマンドをスケジュールに入れ、着信ステータスを検索し、セッションデータベースのエントリを検索し、開示されたプロセッサをプログラムすることができるようにする機能を提供し、ソケット直接アーキテクチャ、TCP/IP の完全終了、IP ストレージオフロードのような機能を RDMA を使用して、又は使用しないで実行することができるようにする。図 27 により詳細に見られるホストインタフェースコントローラ 1710 は、構成レジスタ、メモリ対メモリの直接データ転送のための DMA エンジン、及び上述のいくつかのタスクを実行するホストコマンドブロックをホストインタフェース処理コントローラ及びホスト割り込みコントローラと共に提供する。ホスト入力及び出力待ち行列 1707 及び 1709 は、着信パケットと発信パケットに待ち行列を提供する。ストレージフローと RDMA コントローラブロック 1708 は、開示されたプロセッサに対してホストがコマンドを待ち行列に入れるために必要な機能性を提供し、このプロセッサは、次にこれらのコマンドを取得してそれらを実行し、コマンドが終了するとホストプロセッサに割り込む。ブロック 1708 の RDMA コントローラ部分は、遠隔の直接メモリアクセスを可能にするために必要な様々な機能を提供する。それは、RDMA 領域、アクセスキー、及び仮想アドレス変換機能性のような情報を含むテーブルを有する。このブロック内部の RDMA エンジン、データ転送を実行し、受け取った RDMA コマンドを変換し、実行が許可されている場合に処理を実行する。ブロック 1708 のストレージフローコントローラはまた、ターゲットとイニシエータの間でデータ転送が発生した時に、スケジュールに入れられた様々なコマンドの進行状態の追跡を保存する。ストレージフローコントローラは、実行するためのコマンドをスケジュールに入れ、コマンド終了情報をホストドライバに提供する。以上は RDMA 機能と考えられ、説明したように又は個々のプロセッサを実施することにより、設計者の選択次第で実施することができる。また、本特許の精神及び範囲から逸脱することなく、これらの説明に対して付加的な機能を追加するか又は排除することができる。

#### 【 0 0 4 8 】

このプロセッサの制御プレーンプロセッサブロック 1711 は、ICMP プロトコルによるエラー処理、ネーム解決、及びアドレス解決プロトコルを含むことができる TCP/IP 及びノ又はストレージプロトコルに対する比較的低速のルーティング機能を提供するのに使用され、それはまた、プログラムされてセッションコントローラ / 接続マネージャとしてのセッション開始 / 取外し作動、ログイン及びパラメータ交換などを実行することができる。この制御プレーンプロセッサは、システム開発者に制御プレーンプロセッサの選択を提供するためにオフチップにするか、又は統合されたソリューションを提供するためにオンチップにすることができる。制御プレーンプロセッサがオフチップの場合、インタフェースブロックがそこで作り出されるか又は統合され、このプロセッサが制御プレーンプロセッサとインタフェースで接続することができるようになり、データとコマンドの転送を実行する。内部バス構造及び機能ブロック相互接続は、性能やダイのコスト要件などに関する全ての詳細図に対して図示と異なることがあり、それらは、本発明の精神及び

10

20

30

40

50

範囲から逸脱しないものである。

図17のブロックに関して上述した機能は、以下により詳細に説明するように、開示されたプロセッサの様々な処理リソースによるインストリーム処理を使用して、パケットが最低限の待ち時間で入力から出力まで通過することができるパケットストリームアーキテクチャを実行可能にする。

#### 【0049】

図18は、一般的に、詳細には図17の1701で表される入力待ち行列とコントローラブロックを示す。このブロックの核になる機能性は、ブロック1801及び1802(i)から1802(n)において複数の入力ポートであるポート1からNからの着信パケットを受け入れ、パケットが分類子、スケジューラ、及びスケジューラI/Fブロック1807~1814を通した更に別の処理に対してデキューされた入力パケット待ち行列(ブロック1810)に対する固定又はプログラマブル優先度を使用して、それらを待ち行列に入れることである。入力待ち行列コントローラは、各入力ポート(マルチポート実行のポート1からポートN)とインタフェースで接続し、入力パケット待ち行列上1810に対してパケットを待ち行列に入れる。パケット待ちコントローラとマーカブロック1804は、ポート速度、ポートのネットワークインタフェース、ポート優先度、及び他の適切な様々な特徴に基づいて、固定された優先度機能を提供することができ、又は、別々のインタフェースに対して別々のポリシーが適用されるようにプログラムすることができる。ラウンドロビンや重み付けラウンドロビンなどのような様々なモードの優先度をプログラムすることができる。入力パケットデキューコントローラ(1812)は、パケットをデキューし、それらをスケジューラI/F1814を通じてパケットスケジューラ(図17のブロック1702)に供給する。スケジューラは、図17の分類エンジン1703によってパケットが分類された状態で、SANパケットプロセッサ1706(a)~1706(n)に対してパケットをスケジュールに入れる。暗号化されたパケットは、最初に暗号化される時に分類することができ、実施がセキュリティ処理を含む場合、認証及び/又は暗号解読のために、図18の安全なパケットインタフェースブロック1813により、図17のセキュリティエンジン1705に転送されるが、そうでない場合は、セキュリティインタフェースは存在しないことがあり、同様の機能を実行するために外部セキュリティプロセッサが使用されるであろう。クリアなパケットインタフェース(ブロック1811)から暗号解読されたパケットは、次に、パケットがクリアなパケットと同じパスを追従するブロック1812を通じて入力待ち行列に供給される。

#### 【0050】

フラグメント化したIPパケットは、フラグメント化パケット保存制御バッファ(ブロック1806)内のオンチップに保存することができるか、又は内部又は外部メモリに保存することができる。最後のフラグメントが到着すると、ブロック1806のフラグメント化コントローラは、図17の分類エンジン及びスケジューラと協働して、これらのフラグメントを組み合わせて完全なパケットを組み立てる。フラグメント化したパケットが組み合わせられて完全なパケットを形成した状態で、パケットは、ブロック1804を通じて入力待ち行列に対してスケジュールに入れられ、次に、パケットデキューコントローラ(ブロック1812)によって処理され、このプロセッサの様々な他の処理ステージに通される。図18の入力待ち行列コントローラは、ブロック1808のパケット属性及びタグアレイに保存された、分類エンジンからの他の任意のセキュリティ情報と共にパケット開始、サイズ、バッファアドレスのようなパケット記述子フィールドを使用するブロック1809の属性マネージャによって管理されるパケットタグ/記述子を各着信パケットに割り当てる。パケットタグ及び属性は、プロセッサのスケジューラと他の構成要素により、インタフェース807、1811、1813、及び1814を通じて、プロセッサを通過するパケットのフローを効率的に制御するために使用される。

#### 【0051】

図19は、図17のパケットスケジューラとシーケンサ1702をより詳細に示す。このブロックは、パケットのスケジューラ設定とこのプロセッサの実行リソースに対するタ

10

20

30

40

50

スクとを担い、従って、負荷バランサとしても作用する。スケジューラは、入力待ち行列コントローラ1901から、ヘッダの待ち行列(ブロック1902)のデータのヘッダを検索し、それらを図17の分類エンジン1703に通し、このエンジンは、次に、残りのプロセッサによって使用される分類子待ち行列(ブロック1909)の分類結果を戻す。分類エンジンは、最初はヘッダと共に呈示することができるが、深いデータ検査もプログラムされる場合は、分類エンジンは、分類の後でスケジューラにルーティングする完全なデータを受け取ることができる。スケジューラは、分類エンジンを通じてデータの実行を管理する分類コントローラ/スケジューラ(ブロック1908)を有する。図19のこのブロックは、フラグメント化したデータ又は安全なデータの場合は、そのようなデータに対して適切な動作を行うために、例えば、図17のセキュリティエンジン

10

に対して暗号化されたデータをスケジュールに入れるために、入力待ち行列コントローラ(ブロック1901)にコマンドを供給する。スケジューラ状態制御及びシーケンサ(ブロック1916)は、プロセッサ内部でアクティブな様々な処理/動作の状態情報を受け取り、次の動作のセットに対する命令を提供する。例えば、スケジューラは、ブロック1903の入力データ待ち行列からデータを検索し、分類子から受け取った分類の結果に基づいて、適切なリソース待ち行列に対してこれらのデータをスケジュールに入れるか、又はデータをデータメモリ(ブロック1913又は1704から1906)に導き、スケジュール設定時又はその後動作を実行するために、適切なリソースが必要とする時にデータを検索するために使用することができるデータ記述子/タグを

20

作り出す。状態制御及びシーケンサブロック1916は、分類結果(ブロック1914)を有するデータが、データが動作のためにスケジュールに入れられる時に検索されるデータメモリ(ブロック1913)に保存されるように命令/誘導する。状態コントローラ及びシーケンサは、動作のためのデータを受け取るべき実行リソースを識別し、コマンドを作り出し、このコマンドをデータタグと共に図19のリソース待ち行列ブロック1917(制御プレーン)、1918(ポートi~ポートn)、1919(バイパス)、及び1920(ホスト)に割り当てる。優先度セクタ1921は、割り当てられた優先度に基づいて、それぞれの待ち行列からコマンドとデータタグを検索し、これをデータフェッチとコマンドコントローラ(ブロック1922)に転送するプログラムブロックである。このブロックは、分類結果と共にデータメモリストア1913からデータを検索し、例えば、1926でリソースが動作待機中などの時に、高性能プロセッサ

30

のコマンド上の適切なリソースとデータバスに対してデータ転送をスケジュールに入れる。それぞれの受け取り側のコマンドバスインタフェースコントローラ1905のようなバスインタフェースブロックは、コマンドを解釈してデータと動作のための分類タグを受け入れる。これらの実行エンジンは、何時データ動作が終了し、何時データが最終の宛先(ホストバスインタフェース、又は出力インタフェース、又は制御プレーンインタフェースなどのいずれか)に対してスケジュールに入れられるかをスケジューラに通知する。これによって、スケジューラは、ブロック1904のリタイアメントエンジンの助けを借りて、データをその状態から退去させることができ、リソース割当テーブル(ブロック1923)のこのセッションのためにリソースエントリを解放する。リソース割当

40

テーブルは、例えば、SANデータプロセッサエンジンのバッファに入れられたセッションデータベースキャッシュエントリや発信元内で実行されている現在のデータの接続IDなどのこれらのリソースの内部状態の現在の状態により、受け取ったデータを特定のリソースに割り当てるためにシーケンサによって使用される。従って、順序付けられた実施に依存するデータは、最初に同じリソースに割り当てられ、プロセッサのセッションメモリの現在のDB状態を使用してメモリのトラフィックと性能を改善し、従って、新しいセッションのエントリを検索しなくてもよい。シーケンサはまた、データを待ち行列に入れるためにメモリコントローラ(ブロック1906)とのインタフェースを有し、これらのデータは、フラグメント化されたデータであり、及び/又は、最大回線速度性能を維持するために割り当てられたものよりも多くの時間が掛かり、データ上で実行された特定のアプリケーションが原因となることがある下流のデータ処理の妨害のため

50

にスケジューラの待ち行列が渋滞する場合のものであり、又は、他の任意の下流のシステムが一杯になって回線速度が維持できない場合のものである。

【 0 0 5 2 】

分類エンジンが入力待ち行列からパケットを受け取る図 1 7 に関連して上述したように、分類子がスケジューラの前で実施される場合、アイテム 1 9 0 1、1 9 0 2、1 9 0 8、1 9 0 9、及び 1 9 1 0 は、個々の設計次第で分類子内にあるか又は必要ないこともある。分類子からスケジューラブロック 1 9 0 3、1 9 0 7、1 9 1 4、及び 1 9 1 5 まで / からの適切な連結は、そのようなシナリオで作り出すことができ、分類子は、図 1 8 の入力待ち行列ブロックに直接連結される。

【 0 0 5 3 】

図 2 0 は、通常は図 1 7 の 1 7 0 3 で表されるパケット分類エンジンを示す。パケットのそれらの様々な属性への分類は、非常に計算集約的な作動である。分類子は、パケット形式、例えば、I P、I C M P、T C P、U D P のようなプロトコル形式、ポートアドレス、発信元及び宛先フィールドなどを識別するために、受け取ったパケットの様々なフィールドを試験するプログラマブルプロセッサとすることができる。分類子は、ヘッダ又はペイロード内の特定フィールド又は一組のフィールドを試験するために使用することができる。ブロック図は、内容アドレスメモリベースの分類子を示す。しかし、先に説明したように、これもまた、プログラマブルプロセッサとすることができる。主な違いは、性能とエンジンの実施の複雑さである。分類子は、図 2 0 の入力待ち行列（ブロック 2 0 0 5 及び 2 0 0 4 ）からスケジューラを通じて入力パケットを取得する。入力バッファ 2 0 0 4 は、分類される必要があるパケット / 記述子及び / 又はパケットヘッダを待ち行列に入れる。次に、分類シーケンサ 2 0 0 3 は、待ち行列内の利用可能なパケットをフェッチし、プログラムされるか又はプログラムすることができるグローバルフィールド記述子の組（ブロック 2 0 0 7 ）に基づいて、適切なパケットフィールドを抽出する。次に、分類子は、内容アドレスメモリ（C A M ）アレイ（ブロック 2 0 0 9 ）にこれらのフィールドを通し、分類を実行する。フィールドは、C A M アレイを通過するので、これらのフィールドの符合は、比較されるべきフィールドの次の組と、場合によってはそれらのビットフィールド位置とを識別する。C A M アレイでの符合は、アクション / イベントタグをもたらし、このタグは、結果コンパイラ（ブロック 2 0 1 4 ）から収集され（「コンパイル」が「収集」の意味で使用される場合）、また、特定の C A M 条件又は規則の符合に関連するメモリアレイ（ブロック 2 0 1 3 ）のデータを更新することができるアクションとしても作用する。これは、演算論理装置（A L U ）作動（ブロック 2 0 1 7 ）の実行を含むことができ、これは、例えば条件符合の増分又は減分などのこのフィールド上の実行リソースの一例と考えられる。C A M アレイは、ホスト又は制御プレーンプロセッサインタフェース 1 7 1 0 及び 1 7 1 1 を通してプログラミングのためにアクセス可能なデータベース初期化ブロック 2 0 1 1 を通じて、比較すべき次のフィールドを含むフィールド、それらの期待値、及び符合に対するアクションを使用してプログラムされる。分類がリーフノードに達した状態で、分類は終了し、同じ分類タスクを実行するのを回避するために I P プロセッサの他のエンジンによって次に使用することができる、トラバースされたパスを識別する分類タグが生成される。例えば、分類タグは、フロー又はセッション I D、T C P / U D P / I C M P のようなプロトコル形式の表示、処理、バイパス、パケットのドロップ、セッションのドロップなどをするか否かを示す値を含むことができ、又は実行リソースがパケット処理を開始するための特定のファームウェアコードのルーチンポイントを含むことができ、又はトラバースされた分類パスの署名などを含むこともできる。分類タグフィールドは、プロセッサの実行と機能性に基づいて選択される。分類子リタイヤメント待ち行列（ブロック 2 0 1 5 ）は、分類されたパケットのパケット / 記述子と分類タグを保持し、スケジューラによって検索されるのを待機する。分類データベースは、データベース拡張インタフェースとパイプライン制御論理ブロック 2 0 0 6 とを使用して拡張することができる。これによって、より大きな分類データベースのための拡張性を必要とするシステムの構築が可能になる。アクションインタープリタ、A L U、及び範囲符合ブロック

10

20

30

40

50

2012を有する分類エンジンはまた、ストレージ/ネットワークポリシー/特定のポリシーが満足された場合に取りの必要があるアクションをプログラムする機能を提供する。ポリシーは、規則及びアクションテーブルの形態で実行することができる。ポリシーは、分類テーブルと共にホストインタフェースを使用して分類エンジン内でコンパイルされ、プログラムされる。データベースインタフェースとパイプライン制御2006は、分類/ポリシーエンジンを拡張するためにコンパニオンプロセッサと連結するように実施することができる。

#### 【0054】

図21は、通常は図7の1706(a)から1706(n)で示すようなSANパケットプロセッサを示す。パケットプロセッサは、特別に設計されたパケットプロセッサか、  
 又はARM、MIPS、StrongARM、X86、PowerPC、Pentium  
 プロセッサのような任意の適切なプロセッサか、又は本明細書に説明した機能を満足させる他の任意のプロセッサとすることができる。これは、本特許の様々な節においてパケット  
 プロセッサ複合システムとも呼ぶ。このパケットプロセッサは、パケット処理のための  
 ターゲット命令、又はTCP/IPエンジン(ブロック2102)、又はIPストレージ  
 エンジン(ブロック2103)、又はその組合せを備えた、一般的にRISCマシンである  
 パケットエンジン(ブロック2101)を有する。これらのエンジンは、パケットエン  
 ジンに対する共有エンジンとして、又は独立したエンジンとして構成することができる。

図22は、パケットエンジンをより詳細に示す。パケットエンジンは、通常は上述のよう  
 にパケット処理マイクロルーチン及びパケット及び中間ストレージを保持するために使用  
 される命令メモリ(ブロック2202)及びデータメモリ(ブロック2206)(両方と  
 もRAMである)を備えたRISCマシンである。命令メモリ2202は、本特許の全て  
 のこのようなメモリと同様にRAM又は他の適切なストレージであり、パケット処理中  
 に実行されるコードを使用して初期化される。パケット処理コードは、割り当てられたメモ  
 リ内に適合する緊密なマイクロルーチンとして編成される。命令復号器及びシーケンサ  
 (ブロック2204)は、命令メモリ2202からの命令をフェッチしてそれらを復号化し  
 、ALU(ブロック2208)内に包含された実行ブロックを通じてそれらを順序付けす  
 る。このマシンは、単純なパイプラインエンジンか、又はパケット志向命令セットを提供  
 するように設計されたより複雑な深いパイプラインマシンとすることができる。DMAエ  
 ンジン(ブロック2205)及びバスコントローラ(ブロック2201)により、パケッ  
 トエンジンがデータパケットを図19のスケジューラから移動し、ホストインタフェース  
 を作動のためにデータメモリ2206に移動させることが可能になる。DMAエンジンは  
 、パケット/データをメモリ/パケットメモリへ/から保存/検索するために複数のメモ  
 リ記述子を保持することができる。これによって、メモリアクセスが、プロセッサエン  
 ジン作動と平行して発生するようになる。DMAエンジン2205はまた、データパケッ  
 トをTCPとストレージエンジン2210及び2211の間で移動させるために使用す  
 ることができる。パケットの実行が終了した状態で、抽出されたデータ又は新しく生成された  
 パケットは、メディアインタフェース又はホストインタフェースのいずれかに向けて出力  
 インタフェースに転送される。

#### 【0055】

図23は、通常は図22の2210に見られるプログラマブルTCP/IPパケットプロ  
 セッサエンジンをより詳細に示す。このエンジンは、一般的に、様々なTCP/IP志  
 向命令及び実行エンジンと共に共有RISC命令を有するプログラマブルプロセッサであ  
 るが、本特許に説明する適切な実行エンジンを備えたマイクロコード化又は状態機械駆動  
 プロセッサとすることもできるであろう。TCPプロセッサは、TCPチェックサムの確認、  
 及びプロセッサ上のこれらの命令の実行による新しいチェックサム生成のためのチェ  
 ックサムブロック2311を含む。チェックサムブロックは、パケットバッファメモリ2  
 309(データRAMがそのようなメモリの一例である)からデータを抽出し、チェ  
 ックサムの生成又は確認を行う。パケット検索インタフェースブロック2310は、実行エン  
 ジンと命令シーケンサ2305を助け、様々なデータパケットフィールド又は全データパ

10

20

30

40

50

ケットにアクセスを提供する。分類タグインタープリタ 2313 は、そのような実行が選択された場合、分類の結果に基づいてプログラムフローを誘導するために命令復号器 2304 によって使用される。プロセッサは、例えば、次の予想されるシーケンス番号を検索し、受け取ったものが合意されたスライドウィンドウの範囲内にあるか否か、また、パケットが属する接続に対してどのスライドウィンドウが TCP プロトコルの公知の部分であるかを見るために、特定のシーケンスと TCP/IP データ順序付け計算で使用するための分割 (ブロック 2315) を含むウィンドウ作動を提供する。この要素 2315 は、図 24 の 2413 で表すような分割コントローラを含むことができる。代替的に、当業者は、本特許の教示を使用して、この図 23 の TCP/IP プロセッサ上のどこにでも分割コントローラを容易に実施することができる。このプロセッサは、ハッシュエンジン (ブロック 2317) を提供し、このエンジンは、パケットの特定フィールドに対してハッシュ作動を実行するために使用され、パケットに対して正しいセッションエントリを取得するために必要なハッシュテーブルウォークを実行する。プロセッサはまた、データの発信元及び宛先のためのポインタレジスタ、コンテキストレジスタセット、及び汎用レジスタファイルに加えて TCP 状態を保持するレジスタと共に、TCP 処理のための様々な共通に使用されるヘッダフィールドを抽出するレジスタファイル (ブロック 2316) を含む。TCP/IP プロセッサは、パケットの実行に対して複数のコンテキストを有することができ、そのために、例えばメモリアクセスのような何らかの理由で所定パケットの実行が機能停止すると、別のコンテキストが喚起され、プロセッサは、ごく僅かの効率ロスで別のパケットストリームの実行を継続することができる。TCP/IP プロセッサエンジンはまた、ローカルセッションキャッシュ (ブロック 2320) を維持し、このキャッシュは、最も最近使用された又は最も頻繁に使用されたエントリを保持し、グローバルセッションメモリから検索する必要なしに局所的に使用することができる。ローカルセッションキャッシュは、パケットプロセッサとすることができる TCP/IP プロセッサの内部メモリと考えることができる。セッション又はグローバルメモリから付加的なものを検索することなく、より多くのエントリが使用されて内部メモリに保存することができるほど、処理がより効率的になるのは勿論である。図 19 のパケットスケジューラは、TCP/IP リソース毎にキャッシュされた接続 ID を通知され、従って、同じパケットプロセッサ複合システムに対して同じセッションに属するパケットをスケジューラに入れることができる。特定の接続に対してパケットプロセッサがセッションエントリを保持しない場合、TCP セッションデータベース検索エンジン (ブロック 2319) は、セッションマネージャ (ブロック 2321) 及びハッシュエンジンと協働し、グローバルセッションメモリからメモリコントローラインタフェース (ブロック 2323) を通じて対応するエントリを検索する。ハッシュエンジンと共に作用して、対応するセッションエントリ又はそのフィールドをセッションデータベースに対して保存/検索するためのセッション識別子を生成するセッションエントリ又はセッションエントリのフィールドのアクセスを可能にするセッションマネージャ内の論理回路のような手段がある。これは、それらのフィールド又はエントリをパケット処理の結果として更新するために使用することができる。新しいエントリがフェッチされると、配置されているエントリは、グローバルセッションメモリに保存される。ローカルセッションキャッシュは、キャッシュ原理に排他的に従うことができ、従って、複数プロセッサの複合システムは、セッションの状態に損害を与えるいかなる競合条件ももたらすことはない。MESI プロトコルのような他のキャッシュプロトコルは、類似の結果を達成するのに使用することができる。プロセッサ複合システムでセッションエントリがキャッシュされ、別のプロセッサ複合システムがそのエントリを必要とする場合、アルゴリズムに基づいて独占的アクセス又は適切なキャッシュ状態を使用し、このエントリは、新しいプロセッサに転送される。セッションエントリはまた、グローバルセッションメモリに書き込まれる場合がある。TCP/IP プロセッサはまた、作動されている接続に対して TCP 状態をウォークスルーするために、TCP 状態機械 (ブロック 2322) を含む。この状態機械は、新しく受け取ったパケットから、状態に影響する適切なフィールドと共にセッションエントリに保存された同じ情報を受け取る。これによ

10

20

30

40

50



って、状態移行があり、情報がセッションテーブルエントリで更新されている場合は、状態機械は、次の状態を生成することができるようになる。TCP/IPプロセッサはまた、フレーム情報を抽出して順序不揃いのパケット実行の作動を行うために使用されるフレームコントローラ/順序不揃いマネージャブロック2318を含む。このブロックはまた、図24の2417で表されるRDMA機構を含むことができるが、非ストレージデータ転送のために使用される。当業者はまた、本特許の教示を使用して、TCP/IPプロセッサのどこにでもRDMA機構を実施することができる。このアーキテクチャは、上部層フレーム指示機構を作り出し、この機構は、プログラマブルフレームコントローラによって使用されるフレーム指示キー又は他のキーとしてパケットCRCを使用することができ、パケットの順序が狂って到着した時でも埋め込まれたPDUを抽出し、それらが最後のバッファの宛先に導かれるようにする。このユニットは、セッションデータベースと対話し、順序が狂って到着した情報を処理し、この情報は記録され、従って、中間セグメントが到着した状態で再送信が回避される。パケットがTCP/IPプロセッサを通じて処理された状態で、パケットがストレージデータ転送に属し、特定の実施がストレージエンジンを含む場合は、それは、作動のためにストレージエンジンに配信され、そうでなければ、最終のバッファの宛先に対するDMAの処理のために、パケットは、ホストプロセッサのインタフェースか、又はブロック1708のストレージフロー/RDMAコントローラに通される。パケットはまた、パケット上の任意の付加的な処理のために、パケットプロセッサブロックに転送することができる。これは、TCP/IPプロセッサ及びストレージプロセッサにより、処理の前か後にパケット上で実行することができるアプリケーション及びカスタムアプリケーションコードを含むことができる。ホストから出力メディアインタフェースへのデータ転送はまた、TCP/IPプロセッサを通過し、データ周辺に作り出されるべき適切なヘッダを形成し、フレームコントローラ及び/又はストレージプロセッサと協働し、適切なデータ分割を実行し、同じくセッション状態を更新する。このデータは、作動のためにスケジューラによってパケットプロセッサに対してスケジュールに入れられたホストコマンド又は受信ネットワークパケットの結果として検索することができる。内部バス構造及び機能性ブロック相互接続は、性能やダイコスト要件などに関して説明したものと異なる場合がある。例えば、ホストコントローラインタフェース2301、スケジューラインタフェース2307、及びメモリコントローラインタフェース2323は、バスコントローラの一部とすることができ、このバスコントローラは、スケジューラ、又はストレージフロー/RDMAコントローラ、又はホスト、又はセッションコントローラ、又は、以下に限定されるものではないが、キュリティプロセッサ、メディアインタフェースユニット、ホストインタフェース、スケジューラ、分類プロセッサ、パケットバッファ又はコントローラプロセッサ、又は以上の任意の組合せのような他のリソースへの又はからのデータパケット又は状態情報又はコマンド又はその組合せの転送を可能にするものである。

#### 【0056】

図24は、図22のIPストレージプロセッサをより詳細に示す。ストレージエンジンは、一般的に、通常のRISC様のパケット処理用命令セットと共にIPベースのストレージを意図した命令セットを有するプログラマブルエンジンである。IPストレージプロセッサエンジンは、CRC作動を実行するためのブロック2411を含む。このブロックは、CRCの生成と確認を可能にする。IPストレージを使用して着信パケットは、TCP/IPエンジンからDMA(ブロック2402及び2408)を通過して、データメモリ(そのようなメモリの例がデータRAMである)(ブロック2409)に転送される。実施がTCP/IPエンジン又はパケットプロセッサエンジン又はその組合せを含まない場合は、パケットを例えばスケジューラから直接受け取ることができる。接続に関連するTCPセッションデータベース情報は、必要に応じてローカルセッションキャッシュから検索することができるか、又はTCP/IPエンジンからパケットを使用して受け取ることができる。ストレージPDUは、PDU分類子エンジン(ブロック2418)に提供され、このエンジンは、PDUを適切なコマンドに分類し、このコマンドは、次に、適切なス

10

20

30

40

50

トレイジコマンド実行エンジン（ブロック 2 4 1 2）を呼び出すために使用される。コマンド実行は、R I S C 又はピア又は命令セットを使用し、又は専用ハードウェアエンジンを使用して達成することができる。コマンド実行エンジンは、P D U で受け取ったコマンドを実行する。受け取った P D U は、保留書込コマンド又は I P ストレージプロトコルによって要求された他のコマンドに対して、読取コマンドデータ又は R 2 T を含むことができる。これらのエンジンは、書込データをホストインタフェースから検索するか、又は読取データを宛先バッファに誘導する。ストレージセッションデータベースは、プロセッサによって供給された最近の又は頻繁な接続に対してローカルメモリ（ブロック 2 4 2 0）として見なされるものにローカルにキャッシュされる。コマンド実行エンジンは、コマンドを実行し、ストレージ状態機械（ブロック 2 4 2 2）及びセッションマネージャ（ブロック 2 4 2 1）と協働してストレージデータベースエントリの更新を行う。接続 I D は、セッションを識別するために使用され、セッションがキャッシュに存在しない場合は、ストレージセッション検索エンジン（ブロック 2 4 1 9）により、図 1 7 のグローバルセッションメモリ 1 7 0 4 から検索される。イニシエータからターゲットへのデータ転送のために、プロセッサは、分割コントローラ（ブロック 2 4 1 3）を使用し、パスマ T U などのような様々なネットワーク制約毎にデータユニットをセグメントに分割する。分割コントローラは、発信 P D U が接続に対して最適なサイズになることを保証するように試みる。要求されたデータ転送が分割の最大有効サイズよりも大きい場合は、分割コントローラは、データを複数のパケットに圧縮し、シーケンスマネージャ（ブロック 2 4 1 5）と協働してシーケンス番号を適切に割り当てる。分割コントローラ 2 4 1 3 はまた、図 2 3 の T C P / I P プロセッサ内で実行することができる。換言すれば、分割コントローラは、このプロセッサが T C P / I P 作動のために使用されてストレージ作動でない時には、図 2 3 のシーケンス/ウィンドウ作動マネージャ 2 3 1 5 の一部である。当業者は、本特許の教示を使用して、分割コントローラを T C P / I P プロセッサに組み込むための代替実施形態を容易に示唆することができる。図 2 4 のストレージプロセッサ（又は、図 2 3 の T C P / I P プロセッサ）はまた、この R D M A 機構を使用して実施されるストレージ又はネットワークデータ転送のために P D U 内で受け取った遠隔の直接アクセス命令を変換する R D M A エンジンを含むこともできる。例えば、図 2 4 では、それは R D M A エンジン 2 4 1 7 である。図 2 3 の T C P / I P プロセッサでは、R D M A エンジン、フレームコントローラ及び順序不揃いマネージャ 2 3 1 8、又は他の適切な要素の一部とすることができる。接続の両端がデータ転送の R D M A モードに合致する場合は、実質的なホストの介入なしにターゲット及びイニシエータ間のデータ転送をスケジュールに入れるために R D M A エンジンが使用される。R D M A 転送状態は、セッションデータベースエントリに維持される。このブロックは、データの周囲で階層化された R D M A のヘッダを作り出し、R D M A 使用可能の接続上で受け取った受信パケットからこれらのヘッダを抽出するためにも使用される。R D M A エンジン、メッセージ/命令を通過させることにより、ストレージフロー/R D M A コントローラ（1 7 0 8）及びホストインタフェースコントローラ（1 7 1 0）と協働し、実質的なホストの介入なしに大量のブロックデータの転送を達成する。I P プロセッサのストレージフロー/R D M A コントローラブロック（1 7 0 8）の R D M A エンジン、要求された作動に対して保護検査を実行し、また、R D M A 領域識別子からホスト空間の物理的又は仮想アドレスへの変換も提供する。この機能性はまた、選択された実行に基づく S A N パケットプロセッサのストレージエンジンの R D M A エンジン（ブロック 2 4 1 7）によって提供することができる。2 4 1 7 と 1 7 0 8、及び他の同様のエンジンの間の R D M A 機能の分布は、本特許を使用して当業者が行うことができる実施の選択肢である。発信データは、P D U クリエータ（ブロック 2 4 2 5）によって標準ベースの P D U にパッケージ化される。P D U のフォーマット設定はまた、パケット処理命令を使用することによって達成することができる。図 2 4 のストレージエンジンは、図 2 3 の T C P / I P エンジン及び図 1 7 のパケットプロセッサエンジンと協働し、両方の方向、すなわちイニシエータからターゲット、ターゲットからホスト、及びその逆にデータ及びコマンド転送を伴う I P ストレージ作動を実行する。換言すれば

10

20

30

40

50

、ホストコントローラインタフェース2401及び2407は、コマンド又はデータ又はその組合せをホストプロセッサへ保存し、又はそれから検索する。これらのインタフェースは、直接又は中間の接続を通じてホストに接続される。2つの装置として示されているが、インタフェース2401及び2407は、単一の装置として実施することができる。これらのブロックを通るデータフローは、その転送の方向に基づいて異なるであろう。例えば、コマンド又はデータがホストからターゲットに送られる場合、ストレージ処理エンジンは、PDUをフォーマット設定するために最初に呼び出されることになり、次に、このPDUがTCPプロセッサに通され、PDUを有効なTCP/IPセグメントにパッケージ化する。しかし、受け取ったパケットは、ストレージプロセッサエンジンに対してスケジュールに入れられる前に、TCP/IPエンジンを通じて進むことになる。内部バス構造及び機能性ブロック相互接続は、性能やダイコスト要件などのために図示のものとは異なる場合がある。例えば、図23と同様に、ホストコントローラインタフェース2401及び2407、及びメモリコントローラインタフェース2423は、バスコントローラの一部とすることができ、このバスコントローラは、スケジューラ、又はホスト、又はストレージフロー/RDMAコントローラ、又はセッションコントローラ、又は、以下に限定されるものではないが、セキュリティプロセッサ、又はメディアインタフェースユニット、ホストインタフェース、スケジューラ、分類プロセッサ、パケットバッファ、又はコントローラプロセッサ、又は以上の任意の組合せのような他のリソースへ又はそれからのデータパケット又は状態情報又はコマンド又はその組合せの転送を可能にするものである。

【0057】

一例として図24のiSCSIプロセッサのようなIPストレージプロセッサにより、図23のTCP/IPプロセッサを組み込まないチップ上でストレージが行われる用途では、TCP/IPインタフェース2406は、IPストレージプロセッサによって処理するIPストレージパケットをスケジュールに入れるためのスケジューラに対するインタフェースとして機能するであろう。本特許の開示を鑑みる時、同様の変形は、十分に当業者の知識の範囲内である。

【0058】

図25は、図17の出力待ち行列コントローラブロック1712をより詳細に示す。このブロックは、図16のネットワークメディアの独立インタフェース1601に送る必要があるパケットを受け取る。パケットは、送出される前に暗号化する必要があるか否かを示すためにタグ付けすることができる。コントローラは、待ち行列2511及びセキュリティエンジンインタフェース2510を通じてセキュリティエンジンに保護する必要があるパケットを待ち行列に入れる。暗号化されたパケットは、セキュリティエンジンから受け取られ、ブロック2509で待ち行列に入れられ、それらの宛先に送られる。このような機構がサポートされている場合は、出力待ち行列コントローラは、それらの各サービス品質(QoS)待ち行列にパケットを割り当てることができる。プログラマブルパケット優先度セレクタ(ブロック2504)は、送られる次のパケットを選択し、そのパケットを適切なポート、すなわち、ポート1...ポートNに対してスケジュールに入れる。ポートに付随するメディアコントローラブロック1601は、パケットを受け入れ、それらをその宛先に送る。

【0059】

図26は、一般的に、図17の1708で表すストレージフロー/RDMAコントローラブロックをより詳細に示す。ストレージフロー及びRDMAコントローラブロックは、コマンド(ストレージ、又はRDMA、又はソケット直接、又はその組合せ)をこのプロセッサに対して待ち行列に入れるためにホストに対して必要な機能性を提供し、このプロセッサは、次に、これらのコマンドを取り出してそれらを実行し、主にコマンドが終了する時にホストプロセッサに割り込む。新しいアクティブなコマンド待ち行列(ブロック2611及び2610)と終了待ち行列(ブロック2612)とは、部分的にオンチップにすることができ、又は部分的にホストメモリ領域又はIPプロセッサに付随するメモリに入れることができ、それらからコマンドがフェッチされるか又は終了ステータスが保存さ

10

20

30

40

50

れる。RDMAエンジン(ブロック2602)は、遠隔の直接メモリアクセスを使用可能にするのに必要な様々な機能を提供する。それは、RDMA領域及びアクセスキーのような情報と仮想アドレス変換機能性とを含むRDMA検索テーブル2608のようなテーブルを有する。このブロック2602内部のRDMAエンジンは、データ転送を実行し、受け取ったRDMAコマンドを解釈して、可能であればトランザクションを実行する。ストレージフローコントローラはまた、ターゲットとイニシエータの間でデータ転送が発生した時にスケジュールに入れられた様々なコマンドの進行状態の追跡を維持する。ストレージフローコントローラは、実行のためのコマンドをスケジュールに入れ、また、ホストドライバにコマンド終了情報も提供する。ストレージフローコントローラは、ホストからの新しい要求が入っており、同様にアクティブなコマンドがアクティブコマンド待ち行列に保持されたコマンド待ち行列を提供する。ブロック2601のコマンドスケジューラは、接続が存在しないターゲットのために受信した新しいコマンドを、新しい接続を開始するためにスケジューラに割り当てる。スケジューラ1702は、図17の1711で表される制御プレーンプロセッサを使用して、一般的に図15と図17の1704で表されるセッションキャッシュに接続エントリが移動される点で接続を確立し、ストレージフローコントローラブロック2601内の状態コントローラは、新しいコマンドをアクティブコマンドに移動して、そのコマンドを適切な接続に関連付ける。ブロック2610内のアクティブコマンドは、パケットプロセッサによる作動のために、検索されてスケジューラ(ブロック1702)に送られる。コマンドステータスに対する更新は、再度フローコントローラに提供され、このコントローラは、次に、それをブロック2603を通じてアクセスされるコマンド状態テーブル(ブロック2607)に保存する。2601のシーケンサは、コマンドスケジューラ設定に対してプログラマブル優先度を適用し、従って、アクティブコマンドと新しいコマンドからスケジュールに入れられるべき次のコマンドを選択する。フローコントローラはまた、着信コマンド(ブロック2613)のための新しい要求の待ち行列を含む。新しい要求は、ホスト上でホストドライバによって適切な処理とバッファ予約が行われた状態で、アクティブコマンド待ち行列に転送される。実行のためにコマンドがスケジュールに入れられると、状態コントローラ2601は、ホストデータ事前取り出しマネージャ(ブロック2617)により、ホストインタフェースブロック2707のDMAエンジンを使用して、ホストメモリからデータの事前取り出しを開始し、従って、コマンドが実行される時に、パケットプロセッサ複合システムに対して提供する準備が完了したデータを保存する。出力待ち行列コントローラ(ブロック2616)は、ホストコントローラインタフェース(ブロック2614)と協働して、データ転送を可能にする。ストレージフロー/RDMAコントローラは、解決されたターゲット/イニシエータと、迅速な検索のため及びアクティブ接続にコマンドを関連付けるために確立された接続とを関連付けるターゲット-イニシエータテーブル(ブロック2609)を維持する。コマンドシーケンサはまた、実行されているコマンドがRDMAコマンドの場合か、又はストレージ転送が接続開始時にRDMA機構を通じて実行されるように交渉された場合は、RDMAエンジンと協働することができる。RDMAエンジン2602は、上述のように、複数RDMA領域、アクセス制御キー、及び仮想アドレス変換ポインタを受け入れるための機能性を提供する。ホストアプリケーション(ユーザアプリケーション又はOSカーネル機能、ストレージ、又は、ウェブページ又はビデオファイルなどのダウンロードのような非ストレージとすることができる)は、それが、開示されたプロセッサを用い、関連ホストドライバによって提供されたサービスを通じてRDMAトランザクションでの使用を希望するメモリ領域を登録する。これが行われた状態で、ホストアプリケーションは、その遠隔端末のピアにこの情報を通信する。この時点で、遠隔マシン又はホストは、実質的にホストの介在なしに、両端のRDMAブロックによって供給されたRDMAコマンドを実行することができる。RDMA転送は、領域からの読取、特定オフセットを有するある一定数のバイト、又は属性が類似の書込のような作動を含むことができる。RDMA機構はまた、2つの端末ノード間での通信パイプの作成のために有用な送信機能性を含むことができる。これらの機能は、クラスター内のサーバ、又は、更に可能性があるのは、2つの

10

20

30

40

50

異なるクラスターのサーバ内のサーバ、又は他のそのようなクラスター化システムで実行される2つのアプリケーションのバッファ間で大量のデータ転送が要求される、アプリケーションのクラスター化において有用である。ストレージデータの転送は、それが実質的にホストの介在なしに大きなブロックのデータ転送を可能にするので、RDMA機構を使用しても達成される。両端のホストは、RDMA転送を行なうこと、及び共有されたアクセス制御キーを通じてメモリ領域と認可を割り当てることに同意するために最初に関わり合っている。更に、2つのノード間のデータ転送は、利用可能なバッファスペースとバッファ転送クレジットが2つの端部ノードによって維持されている限りは、ホストプロセッサの介入がなくても継続することができる。ストレージデータ転送プロトコルは、RDMAプロトコルを使用することに同意し、それを両端で使用可能にすることにより、RDMAの上で実行されるであろう。図26のストレージフローコントローラ及びRDMAコントローラは、次に、RDMAコマンドを使用して、ストレージコマンド実行とデータ転送を遂行することができる。予想されたデータ転送が終了すると、終了待ち行列2612を使用して、ストレージコマンド終了ステータスがホストに通信される。ネットワークから到着した着信データパケットは、図17のパケットプロセッサ複合システムによって処理され、次にPDUが抽出され、ストレージ/RDMAデータパケットの場合、図26のフローコントローラに呈示される。これらは、次に着信待ち行列ブロック2604に割り当てられ、受信バッファのメモリ記述子の検索によって端末宛先バッファに転送され、次に、ホストインタフェースブロック2707のDMAエンジンを使用してDMAを実行する。RDMAコマンドはまた、RDMA初期化のように保護キー検索及びアドレス変換を通過することができる。

10

20

以上の内容はまた、RDMA能力、又はRDMA機構、又はRDMA機能の一部と考えることができる。

#### 【0060】

図27は、図17の1710で表されるホストインタフェースコントローラをより詳細に示す。ホストインタフェースブロックは、ホストバスに対して物理的なインタフェースを提供するホストバスインタフェースコントローラ(ブロック2709)を含む。ホストインタフェースブロックは、スイッチ又はゲートウェイ又は同様な構成に組み込まれる時は、システムアーキテクチャに依存し、ファブリックインタフェースとして又はメディア独立インタフェースとして実行することができる。ブロック2708の処理コントローラ部分は、様々なバス処理を実行し、それらのステータスを維持して終了するために要求された処理を行う。ホストコマンドユニット(ブロック2710)は、ホストによって配送されたコマンドを実行するために、ホストバス構成レジスタと1つ又はそれ以上のインタープリタとを含む。ホストドライバは、ホスト出力待ち行列インタフェース2703上でこれらのコマンドをこのプロセッサに提供する。コマンドは、構成レジスタやスケジュールDMA転送の設定、必要に応じてDMA領域及び許可の設定、セッションエントリや検索セッションデータベースの設定、及びRDMAエンジンの構成などのような様々な機能を提供する。ストレージ及び他のコマンドはまた、IPプロセッサによる実行のためにこのインタフェースを使用して転送することができる。

30

40

#### 【0061】

図28は、図17のセキュリティエンジン1705をより詳細に示す。図示のセキュリティエンジンは、例えば、IPSECのような基準によって要求されるような認証及び暗号化及び暗号解読サービスを提供する。セキュリティエンジンによって提供されるサービスは、複数の認証及びセキュリティアルゴリズムを含むことができる。セキュリティエンジンは、上述のように、プロセッサをオンボードにすることができ、又は別のシリコンチップの一部とすることができる。IPセキュリティサービスを提供する外部セキュリティエンジンは、着信パケットのパケット処理の最初のステージの1つとして、かつ、出ていくパケットの最後のステージの1つとして、データフロー内の同様な位置に配置される。図示のセキュリティエンジンは、セキュリティ基準として適合した非常にハードウェア性

50

能効率的なアルゴリズムの最新暗号化基準 (AES) ベースの暗号化及び暗号解読サービスを提供する。このブロックは、一例としてDESや3DESのような他のセキュリティ機能を同様に提供することができる。セキュリティと認証のためにサポートされたアルゴリズムと機能は、シリコンコストと開発コストから影響されたものである。選択されたアルゴリズムはまた、IPストレージ基準によって要求されたものである。認証エンジン (ブロック2803) は、利用可能なアルゴリズムの一例としてSHA-1を含むように示されている。このブロックは、IPセキュリティ基準で規定されたようなメッセージのダイジェストと認証機能を提供する。セキュリティ及びメッセージ認証サービスが要求される時、データは、これらのブロックを通過する。ターゲットに向けて出ていく途中のクリアなパケットは暗号化され、次に、必要に応じて適切なエンジンを使用して認証される。受け取った安全なパケットは、反対の順序で同じ段階を進む。安全なパケットは認証され、次に、このブロックのエンジン2803及び2804を使用して暗号解読される。セキュリティエンジンはまた、接続のために確立されたセキュリティコンテキストメモリ (ブロック2809) にセキュリティアソシエーションを維持する。セキュリティアソシエーションは、(安全なセッションのインデックス、セキュリティキー、使用されるアルゴリズム、セッションの現在の状態などを含むことができ)、メッセージ認証及び暗号化/暗号解読サービスを実行するために使用される。メッセージ認証サービス及び暗号化/暗号解読サービスを互いに独立して使用することも可能である。

10

## 【0062】

図29は、全体的に図17の1704で表されるセッションキャッシュ及びメモリコントローラ複合システムをより詳細に示す。メモリ複合システムは、キャッシュ又はメモリ又はその組合せとして、本特許ではセッション/グローバルセッションメモリ又はセッションキャッシュと呼ばれるTCP/IPセッションデータベースのためのキャッシュ/メモリアーキテクチャを含む。セッションキャッシュ検索エンジン (ブロック2904) は、特定のセッションキャッシュエントリを検索するための機能性を提供する。この検索ブロックは、提供されたフィールドの外にハッシュインデックスを作り出すか、又はハッシュキーを受け入れてセッションキャッシュエントリを検索することができる。キャッシュアレイ内にハッシュインデックスとのタグの符合がない場合、検索ブロックは、このキーを使用して外部メモリからセッションエントリを見つけ、現在のセッションキャッシュエントリをそのセッションエントリと取り替える。それは、要求のパケットプロセッサ複合システムにセッションエントリフィールドを提供する。ローカルプロセッサ複合システムキャッシュ内に存在するキャッシュエントリは、グローバルキャッシュに共有されるようにマーク付けされる。従って、任意のプロセッサがこのキャッシュエントリを要求すると、それは、グローバルキャッシュと要求しているプロセッサとに転送され、グローバルキャッシュにそのようにマーク付けされる。セッションメモリコントローラはまた、退去させられたキャッシュエントリをこのブロック内部のグローバルキャッシュに移動させる役割を果たす。従って、セッションエントリの任意の要求者に対して、最新のセッション状態だけがいつでも利用可能である。セッションキャッシュが満杯の場合、新しいエントリにより最低使用頻度のエントリを退去させることができる。セッションメモリは、所定の処理技術で使用することができるシリコンスペースに応じて、単回路又は複数回路キャッシュ又はハッシュインデックスメモリ、又はその組合せにすることができる。ネットワークスイッチ又はルータのためのネットワークングアプリケーションでは、一般的にパケット間で利用可能な基準特性の局所性があまり多くなく、従って、キャッシュの使用はキャッシュミスのために多くの性能改善をもたらさないという点で、セッションデータベースエントリを保存するためのキャッシュの使用は独特である。しかし、ストレージ処理は、2つの端部システム間の持続時間が長い処理であり、大量のデータを交換することができる。クラスター化又はメディアサーバのように2つのノード間で大量のデータ転送が発生するこのシナリオ又はケースでは、キャッシュベースのセッションメモリアーキテクチャは、オフチップメモリからのデータ転送を大幅に低減することにより、性能的な恩典を達成することになる。セッションキャッシュのサイズは、利用可能なシリコンダイ面積の

20

30

40

50

関数であり、トレードオフによって性能に影響を及ぼすことがある。メモリコントローラブロックはまた、パケット、パケットのフラグメント、又はメモリ内のオペレーティングデータを保存する必要がある他のブロックにサービスを提供する。メモリインタフェースは、サポートする必要のある予想データ帯域幅に依存して、単一又は複数の外部メモリコントローラ（ブロック2901）を提供する。これは、DRAM又はSRAM又はRDRAM、又は他の動的又は静的RAM、又はその組合せのための1つ又は複数の2倍データ速度コントローラとすることができる。図は、複数コントローラを表すが、その数は必要な帯域幅とコストに応じて変動する。メモリ複合システムはまた、セッションデータベースメモリブロックによって維持された再発信待ち行列上でセッション自体を待ち行列に入れるセッションに対する再発信タイムアウトに使用するためのタイマ機能性を提供することができる。

10

【0063】

図30は、分類エンジンのデータ構造の詳細を示す。これは、分類エンジンのデータ構造を編成するための1つの方法である。分類データベースは、ツリー構造（ブロック3001）として表され、ツリー内のノード（ブロック3003）とノードに関連したアクション（ブロック3008）とは、分類エンジンがツリーの下方にウォークし、特定のノード値の比較を行うことを可能にする。ノード値とそれらが呈示するフィールドは、プログラム可能である。アクションフィールドは、フィールドが特定のノード値に符合する時に抽出される。アクションアイテムは、次の段階を規定し、この段階は、新しいフィールドの抽出と比較を含むことができ、このノード値の対に関連する特定のデータフィールド上でALU作動のような他の作動を実行するか、又は端末ノードを表示することができ、その時点で特定パケットの分類が終了する。データ構造は、それがパケットスケジューラから受け取るパケットを分類するために分類エンジンによって使用される。符合する値を使用して検索されたアクションアイテムは、パケットの異なるフィールドを繰り返す一方で、結果コンパイラによって使用され、一般的にパケットのヘッダの前でパケットに接続される分類タグを生成する。分類タグは、次に、プロセッサの残りによって基準として使用され、分類の結果に基づいて取る必要があるアクションを決める。プログラマブルな特性を有する分類子は、分類ツリー構造のシステム変更を可能にし、分類の必要性の異なるシステムでプロセッサの使用を可能にする。分類エンジンはまた、分類のツリーノード値アクションの構造の一部としてプログラムすることができるストレージ/ネットワークポリシーの作成を可能にし、IPベースのストレージシステムにおいて非常に強力な機能を提供する。ポリシーは、このプロセッサを使用するシステムの管理を増強し、特定のポリシー又は規則が符合した時又は侵害された時の強化機能を可能にする。分類エンジンは、特定のシステム制約によって要求された時に、外部の要素を使用して分類データベースの拡張を可能にする。ツリーとノードの数は、シリコン面積と性能トレードオフとに基づいて決まる。データ構造の構成要素は、分類エンジンの様々なブロックに維持され、分類シーケンサによって使用され、構造によってパケット分類の方向付けが行われる。分類データ構造は、ターゲットソリューションに応じて、表示されたものよりも多いか又は少ないフィールドを要求することができる。従って、分類のコアの機能性は、基本となるアーキテクチャから離れることなく、より少ない要素と構造を使用して達成することができる。分類処理は、プログラムに従ってツリーとノードを通してウォークする。特定ノードアクションにより、分類の残りのフィールドのために新しいツリーの使用を可能にすることができる。従って、分類処理は、ツリールートで始まり、ノードを通して進行し、最後はリーフノードに到達する。

20

30

40

【0064】

図31は、イニシエータとターゲットの間の読取作動を示す。イニシエータは、処理を開始するためにターゲットにREADコマンド要求（ブロック3101）を送信する。これは、アプリケーション層の要求であり、特定のSCSIプロトコルコマンドにマップされ、IPベースのストレージネットワーク内でREADプロトコルデータユニット（ブロック3102）として移送されるものである。ターゲットは、要求されたデータ（ブロッ

50

ク 3 1 0 3 ) を準備し、最大転送ユニット制約を満足するように分割された読取応答 P D U ( ブロック 3 1 0 5 ) を提供する。イニシエータは、次に、 I P パケットからデータを検索 ( ブロック 3 1 0 6 ) し、それは、次に、この作動のために割り当てられた読取バッファに保存される。全てのデータが転送された状態で、ターゲットは、コマンド終了の応答を行ってステータスを感知する ( ブロック 3 1 0 7 ) 。完全な転送が終了する ( ブロック 3 1 0 9 ) と、イニシエータは、次にコマンドを退去させる。ターゲットに何らかのエラーがあり、何らかの理由でコマンドが中途終了している時は、イニシエータによって単独で回復の手順を開始することができる。このトランザクションは、プロトコルの P D U としての i S C S I のような I P ベースのストレージプロトコル上のデータ移送を有する標準的 S C S I の R E A D トランザクションである。

10

## 【 0 0 6 5 】

図 3 2 は、図 3 1 に示すトランザクションの受信「 R E A D P D U 」の 1 つに対する本発明の I P プロセッサ内部のデータフローを示す。内部データフローは、イニシエータ端末上の I P プロセッサによって受信した読取データ P D U に対して示される。この図は、パケットが通過する作動の様々なステージを表す。ステージは、パケットがトラバースするパイプラインステージと考えることができる。トラバースするパイプのステージの数は、受け取ったパケットの種類に依存する。図は、確立された接続上で受け取ったパケットに対するパイプステージを表す。パケットは、以下の主要なパイプステージを通過してトラバースする。

## 【 0 0 6 6 】

20

1 . ブロック 3 2 0 7 に示す主要な段階を有するブロック 3 2 0 1 の受信パイプステージ。パケットは、メディアアクセスコントローラによって受信される。パケットが検出され、プリアンプル/トレーラーが除去され、パケットは、第 2 層ヘッダとペイロードと共に抽出される。これは、第 2 層確認、及び全てのエラー検出が意図する受取人に対して発生するステージである。確立されたポリシー毎に適用されるサービス品質検査がある場合もある。パケット確認がクリアされた状態で、パケットは、入力待ち行列に入れられる。

2 . ブロック 3 2 0 8 に示す主要な段階を有するブロック 3 2 0 2 のセキュリティパイプステージ。パケットは、入力待ち行列から分類エンジンまで動かされ、そこでセキュリティ処理のための迅速な判断が行われ、パケットがセキュリティ処理を通過する必要がある場合、セキュリティパイプステージに入る。パケットがクリアなテキストで受信されて認証を必要としない場合、セキュリティパイプステージは省略される。セキュリティパイプステージはまた、セキュリティエンジンが I P プロセッサと統合されていない場合は、省略することができる。パケットは、第 1 にこの接続に対するセキュリティアソシエーションがメモリから検索され、パケットが選択されたメッセージ認証アルゴリズムを使用して認証されるセキュリティエンジンの様々なステージを通過する。パケットは、次に、セッションのために確立されたセキュリティキーを使用して暗号解読される。パケットがクリアなテキストになった状態で、それは、入力待ち行列コントローラに対して再度待ち行列に入れられる。

30

3 . ブロック 3 2 0 9 に示す主要な段階を有するブロック 3 2 0 3 の分類パイプステージ。スケジューラは、入力待ち行列からクリアなパケットを検索し、分類のためにパケットをスケジュールに入れる。分類エンジンは、第 3 層及びより高い層の分類のためにパケットから関連フィールドを抽出するような様々なタスクを実行し、 T C P / I P / ストレージプロトコルなどを識別し、それらの分類タグを生成し、分類エンジンでプログラムされたポリシーに応じて、パケットを排除したりバイパスのためにパケットにタグ付けするようなアクションを取ることができる。分類エンジンはまた、実行を容易にするためのパケットヘッダ及びペイロードのマーク付けと共に、セッション又はこのエンジンが属するフローでパケットにタグ付けすることができる。分類エンジンのプログラムに依存して、列挙されたタスクのいくつかを実行することができ、又は実行されなくてもよく、他のタスクが実行されてもよい。分類が実行されると、パケットに分類タグが追加され、処理するためにスケジューラに対してパケットが待ち行列に入れられる。

40

50



4. ブロック3210に示す主要な段階を有するブロック3204のスケジュールパイプステージ。分類パケットが、分類エンジン待ち行列から検索され、処理のためにスケジューラ内に保存される。スケジューラは、パケットヘッダから発信元及び宛先フィールドのハッシュを実行し、分類子によって実行されない場合にパケットが属するフローを識別する。フローの識別が行われた状態で、パケットは、フローの従属性に基づいて実行リソース待ち行列に割り当てられる。新しいパケットを受け入れるためにリソースが利用可能になると、待ち行列の次のパケットが、そのリソースに対する実行のために割り当てられる。

5. ブロック3211に示す主要な段階を有するブロック3205の実行パイプステージ。パケットは、このパケットを実行するためにリソースが利用可能になると、実行パイプステージに入る。パケットは、パケットを実行するために想定されたパケットプロセッサ複合システムに転送される。プロセッサは、パケットに接続された分類タグを見て、パケットのために要求される処理段階を決める。これがIPベースのストレージパケットの場合、このセッションのためのセッションデータベースエントリが検索される。ローカルセッションキャッシュが既にセッションエントリを保持する場合は、データベースアクセスは要求されない。フローに基づいてパケット割当が行われた場合は、グローバルセッションメモリからセッションエントリを検索する必要はない。パケットプロセッサは、次に、TCPエンジン/ストレージエンジンを開始し、それらの作動を実行する。TCPエンジンは、チェックサム、シーケンス番号検査、必要なCRC作動を使用したフレーム指示検査、及びTCP状態更新を含む様々なTCP検査を実行する。次に、ストレージPDUが抽出され、実行のためにストレージエンジンに割り当てられる。ストレージエンジンは、PDU内のコマンドを変換し、この特定のケースでは、それをアクティブなセッションに対する読取応答として識別する。次に、それは、ペイロードの整合性とシーケンスの整合性を確認し、次に、セッションデータベースエントリ内のストレージフロー状態を更新する。宛先バッファのメモリ記述子はまた、セッションデータベースエントリから検索され、抽出されたPDUペイロードはストレージフロー/RDMAコントローラに対して、それらのホストインタフェースブロックはDMAに対して、データは最終のバッファ宛先に対して、待ち行列に入れられる。データは、メモリ記述子と実行するコマンド/作動と共に、フローコントローラに配信することができる。この場合、このアクティブ読取コマンドに対してデータを入れる。ストレージフローコントローラは、そのアクティブコマンドデータベースを更新する。実行エンジンは、パケットが既に退去しており、パケットプロセッサ複合システムがそのシステムの次のコマンドを受け取る準備ができていないことをスケジューラに示す。

6. ブロック3212に示す主要な段階を有するブロック3206のDMAパイプステージ。ストレージフローコントローラがメモリ記述子、コマンド、及びフロー状態の適切な確認を行った状態で、このコントローラは、ホストメモリに対する転送のためにDMAエンジンにデータブロックを送る。DMAエンジンは、そのようなQOS機構がプログラムされるか又は実行されると、優先度ベースの待ち行列を実行することができる。データは、DMAを通じてホストメモリ位置に転送される。これがコマンドの最後の作動の場合、ホストドライバにコマンド実行の終了が表示される。これがコマンドに対する最後の作動であり、コマンドが終了待ち行列に対して待ち行列に入れられている場合、コマンドに対して割り当てられたリソースは、新しいコマンドを受け入れるために解放される。コマンド統計データは、性能分析、ポリシー管理、他のネットワーク管理、又は統計目的のために要求される場合があるので、終了ステータスと共に収集して転送することができる。

【0067】

図33は、イニシエータとターゲットの間の書込コマンド作動を示す。イニシエータは、トランザクションを開始するためにターゲットにWRITEコマンド(3301)を送る。このコマンドは、IPストレージネットワーク上で「WRITE PDU」(ブロック3302)として移送される。レシーバは、新しい要求待ち行列に受け取ったコマンドを入れる。作動の古いコマンドが終了した状態で(ブロック3304)、レシーバは、リ

10

20

30

40

50

ソースを割り当ててコマンド（ブロック 3305）に対応するWRITEデータを受け入れる。このステージで、レシーバは、受け取る予定のデータ量とどの位置から受け取るかの指示を使用して、イニシエータに転送待機中（R2T）PDU（ブロック 3306）を出す。イニシエータは、R2T要求のフィールドを変換し、データパケット（ブロック 3307）を、受け取ったR2T毎にレシーバに送る。イニシエータとターゲットの間のこの交換のシーケンスは、コマンドが終了されるまで続く。成功したコマンドの終了又はエラー条件は、応答PDUとしてターゲットによってイニシエータに通信され、それが次にコマンドを終了させる。イニシエータは、エラーの場合に回復処理を開始するように要求することができる。これは、図33の交換には示されていない。

#### 【0068】

図34は、「R2T PDU」の1つに対する本発明のIPプロセッサ内部のデータフローと、図33に示す書込トランザクションの追従する書込データとを示す。イニシエータは、ネットワークメディアインタフェースを通じてR2Tパケットを受け取る。パケットは、受信、セキュリティ、分類、スケジュール、及び実行を含む図32の「READ PDU」と同様の対応するブロック3415、3416、3409、及び3410内の詳細な主要段階を使用して、全てのステージ（ブロック3401、3402、3403、及び3404）を通過する。セキュリティ処理は、この図には表示されていない。これらのステージに続いてR2Tは、図34に示すDMAステージを使用して、書込データフェッチを開始する（ブロック3405及び3411）。書込データは、次に分割され、実行ステージ（ブロック3406及び3412）を通じてTCP/IPパケットに置かれる。TCP及びストレージセッションDBエントリは、R2Tに回答して転送されたデータを用いてWRITEコマンドに対して更新される。パケットは、次に、出力待ち行列コントローラに対して待ち行列に入れられる。接続セキュリティの合意に基づいて、パケットは、セキュリティパイプステージ（ブロック3407及び3413）に入ることができる。パケットが暗号化されてメッセージ認証コードが生成された状態で、パケットは、宛先への送信のためにメディアインタフェースに対して待ち行列に入れられる。このステージ（ブロック3408及び3414）の間に、パケットは、パケットプロセッサによってまだ為されていなければ、第2層ヘッダ内にカプセル化され、かつ送信される。パイプラインの各ステージで従うべき段階は、上述の「READ PDU」パイプステージと類似であり、この図に示されている書込データパケットステージのための付加的なステージを有する。各ステージで実行される特定の作動は、コマンドの種類、セッションの状態、コマンドの状態、及び設定することができるポリシーに対する他の様々な構成に依存する。

#### 【0069】

図35は、イニシエータとターゲットの間のRDMA機構を使用したREADデータ転送を示す。イニシエータとターゲットは、RDMAデータ転送を開始させる前にRDMAバッファを登録する（ブロック3501、3502、及び3503）。イニシエータは、予想される受け取り側としてのRDMAを使用してREADコマンド（ブロック3510）を出す。このコマンドは、ターゲット（ブロック3511）に移送される。ターゲットは、読み込まれるデータを準備し（ブロック3504）、次に、RDMA書込作動（ブロック3505）を実行し、書込データをイニシエータにおいてホストの介在なしに直接RDMAバッファに保存する。作動の終了は、コマンド終了応答を使用して示される。

#### 【0070】

図36は、READコマンドフローを実行するRDMA書込パケットの内部アーキテクチャデータフローを示す。RDMA書込パケットはまた、ネットワークインタフェース上で受け取った他の任意の有効なパケットと同じパイプステージに従う。このパケットは、受信パイプステージ（ブロック3601及び3607）内の第2層処理を通過し、そこから、それは、セキュリティ処理の必要性を検出するためにスケジューラに対して待ち行列に入れられる。パケットが暗号解読される必要がある場合、それは、セキュリティパイプステージに入る（ブロック3602及び3608）。暗号解読されたパケットは、次に、プログラムされていた分類タスクを実行するために、分類エンジンに対してスケジュール

10

20

30

40

50

に入れられる（ブロック3603及び3609）。分類が終了した状態で、タグ付けされたパケットは、スケジュールパイプステージ（ブロック3604及び3610）に入り、そこで、スケジューラは、フローベースのスケジュール設定に依存して、このパケットをリソース特定の待ち行列に割り当てる。意図されたりソースがこのパケットを実行する準備が完了している場合、それは、そのパケットプロセッサ複合システム（ブロック3605及び3611）に転送され、そこで、全てのTCP/IP確認、検査、及び状態更新が行われ、PDUが抽出される。次に、ストレージエンジンは、RDMAを使用して実施されたストレージPDUに対するストレージフローに属するとしてPDUを識別し、RDMAコマンドを解釈する。この場合、特定のRDMAバッファに書き込むのはRDMAである。このデータは、抽出されて、RDMA領域変換及び保護検査を行うストレージフロー/RDMAコントローラブロックに通され、パケットは、ホストインタフェース（ブロック3606及び3612）を通じてDMAに対して待ち行列に入れられる。パケットがパケットプロセッサ複合システムを通じた作動を終了した状態で、スケジューラに通知され、パケットは、スケジューラ内に保持された状態から退去させられる。DMAステージでRDMAデータ転送が終了した状態で、これがストレージコマンド実行を終了させる最後のデータ転送の場合、そのコマンドは退去させられ、コマンド終了待ち行列に割り当てられる。

10

#### 【0071】

図37は、RDMA読取作動を使用するストレージ書込コマンド実行を示す。イニシエータ及びターゲットは、RDMAコントローラを使用して最初にそれらのRDMAバッファを登録し、次に、同じくバッファをピアに通知する。次に、イニシエータは、ターゲットに書込コマンド（ブロック3701）を出し、それは、IPストレージPDUを使用して移送される。受け取り側は、最初にRDMAバッファを割り当てることによって書込コマンドを実行し、書込を受け取り、次にイニシエータに対してRDMA読取を要求する（ブロック3707及び3708）。イニシエータから書き出されるデータは、次に、RDMA読取応答パケット（ブロック3707及び3708）として供給される。受け取り側は、いかなるホストの介入もなく、パケットを直接RDMAバッファに保存する。読取要求がセグメントサイズよりも大きなデータに対するものだった場合、READ要求に回答して、複数のREAD応答PDUがイニシエータに送られるであろう。データ転送が終了した状態で、イニシエータに終了ステータスが移送され、ホストに対してコマンド終了が示される。

20

30

#### 【0072】

図38は、図37に示すフロー処理の1つのセクションに対するRDMA読取要求及び得られる書込データ転送のデータフローを示す。データフローは、図34の書込データフローと非常に似ている。RDMA読取要求パケットは、受信、分類、スケジュール、及び実行（ブロック3801、3802、3803、3804、3815、3816、3809、及び3810）を含む様々な処理パイプステージを通して流れる。この要求が実行された状態で、それは、RDMA読取応答パケットを生成する。RDMA応答は、最初にシステムメモリから要求されたデータのDMA（ブロック3805及び3811）を行うことによって生成され、次に、実行ステージ（ブロック3806及び3812）を通してセグメントとパケットを生成する。適切なセッションデータベースエントリが更新され、データパケットは、必要に応じてセキュリティステージに進む（ブロック3807及び3813）。安全又はクリアなパケットは、次に、送信ステージに対して待ち行列に入れられ（ブロック3808及び3814）、このステージは、適切な第2層更新を実行してパケットをターゲットに送信する。

40

#### 【0073】

図39は、イニシエータから開始されたストレージコマンドに対するイニシエータコマンドフローをより詳細に示す。図示のように、以下は、コマンドが従う主要段階のいくつかである。

1. ホストドライバは、コマンドをストレージフロー/RDMAコントローラのプロセ

50

ッサコマンド待ち行列に入れる。

2. ホストは、作動のため及びリソースを予約するためにコマンドが正常にスケジュールに入れられたかを通知される。

3. ストレージフロー/RDMAコントローラは、ターゲットに対する接続が確立された場合、作動のためのコマンドをパケットスケジューラに対してスケジュールに入れる。そうでなければ、コントローラは、ターゲットセッション開始を開始し、セッションが確立された状態で、コマンドは、パケットスケジューラに対してスケジュールに入れられる。

4. スケジューラは、このコマンドを受け入れる準備が完了したSANパケットプロセッサの1つにコマンドを割り当てる。

5. プロセッサ複合システムは、セッションエントリのために、要求をセッションコントローラに送る。

6. セッションエントリが、パケットプロセッサ複合システムに供給される。

7. パケットプロセッサは、PDUとしてコマンドを運び、出力待ち行列に対してスケジュールに入れられるパケットを形成する。

8. コマンドPDUは、それをターゲットに送るネットワークメディアインタフェースに与えられる。

これは、イニシエータとターゲットの間で接続が確立された時に、主にイニシエータからターゲットまで大部分のコマンドが追従する高レベルのフローである。

#### 【0074】

図40は、読取パケットデータフローをより詳細に示す。ここでは、読取コマンドは、最初に図39に示すフローと同様のフローを使用し、イニシエータからターゲットに送られる。ターゲットは、図40に示すようなフローに従う読取応答PDUをイニシエータに送る。図示のように、読取データパケットは、以下の主要段階を通過する。

1. 入力パケットは、ネットワークメディアインタフェースブロックから受信される。

2. パケットスケジューラは、入力待ち行列からパケットを検索する。

3. パケットは、分類のためにスケジュールに入れられる。

4. 分類されたパケットは、分類タグと共に分類子から戻る。

5. 分類とフローベースのリソース割当てに基づいて、パケットは、パケット上で作動するパケットプロセッサ複合システムに割り当てられる。

6. パケットプロセッサ複合システムは、セッションキャッシュ内のセッションエントリを検索する(ローカルに存在しない場合)。

7. セッションキャッシュエントリが、パケットプロセッサ複合システムに戻される。

8. パケットプロセッサ複合システムは、TCP/IP作動/IPストレージ作動を実行し、ペイロード内の読取データを抽出する。MDL(メモリ記述子リスト)のような適切な宛先タグを有する読取データが、ホストインタフェース出力コントローラに提供される。

9. ホストDMAエンジンは、システムバッファメモリに読取データを転送する。

これらの段階のいくつかは、安全なパケットフローが表された図32により詳細に与えられ、一方、図40は、クリアテキスト読取パケットフローを表す。本特許に示すこのフロー及び他のフローは、本特許の教示により当業者が容易に実施することができる開示されたプロセッサの適切なリソースを使用することにより、ストレージ及び非ストレージデータ転送に適用可能である。

#### 【0075】

図41は、書込データフローをより詳細に示す。書込コマンドは、図39のフローと同様のフローに従う。イニシエータは、書込コマンドをターゲットに送る。ターゲットは、ターゲットが指定量のデータを受け取る準備が完了していることをイニシエータに示す転送待機中(R2T)PDUによりイニシエータに応答する。イニシエータは、次に、要求されたデータをターゲットに送る。図41は、イニシエータからターゲットまで要求された書込データパケットが従うR2Tを示す。このフローで追従される主な段階は、以下の

通りである。

- 1 . 入力パケットが、ネットワークメディアインタフェースブロックから受信される。
- 2 . パケットスケジューラは、入力待ち行列からパケットを検索する。
- 3 . パケットは、分類のためにスケジュールに入れられる。
- 4 . 分類されたパケットは、分類タグと共に分類子から戻る。
  - a . 分類とフローベースのリソース割当てに依存して、パケットは、パケット上で動作するパケットプロセッサ複合システムに割り当てられる。
- 5 . パケットプロセッサ複合システムは、セッションキャッシュ内のセッションエントリを検索する（ローカルに存在しない場合）。
- 6 . セッションキャッシュエントリが、パケットプロセッサ複合システムに戻される。
- 7 . パケットプロセッサは、「R 2 T P D U」を判断し、ストレージフロー/R D M Aコントローラへの要求で書込データを要求する。
- 8 . フローコントローラは、ホストインタフェースに対してD M Aを開始する。
- 9 . ホストインタフェースは、D M Aを実行し、ホスト入力待ち行列にデータを戻す。
- 1 0 . パケットプロセッサ複合システムは、ホスト入力待ち行列からデータを受け取る。
- 1 1 . パケットプロセッサ複合システムは、有効なP D Uとパケットをデータの周りに形成し、適切なセッションエントリを更新し、出力待ち行列にパケットを転送する。
- 1 2 . パケットは、データパケットを宛先に送信する出力ネットワークメディアインタフェースブロックに転送される。

10

20

#### 【 0 0 7 6 】

図 4 1 のフローは、クリアテキストデータ転送を示す。データ転送が安全である必要がある場合は、フローは、出力データパケットが 1 1 a 及び 1 1 b とラベル付けされた矢印で示すような安全パケットを通るようにルーティングされた図 4 3 に示すフローと同様である。入力 R 2 T パケットはまた、安全である場合に、セキュリティエンジンを通るようにルーティングされるであろう（これは、図示されていない）。

#### 【 0 0 7 7 】

図 4 2 は、パケットが暗号テキストか又は安全である場合の読取パケットフローを示す。このフローは、上述の関連する説明を使用して図 3 2 により詳細に示されている。安全な読取フローとクリアな読取フローの主な相違点は、パケットが分類子によって安全パケットとして最初に分類され、従って、セキュリティエンジンにルーティングされることである。これらの段階は、2 a、2 b、及び 2 c とラベル付けされた矢印で示されている。セキュリティエンジンは、パケットを暗号解読し、メッセージの認証を実行し、矢印 2 d で示すような更なる処理のために入力待ち行列にクリアパケットを転送する。クリアパケットは、次に、スケジューラによって検索され、図 4 2 の 2 e 及び 3 とラベル付けされた矢印で示すように、分類エンジンに供給される。残りの段階と作動は、上述の図 4 0 のものと同様である。

30

#### 【 0 0 7 8 】

図 4 4 は、R D M A バッファ通知フローを示す。このフローは、図 3 9 に示すような他の全てのストレージコマンドフローと非常に似ているように示されている。主な段階で取られる詳細なアクションは、コマンドによって異なる。R D M A バッファ通知とレジストレーションに対して、R D M A 領域 I D が作成され、この領域に対するアドレス変換機構と共に記録される。R D M A レジストレーションはまた、アクセス制御のための保護キーを含み、R D M A 転送のために必要な別のフィールドを含むことができる。コマンドのためのパケットを作り出す段階は、図 3 9 の段階と類似である。

40

#### 【 0 0 7 9 】

図 4 5 は、R D M A 書込フローをより詳細に示す。R D M A 書込は、R D M A 書込を受け取るイニシエータに対する通常の見える。R D M A 書込パケットは、図 4 0 に示す読取 P D U のような同じ主要フローの段階に従う。R D M A 転送は、R D M A アドレス変換と領域アクセス制御キー検査とを伴い、他のセッションエントリ以外に

50

R D M A データベースエントリを更新する。主要なフローの段階は、通常の読取応答 P D U と同じである。

【 0 0 8 0 】

図 4 6 は、R D M A 読取データフローをより詳細に示す。この図は、イニシエータによってターゲットから受け取られる R D M A 読取要求と、イニシエータからターゲットに書き出される R D M A 読取データとを示す。このフローは、ストレージ書込コマンドが従う R 2 T 応答と非常に似ている。このフローでは、ストレージ書込コマンドは、R D M A 読取を使用して達成される。パケットが従う主な段階は、図 4 1 に示す R 2 T 書込データフローと基本的に同じである。

【 0 0 8 1 】

図 4 7 は、セッション作成フローの主要な段階を示す。この図は、イニシエータとターゲットの間のセッション開始で要求される、この低速パス作動のための制御プレーンプロセッサの使用を表す。この機能性は、パケットプロセッサ複合システムを通して実施することが可能である。しかし、ここでは、制御プレーンプロセッサを使用して実施されるものとして示されている。どちらの手法も許容することができる。以下は、セッション作成中の主要な段階である。

- 1 . コマンドが、ホストドライバによってスケジュールに入れられる。
- 2 . ホストドライバは、コマンドがスケジュールに入れられ、ホストによって要求された全ての制御情報が通されることを通知される。
- 3 . ストレージフロー / R D M A コントローラは、セッションが存在しないターゲットにコマンドを送る要求を検出し、従って、それは、その要求を制御プレーンプロセッサに通して移送セッションを確立する。
- 4 . 制御プレーンプロセッサは、出力待ち行列に「 T C P S Y N 」パケットを送る。
- 5 . S Y N パケットは、ネットワークメディアインタフェースに送信され、そこから宛先に送信される。
- 6 . 宛先は、S Y N パケットを受け取った後に、S Y N - A C K 応答を使用して応答し、そのパケットは、ネットワークメディアインタフェースから受け取り次第入力待ち行列に入れられる。
- 7 . パケットが、パケットスケジューラによって検索される。
- 8 . パケットが、分類エンジンに通される。
- 9 . タグ付けされて分類されたパケットが、スケジューラに戻される。
- 1 0 . スケジューラは、分類に基づいてこのパケットを制御プレーンプロセッサに転送する。
- 1 1 . プロセッサは、次に、A C K パケットを使用して出力待ち行列に応答する。
- 1 2 . パケットは、次に、最終の宛先に送信され、従って、セッション確立ハンドシェイクを終了する。
- 1 3 . セッションが確立された状態で、この状態は、ストレージフローコントローラに供給される。セッションエントリは、このように作り出され、次に、セッションメモリコントローラに通される（この部分は、図に示されていない）。

セッションを段階 1 3 のような確立状態にする前に、制御プレーンプロセッサは、ストレージプロトコルの完全なロゲイン段階を実行し、パラメータを交換し、それがストレージデータ転送接続の場合は、それらを特定接続のために記録するように要求されることがある。ロゲインが認証されてパラメータ変換が終了した状態で、セッションは、上述の段階 1 3 で示されたセッション確立状態に入る。

【 0 0 8 2 】

図 4 8 は、セッション取外しフローの主要な段階を示す。このフローの段階は、図 4 7 の段階と非常に似ている。2 つのフローの主な違いは、セッション作成のための S Y N 、 S Y N - A C K 、 及び A C K パケットの代わりに、イニシエータとターゲットの間で、F I N 、 F I N - A C K 、 及び A C K パケットが転送されることである。それ以外は、主要段階は、非常に類似している。ここでの別の主な相違点は、適切なセッションエントリが

10

20

30

40

50

作り出されず、セッションキャッシュとセッションメモリから除去されることである。接続の作動統計データは、図に示されていないが、記録されてホストドライバに供給することができる。

#### 【 0 0 8 3 】

図 4 9 は、ターゲットの観点からのセッション作成とセッション取外しの段階を示す。以下は、セッション作成のために従うべき段階である。

1 . イニシエータからの S Y N 要求が、ネットワークメディアインタフェース上で受信される。

2 . スケジューラは、入力待ち行列から S Y N パケットを検索する。

3 . スケジューラは、分類のためにこのパケットを分類エンジンに送る。

4 . 分類エンジンは、適切なタグを付けて分類されたパケットを戻す。

5 . スケジューラは、S Y N パケットとしての分類に基づいて、このパケットを制御ブレンプロセッサに転送する。

6 . 制御ブレンプロセッサは、S Y N - A C K 肯定応答パケットを使用して応答する。それはまた、イニシエータからの非送信請求データ転送のために適切なバッファスペースを割り当てるようにホストに要求する（この部分は、図示されていない）。

7 . S Y N - A C K パケットが、イニシエータに送られる。

8 . イニシエータは、次に、A C K パケットを使用して S Y N - A C K パケットに肯定応答し、3 方向ハンドシェイクを終了する。このパケットは、ネットワークメディアインタフェースで受信され、第 2 層処理の後で入力待ち行列に入れられる。

9 . スケジューラは、このパケットを検索する。

1 0 . パケットが、分類子に送られる。

1 1 . 分類されたパケットは、スケジューラに戻され、3 方向ハンドシェイクを終了するために、制御ブレンプロセッサに提供されるべくスケジュールに入れられる。

1 2 . コントローラは、A C K パケットを取得する。

1 3 . 制御ブレンプロセッサは、この時点で確立状態での接続を有し、セッションキャッシュ内にエントリを作り出すストレージフローコントローラに送る。

1 4 . ホストドライバは、終了したセッション作成を通知される。

セッション確立はまた、図 4 9 に示されていないが、ログイン段階を伴う場合がある。しかし、ログイン段階とパラメータ変換は、完全に構成されて確立された状態にセッションが入る前に起きる。これらのデータ転送とハンドシェイクは、主として制御ブレンプロセッサによって行うことができる。これらの段階が行われた状態で、上述のフローの残りの段階を実行することができる。

#### 【 0 0 8 4 】

図 5 0 及び 5 1 は、ターゲットサブシステムの書込データフローを示す。図 5 0 は、イニシエータからのデータ書込を受け入れる準備が完了していることをイニシエータに通知するためにターゲットによって使用される R 2 T コマンドフローを示す。イニシエータは、次に、ターゲットで受け取った書込を送り、その内部データフローは、図 5 1 に示されている。2 つの図は、両方で 1 つの R 2 T 及びデータ書込の対を示す。以下は、図 4 0 及び 5 1 の両方で示すように追従される主な段階である。

1 . ターゲットホストシステムは、図 3 3 に示すような書込要求の受信に応答して、書込データを受け入れるために適切なバッファを準備し、ストレージフローコントローラに準備完了したらそれを通知し、転送待機中要求をイニシエータに送る。

2 . フローコントローラは、ホストドライバに対して、要求と D M A のためのバッファポインタとの受領を肯定応答する。

3 . フローコントローラは、次に、実行される R 2 T コマンドをスケジューラに対してスケジュールに入れる。

4 . スケジューラは、コマンドをこのコマンドを実行する準備が完了しているパケットプロセッサ複合システムの 1 つに出す。

5 . パケットプロセッサは、セッションキャッシュコントローラからセッションエント

10

20

30

40

50

リを要求する。

6. セッションエントリが、パケットプロセッサに戻される。

7. パケットプロセッサは、TCPパケットを形成し、R2Tコマンドをカプセル化し、出力待ち行列にそれを送る。

8. パケットは、次に、ネットワークメディアインタフェースに送られ、それは、次に、イニシエータにこのパケットを送る。転送が安全転送である必要がある場合は、セキュリティエンジンを伴うことができるであろう。

9. 次に、図51に示すように、イニシエータは、書込データをターゲットに送ることによってR2Tに応答する。ネットワークメディアインタフェースは、パケットを受け取り、入力待ち行列にそれを入れる。

10. パケットスケジューラは、入力待ち行列からパケットを検索する。

11. パケットは、分類エンジンに対してスケジュールに入れられる。

12. 分類エンジンは、分類タグと共に分類されたパケットをスケジューラに供給する。図示のフローは、暗号化されていないパケットのためであり、従って、セキュリティエンジンは実行されない。

13. スケジューラは、パケットプロセッサ待ち行列に対するフローベースのリソース割当待ち行列に基づいてパケットを割り当てる。パケットは、次に、パケットプロセッサがこのパケットの実行を待機中の時に、パケットプロセッサ復号システムに転送される。

14. パケットプロセッサは、セッションキャッシュエントリを要求する（それがそのローカルキャッシュにまだそれを有していない場合）。

15. セッションエントリが、要求するパケットプロセッサに戻される。

16. パケットプロセッサは、全てのTCP/IP機能を実行し、セッションエントリを更新し、ストレージエンジンは、以前のR2Tに回答して書込コマンドとしてPDUを抽出する。それは、ストレージセッションエントリを更新し、パケットをそれがホストバッファに転送されるようにホスト出力待ち行列にルーティングする。パケットは、ホスト割当の宛先バッファ内へのこのパケットのDMAの実行に使用することができるメモリ記述子又はメモリ記述子リストを使用してタグ付けすることができる。

17. ホストインタフェースブロックは、DMAを実行し、書込データコマンドのこのセグメントを終了する。

#### 【0085】

図52は、ターゲット読取データフローを示す。このフローは、図41に示すイニシエータR2T及び書込データフローに非常に似ている。このフローで従うべき主な段階は、次の通りである。

1. ネットワークメディアインタフェースブロックから入力パケットが受信される。

2. パケットスケジューラは、入力待ち行列からパケットを検索する。

3. 分類のためにパケットがスケジュールに入れられる。

4. 分類されたパケットが、分類タグを付けて分類子から戻る。

a. 分類とフローベースのリソース割当とにより、パケットは、パケット上で作動するパケットプロセッサ複合システムに割り当てられる。

5. パケットプロセッサ複合システムは、セッションキャッシュ内のセッションエントリを検索する（ローカルに存在しない場合）。

6. セッションキャッシュエントリが、パケットプロセッサ複合システムに戻される。

7. パケットプロセッサは、読取コマンドPDUを判断し、フローコントローラに対する要求を用いて読取データを要求する。

8. フローコントローラは、ホストインタフェースに対してDMAを開始する。

9. ホストインタフェースは、DMAを実行し、ホスト入力待ち行列にデータを戻す。

10. パケットプロセッサ複合システムは、ホスト入力待ち行列からデータを受け取る。

。

11. パケットプロセッサ複合システムは、データの周囲に有効なPDUとパケットを形成し、適切なセッションエントリを更新し、出力待ち行列にパケットを転送する。



12. パケットは、宛先にデータパケットを送信する出力ネットワークメディアインタフェースブロックに転送される。

【0086】

上述のフローの説明は、高帯域幅データ転送に伴ういくつかの主なフローの一例である。図示していないが本発明のIPプロセッサによってサポートされるフラグメント化データフロー、複数の異なる種類のエラーを有するエラーフロー、ネーム解決サービスフロー、アドレス解決フロー、及び、ログイン及びログアウトフローなどのようないくつかのフローが存在する。

本発明のIPプロセッサは、限定する意味でなく以下に概説する製造工程を使用して、様々な実現可能な実施形態の選択されたものにおけるハードウェア製品に製造することができる。10  
プロセッサは、特定ターゲット製造工程技術に対する機能性、タイミング、及び他の設計及び製造制約に対するRTLレベル、回路/回路図/ゲートレベル、レイアウトレベルなどのような様々なレベルのチップ設計抽象概念で設計及び検証することができる。適切な物理的/レイアウトレベルでのプロセッサ設計を用いて、ターゲット処理技術においてチップを製造するのに使用されるマスクセットを作り出すことができる。次に、マスクセットを用いて、選択された処理技術に対して使用される段階を通してプロセッサチップが構築される。プロセッサチップは、次に、製造されたプロセッサ製品の品質を保証するのに適切な場合は、試験/包装工程を通過することができる。

上述の事項は本発明の特定の実施形態に関連したものであるが、当業者は、本発明の原理及び精神から逸脱することなく、これらの実施形態に変更を行うことができることを認20  
めるであろう。

【図面の簡単な説明】

【0087】

【図1】階層化されたSCSIアーキテクチャと、イニシエータとターゲットシステムの間25  
に位置するそれぞれの層の間の対話とを示す図である。

【図2】イニシエータとターゲットシステムの間でのiSCSI及びTCP/IPベースの移送を有する階層化されたSCSIアーキテクチャを示す図である。

【図3】ファイバチャンネルのようなハードウェア志向プロトコルを有するソフトウェアベースTCP/IPスタックのOSスタック比較を示す図である。

【図4】他の非IPハードウェア志向プロトコルとの性能パリティを提供するためのハードウェアベースTCP/IP実施によるOSスタックを示す図である。30

【図5】ネットワーキング及びストレージスタックを実行するオペレーティングシステム層を表すホストソフトウェアスタックを示す図である。

【図6】ソフトウェアTCPスタックデータ転送を示す図である。

【図7】本特許に説明するようなホストプロセッサからのTCP/IPオフロードを使用する遠隔の直接メモリアクセスデータ転送を示す図である。

【図8】IPネットワーク上でブロックストレージデータを移送するためのホストソフトウェアSCSIストレージスタック層を示す図である。

【図9】本発明の実施形態の特定のiSCSIストレージネットワーク層スタックの詳細を示す図である。40

【図10】本発明の実施形態のTCP/IPネットワークスタックの機能的詳細を示す図である。

【図11】本発明の実施形態の様々な要素を通じたiSCSIストレージデータフローを示す図である。

【図12】本発明に有用なiSCSIストレージデータ構造を示す図である。

【図13】本発明の実施形態に有用なセッションデータベースエントリのためのTCP/IP伝送制御ブロックデータ構造を示す図である。

【図14】本発明の実施形態に有用なiSCSIセッションデータベース構造を示す図である。

【図15】本発明の実施形態に有用なiSCSIセッションメモリ構造を示す図である。50

【図16】本発明の実施形態に有用なIPネットワークアプリケーションプロセッサの高レベルアーキテクチャブロック図である。

【図17】図16のIPネットワークアプリケーションプロセッサのアーキテクチャブロック図の詳細図である。

【図18】IPプロセッサの一実施形態の入力待ち行列とコントローラを示す図である。

【図19】IPプロセッサの一実施形態に有用なパケットスケジューラ、シーケンサ、及び負荷バランスを示す図である。

【図20】IPストレージプロセッサの一実施形態のポリシーエンジンブロックを含むパケット分類エンジンを示す図である。

【図21】高レベルでのIPプロセッサの一実施形態のSANパケットプロセッサブロックの実施形態を概略で示す図である。

【図22】説明したIPプロセッサのSANパケットプロセッサブロックの実施形態をより詳細に示す図である。

【図23】説明したSANパケットプロセッサの一部として使用することができるプログラマブルTCP/IPプロセッサエンジンの実施形態を示す図である。

【図24】説明したSANパケットプロセッサの一部として使用することができるプログラマブルIPストレージプロセッサエンジンの実施形態を示す図である。

【図25】図17のプログラマブルIPプロセッサの出力待ち行列ブロックの実施形態を示す図である。

【図26】ストレージフローコントローラ及びRDMAコントローラの実施形態を示す図である。

【図27】本発明の実施形態に有用なIPプロセッサのホストインタフェースコントローラブロックの実施形態を示す図である。

【図28】セキュリティエンジンの実施形態を示す図である。

【図29】説明したプロセッサに有用なメモリ及びコントローラの実施形態を示す図である。

【図30】説明した分類エンジンの実施形態に利用可能なデータ構造を示す図である。

【図31】イニシエータとターゲットの間のストレージ読取フローを示す図である。

【図32】説明したプロセッサのパイプラインステージを通じた読取データパケットフローを示す図である。

【図33】イニシエータとターゲットの間のストレージ書込作動フローを示す図である。

【図34】説明したプロセッサのパイプラインステージを通じた書込データパケットフローを示す図である。

【図35】イニシエータとターゲットの間の遠隔DMA(RDMA)機能を使用したイニシエータとターゲットの間のストレージ読取フローを示す図である。

【図36】説明したプロセッサのパイプラインステージを通じてRDMAを使用したイニシエータとターゲットの間の読取データパケットフローを示す図である。

【図37】RDMA機能を使用したイニシエータとターゲットの間のストレージ書込フローを示す図である。

【図38】説明したプロセッサのパイプラインステージを通じてRDMAを使用した書込データパケットフローを示す図である。

【図39】説明したプロセッサのパイプラインステージを通じたイニシエータコマンドフローをより詳細に示す図である。

【図40】説明したプロセッサのパイプラインステージを通じた読取パケットデータフローをより詳細に示す図である。

【図41】説明したプロセッサのパイプラインステージを通じた書込データフローをより詳細に示す図である。

【図42】パケットがチップセットテキストか又はそうでなければ説明されたプロセッサのパイプラインステージを通じた安全なパケットの時の読取データパケットフローを示す図である。

10

20

30

40

50

【図43】パケットがチップセットテキストか又はそうでなければ説明されたプロセッサのパイプラインステージを通じた安全なパケットの時の書込データパケットフローを示す図である。

【図44】説明したプロセッサのパイプラインステージを通じたRDMAバッファ通知フローを示す図である。

【図45】説明したプロセッサのパイプラインステージを通じたRDMA書込フローをより詳細に示す図である。

【図46】説明したプロセッサのパイプラインステージを通じたRDMA読取データフローをより詳細に示す図である。

【図47】説明したプロセッサのパイプラインステージを通じたセッション作成フローの段階を示す図である。

10

【図48】説明したプロセッサのパイプラインステージを通じたセッション取外しフローの段階を示す図である。

【図49】説明したプロセッサのパイプラインステージを通じてターゲットの展望で見たセッション作成及びセッション取外しの段階を示す図である。

【図50】説明したプロセッサのパイプラインステージを通じたターゲットサブシステムのR2Tコマンドフローを示す図である。

【図51】説明したプロセッサのパイプラインステージを通じたターゲットサブシステムの書込データフローを示す図である。

【図52】説明したプロセッサのパイプラインステージを通じたターゲット読取データフローを示す図である。

20

【符号の説明】

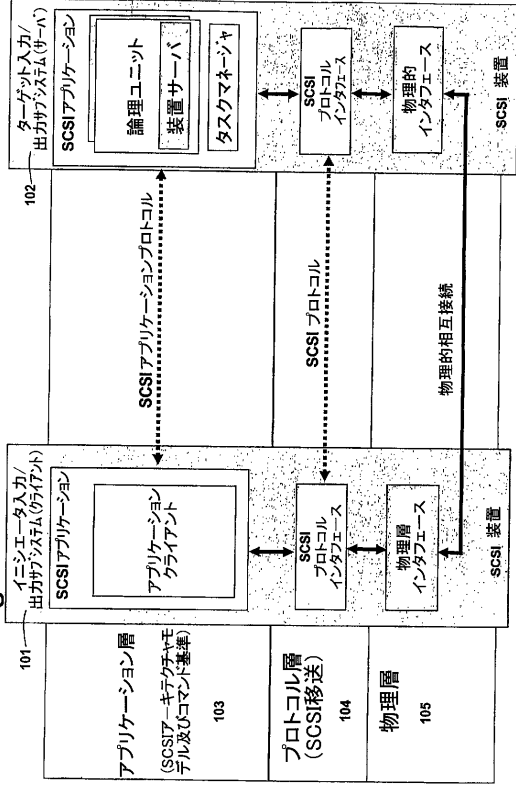
【0088】

- 1701 入力待ち行列
- 1702 パケットスケジューラ
- 1703 分類エンジン
- 1705 セキュリティエンジン
- 1706 ストレージエリアネットワーク(SAN)パケットプロセッサ
- 1707 ホスト入力待ち行列
- 1709 ホスト出力待ち行列
- 1712 出力待ち行列

30

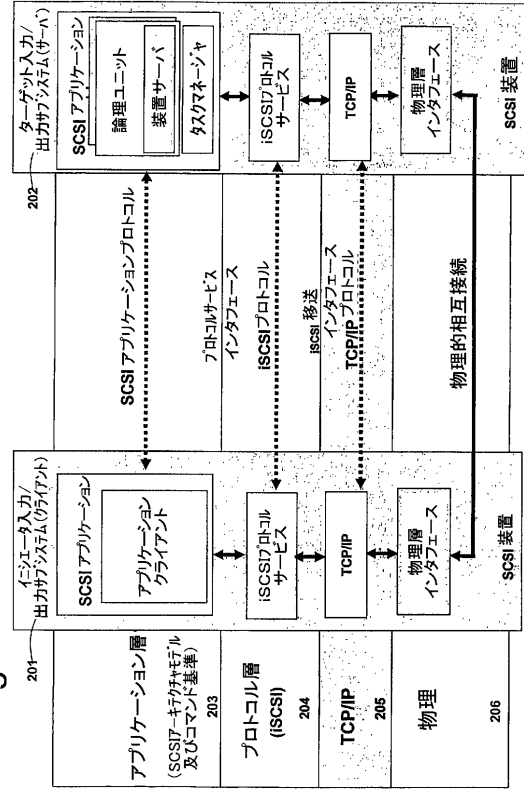
【 図 1 】

Fig. 1 アーキテクチャ層



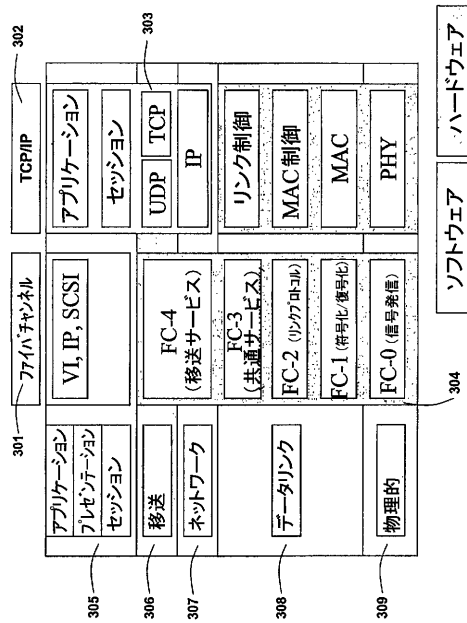
【 図 2 】

Fig. 2 iSCSI上のSCSIアーキテクチャ層



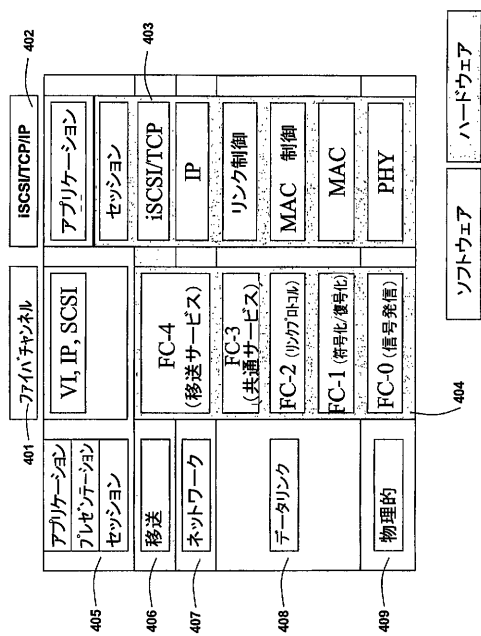
【 図 3 】

Fig. 3 OSI スタック



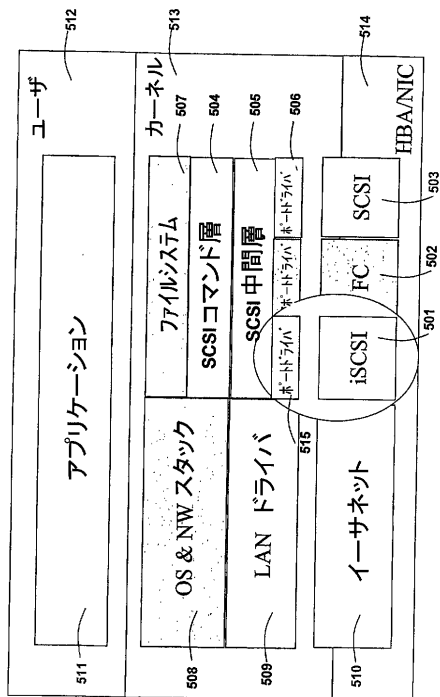
【 図 4 】

Fig. 4 スタック(ハードウェアTCP/IP付き)



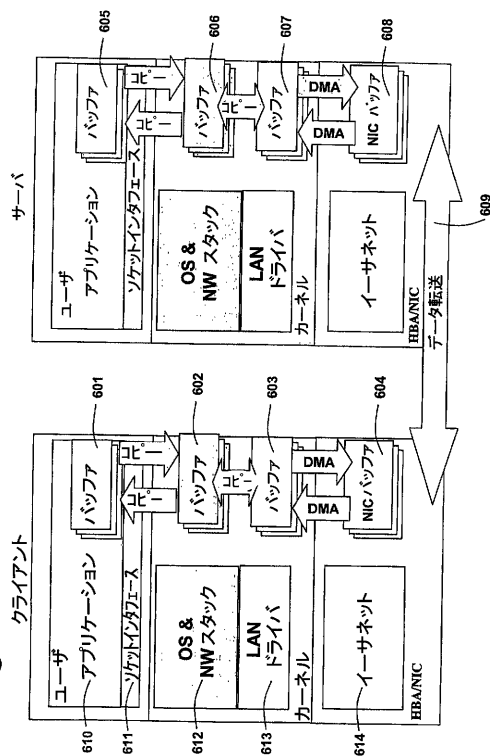
【 図 5 】

Fig. 5 ホストソフトウェアスタック (iSCSI)



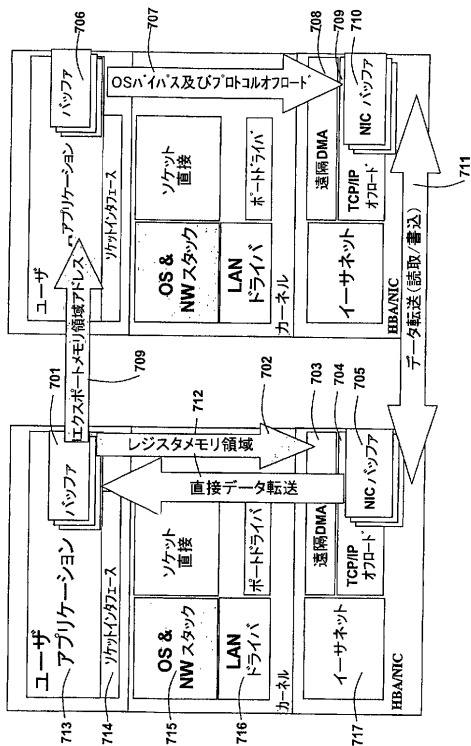
【 図 6 】

Fig. 6 SW TCP スタックデータ転送



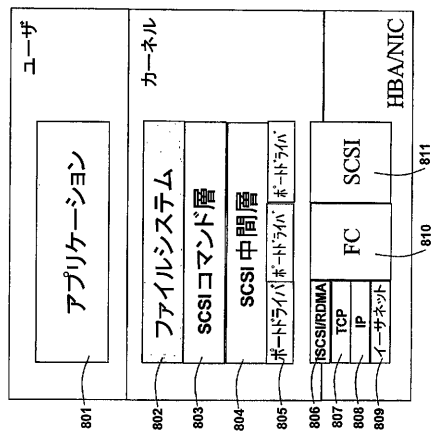
【 図 7 】

Fig. 7 遠隔直接メモリアクセス



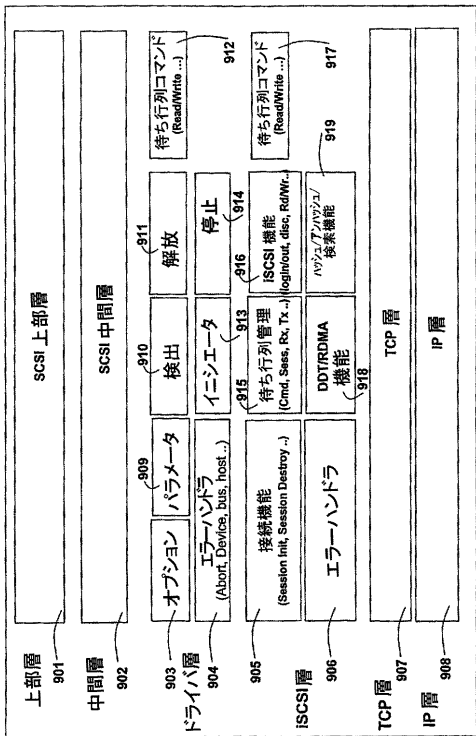
【 図 8 】

Fig. 8 ホストソフトウェアスタック (SCSI)



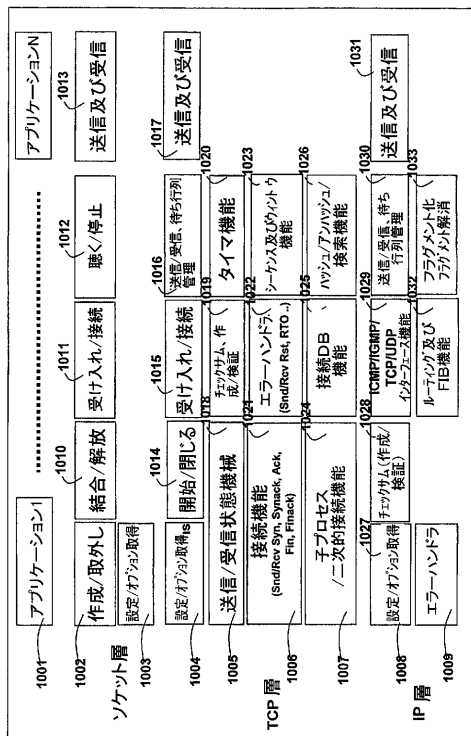
【 図 9 】

Fig. 9 iSCSIスタック



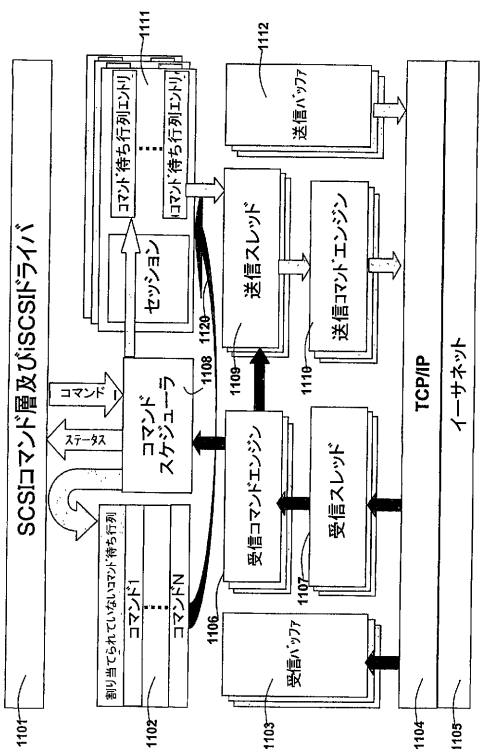
【 図 10 】

Fig. 10 TCP/IPスタック



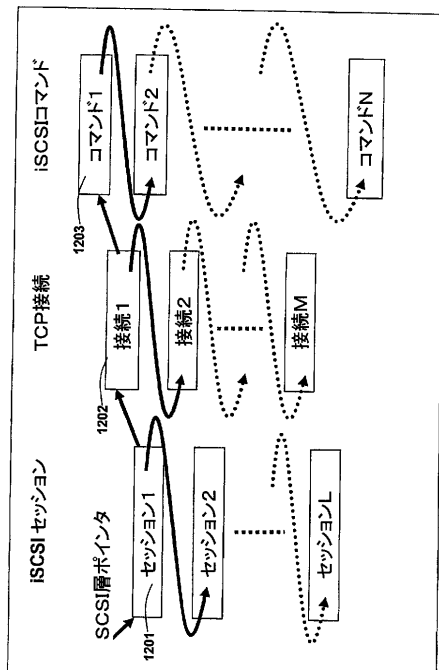
【 図 11 】

Fig. 11 iSCSI データフロー



【 図 12 】

Fig. 12 iSCSIデータ構造



【 図 1 3 】

Fig. 13 TCP/IPセッションDBエントリ

Source_IP	Destination_IP	Source_Port	Destination_Port	プロトコル
Connection_ID	TCP_Window_PTR	Window_Size	Sender_MSS/MTU	Expected_SQN
TCP_State	TCP_SND_SEQ#	TCP_SEQ#	TCP_ACK#	TCP_RCV_SEQ#
Fragment_PTR	Packet_send_PTR (パケット新着)	Packet_Start_Time	Last_PKT_Time	TCP_Parameters (スロースタート、再送信)
Packets_Count	輻輳ウィンドウ	ssthresh	レシーバMSS	他のフィールド
Frame_Key	順序不揃いPTR	他のフレームフィールド		
Hash_Key	NXT_H_Match_PTR	有効/無効		

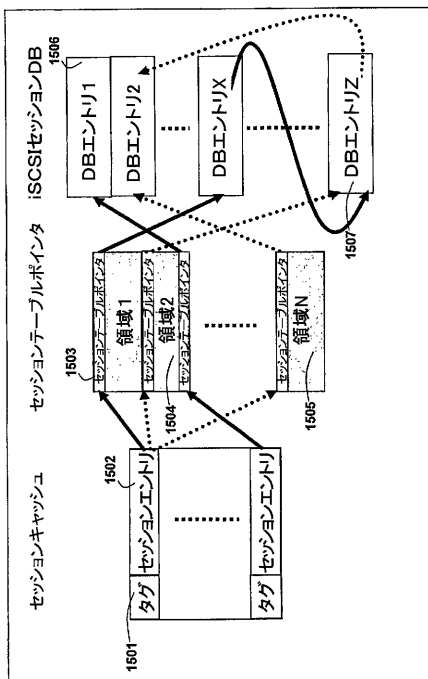
【 図 1 4 】

Fig. 14 iSCSIセッションDBエントリ

iSCSI_State	Data_SEQ#	Command_SEQ#	Task_TAG	R2T_SEQ#
Status_SEQ#	EXP_Status_SEQ#	EXP_CMD_SEQ#	TARGET_MODE (送信請求か否か) (属のターゲット、R2P、...)	Target_Parameters
ISID	TSID	ポータルグループタグ	転送方向	他のフィールド
Connection_ID	コマンドID	RDMA使用可能化	RDMAキー	他のRDMAパラメータ
MDL_List_Pointer	メモリ記述子	ブロックサイズ	ブロック数 (送信又は受信)	UPL_Pointers
Hash_Key	NXT_H_Match_PTR	有効/無効		

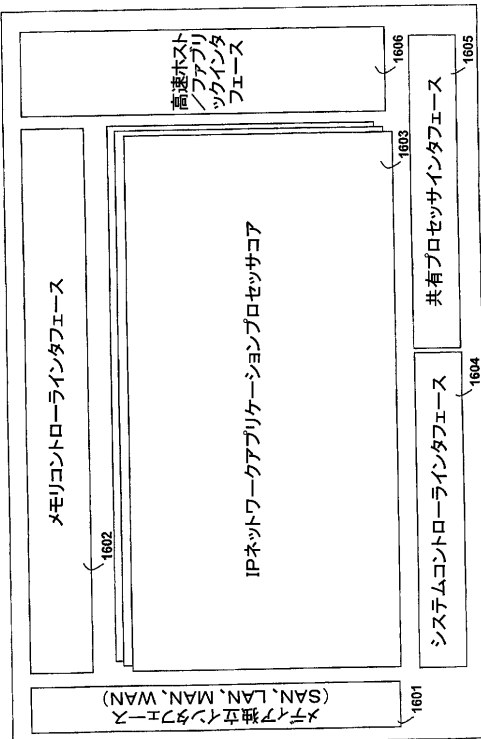
【 図 1 5 】

Fig. 15 iSCSIセッションメモリ



【 図 1 6 】

Fig. 16 IPネットワークアプリケーションプロセッサ



【 図 17 】

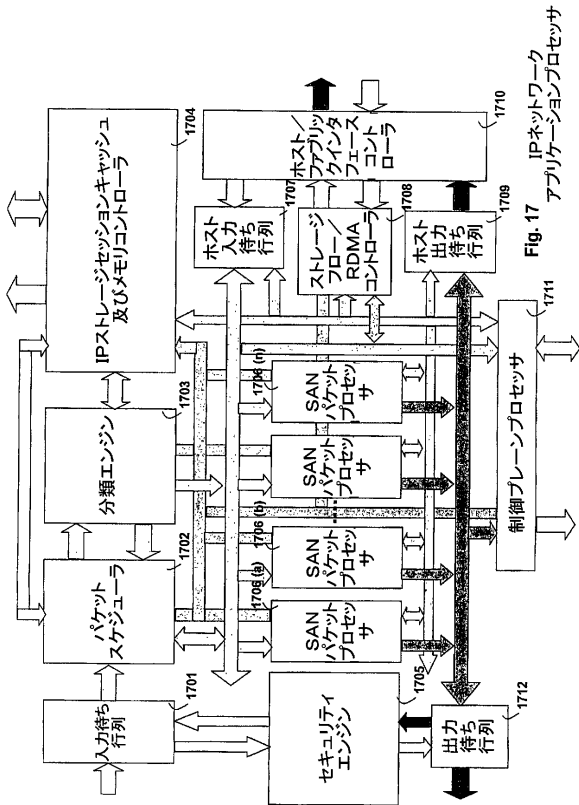


Fig. 17 IPネットワーク アプリケーションプロセッサ

【 図 18 】

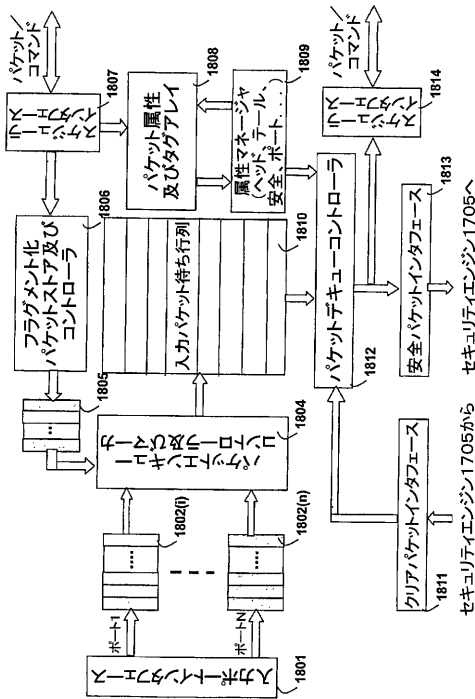


Fig. 18 入力待ち行列及びコントロール

【 図 19 】

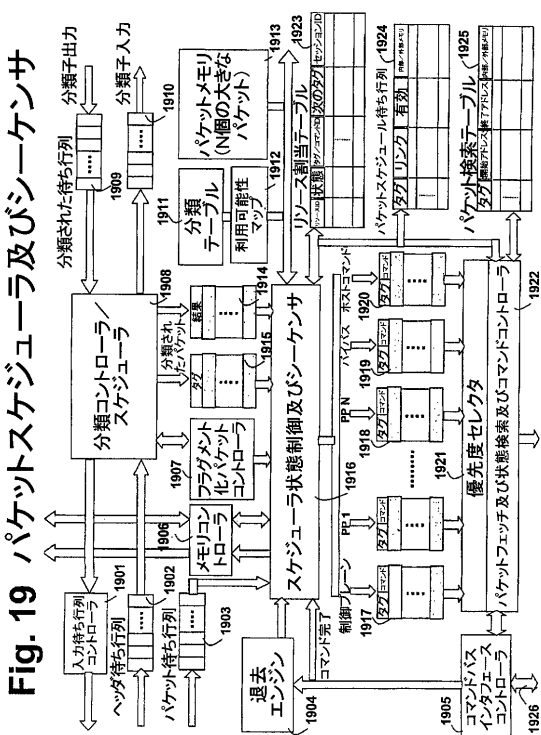


Fig. 19 パケットスケジューラ及びシーケンサ

【 図 20 】

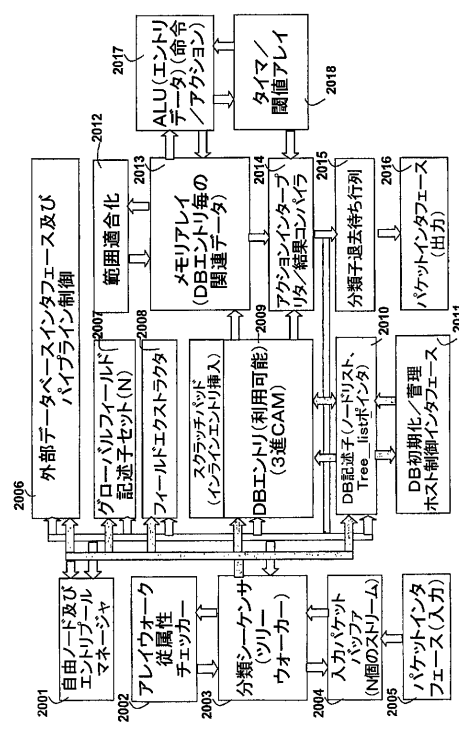
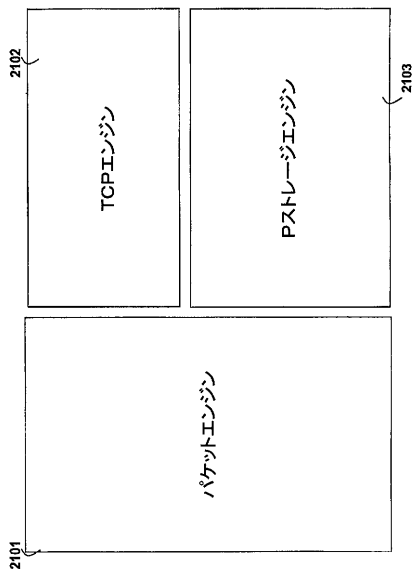


Fig. 20 パケット分類／ポリシーエンジン



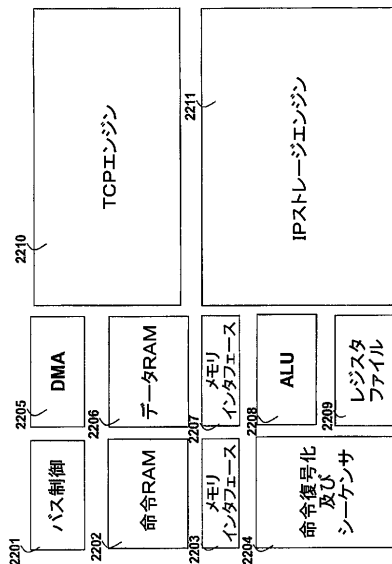
【図 2 1】

Fig. 21 SAN/パケットプロセッサ



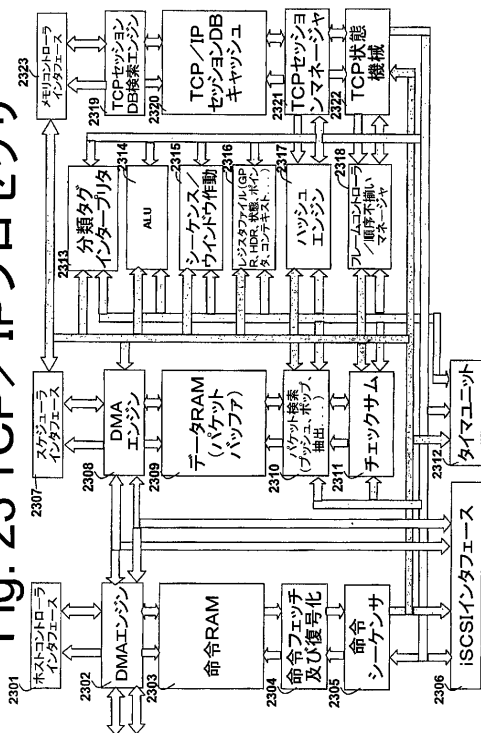
【図 2 2】

Fig. 22 SAN/パケットプロセッサ



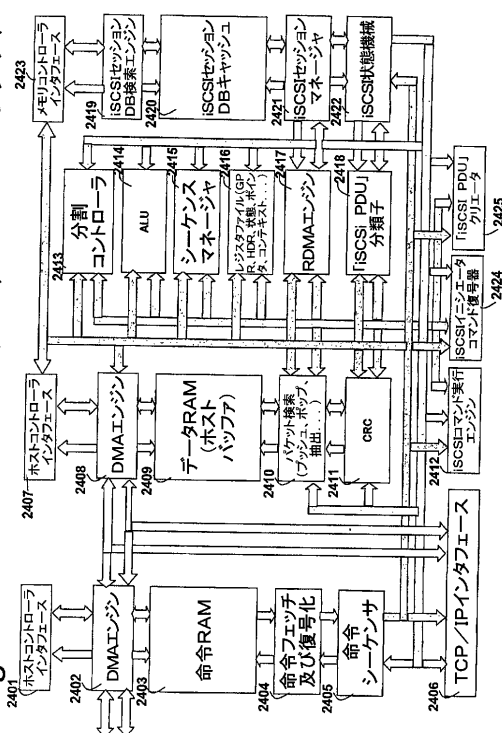
【図 2 3】

Fig. 23 TCP/IPプロセッサ



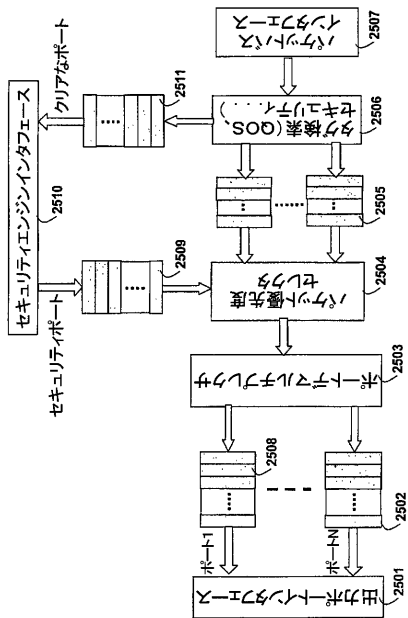
【図 2 4】

Fig. 24 IPストレージエンジン (iSCSIプロセッサ)



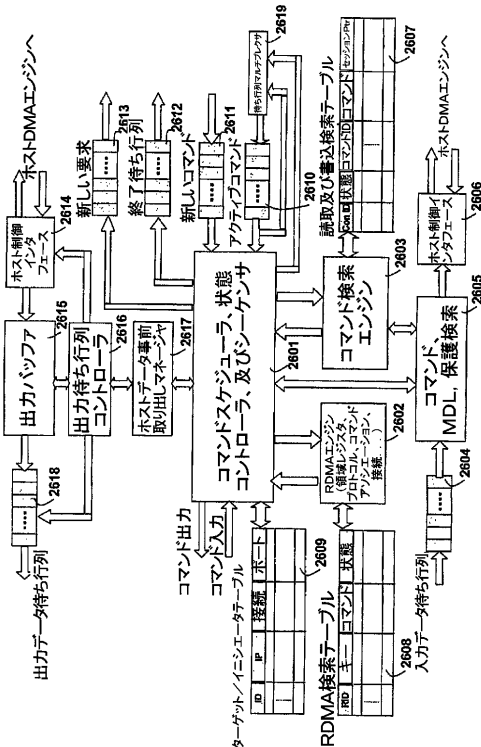
【 図 2 5 】

Fig. 25 出力待ち行列



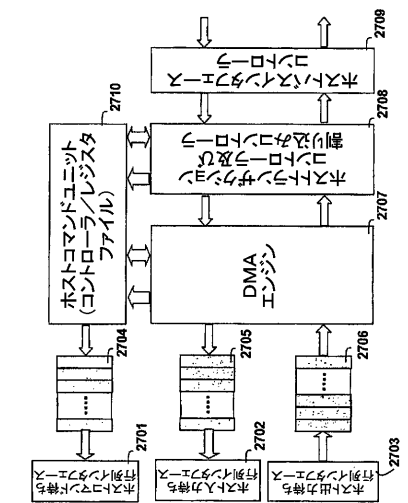
【 図 2 6 】

Fig. 26 ストレージプロセッサ及びRDMAコントローラ



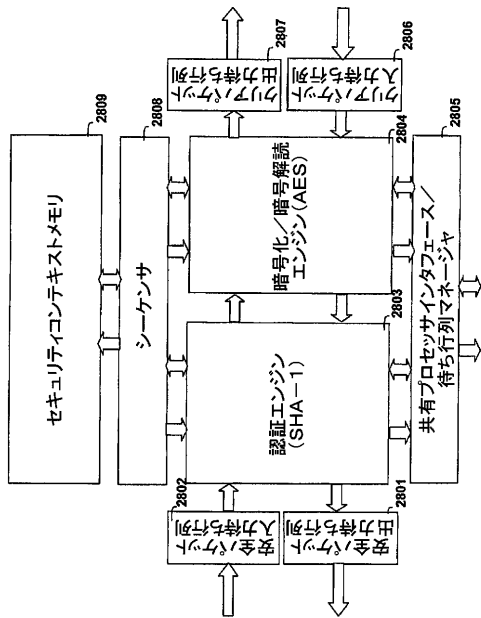
【 図 2 7 】

Fig. 27 ホストインタフェースコントローラ



【 図 2 8 】

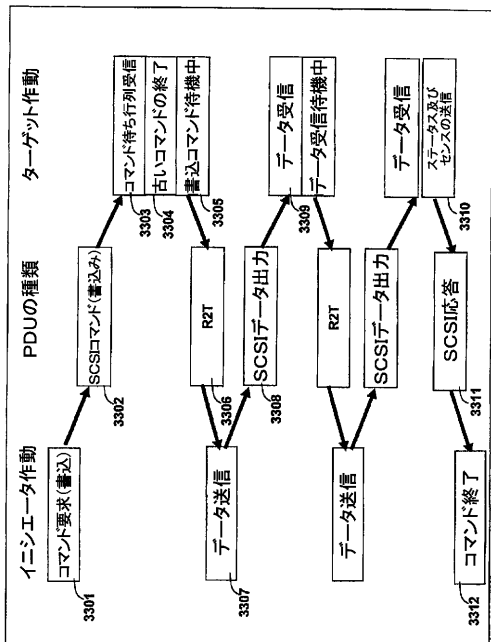
Fig. 28 セキュリティエンジン





【 図 3 3 】

Fig. 33 書込作動



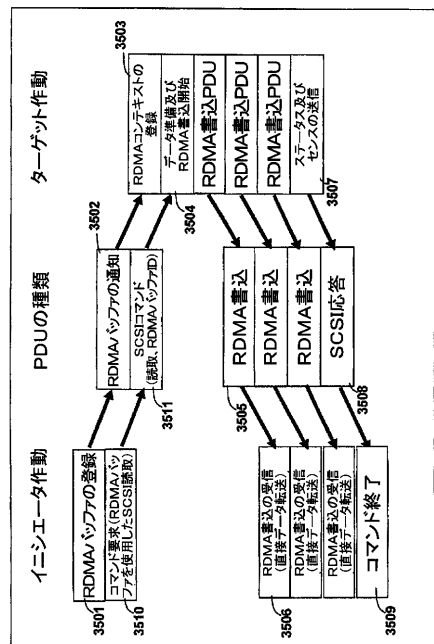
【 図 3 4 】

Fig. 34 書込データパケットフロー

3401	3402	3403	3404	3405	3406	3407	3408
受信	分類	スケジュール	実行	DMA	実行	セキュリティ	送信
検出 パケット抽出 L2確認 待ち行列	L3フィールド L4フィールド TCPマーク付け iSCSI/TCP ID エントリの転送 TCPマーク付け iSCSI/TCP ID エントリの転送 TCP状態更新 iSCSI 状態更新 iSCSI作動識別	スケジュール ハッシュ 従属性 フローベースの リソース割当 待ち行列	実行 パケット転送 iSCSI/TCP ID エントリの転送 TCPマーク付け TCP状態更新	優先度の 待ち行列 システム メモリからの データ転送	ヘッダ準備 チェックサム パケット準備 iSCSI/TCP DBエントリの転送 TCP状態更新 iSCSI状態 更新 パケット転送	セキュリティ 暗号化 メッセージ サイズ 待ち行列	アセンブル L2ヘッダ 待ち行列 送信
3415	3416	3409	3410	3411	3412	3413	3414
R2T読取応答							
書込データ転送							

【 図 3 5 】

Fig. 35 RDMA書込を使用したiSCSI読取



【 図 3 6 】

Fig. 36 読取データパケットフロー(RDMA)

3601	3602	3603	3604	3605	3606
受信	セキュリティ	スケジュール	実行	DMA	
検出 パケット抽出 L2確認 待ち行列	セキュリティ アソシエーション 認証 暗号解読 待ち行列	ハッシュ 従属性 フローベースの リソース割当 待ち行列	実行 パケット転送 iSCSI/TCP ID エントリの転送 TCP状態更新 TCPマーク付け iSCSI/RDMA 状態更新 RDMA 作動識別	優先度ベースの 待ち行列 システムメモリ へのデータ転送 RDMA 作動終了番号 iSCSI 作動終了(最終 RDMAの場合)	3612
3607	3608	3610	3611	3611	3612



【 図 4 1 】

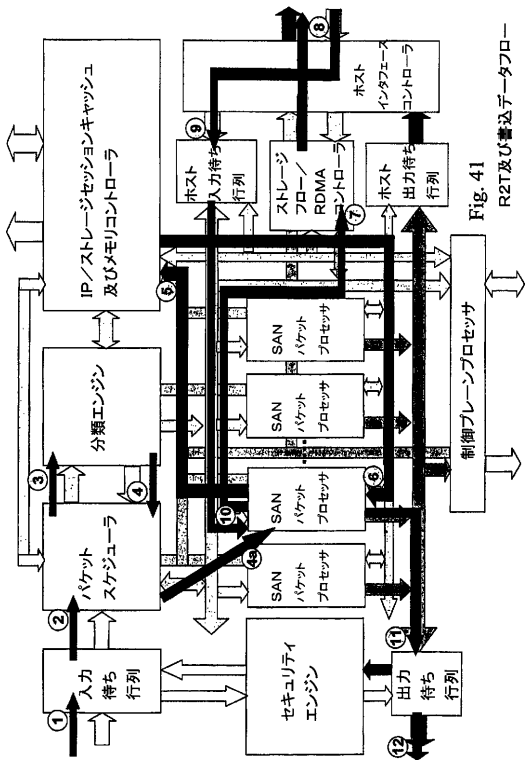


Fig. 41  
R2T及び書込データフロー

【 図 4 2 】

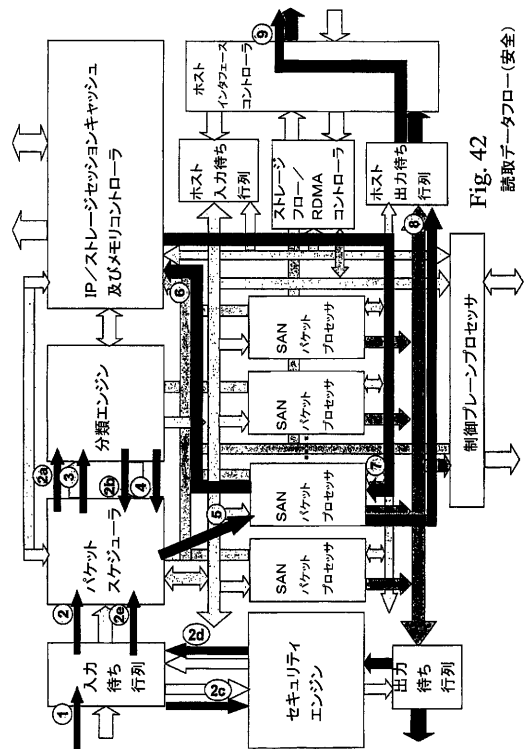


Fig. 42  
読取データフロー(安全)

【 図 4 3 】

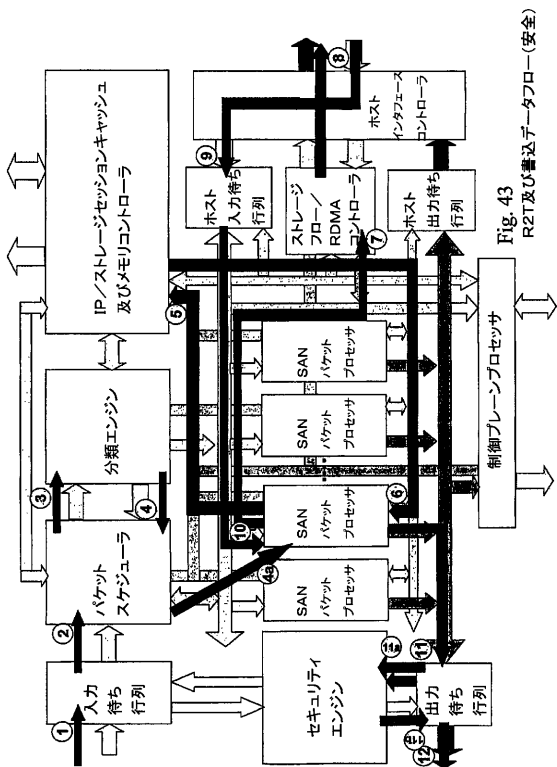


Fig. 43  
R2T及び書込データフロー(安全)

【 図 4 4 】

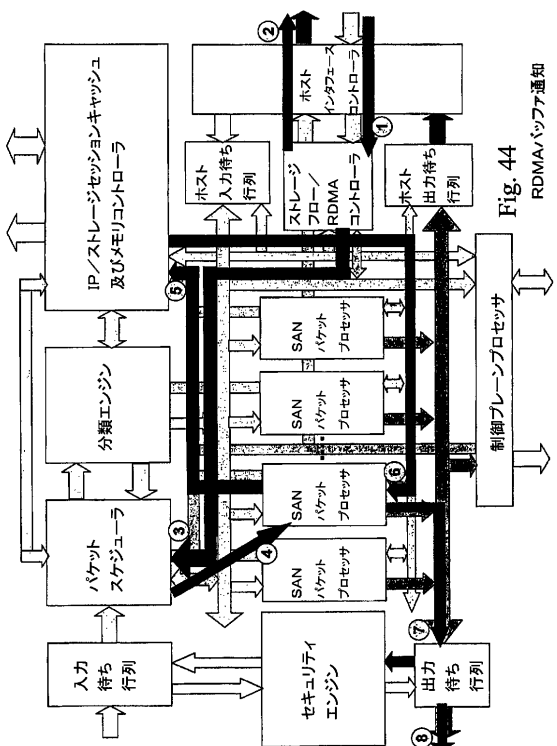


Fig. 44  
RDMAハップア通知

【 図 45 】

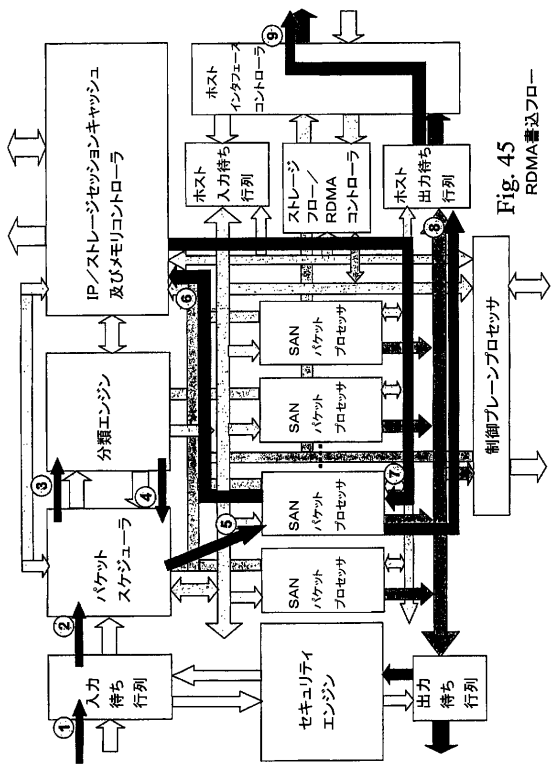


Fig. 45  
RDMA書き込みフロー

【 図 46 】

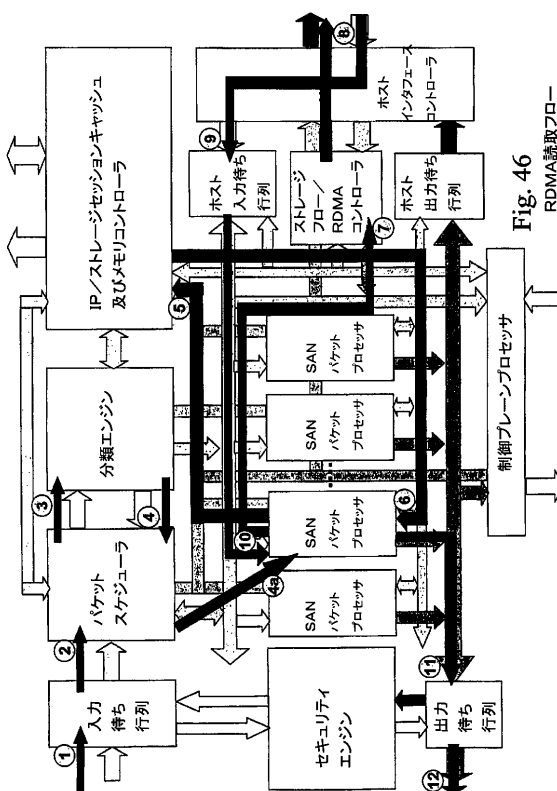


Fig. 46  
RDMA読取りフロー

【 図 47 】

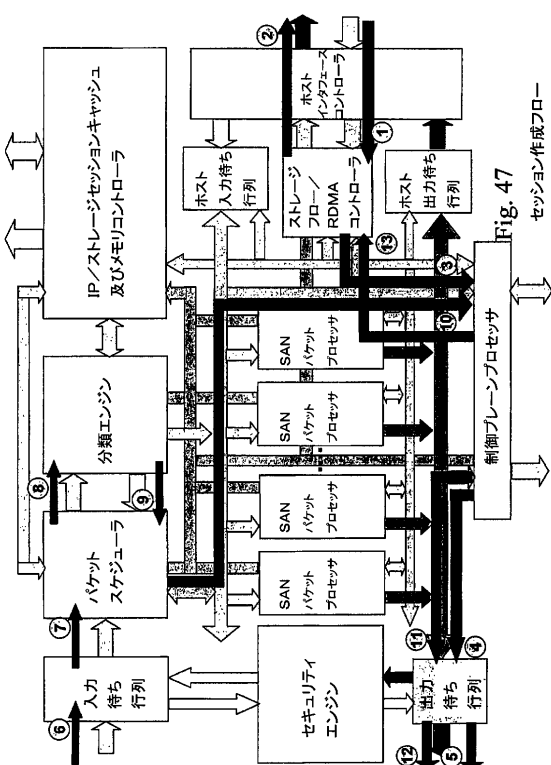


Fig. 47  
セッション作成フロー

【 図 48 】

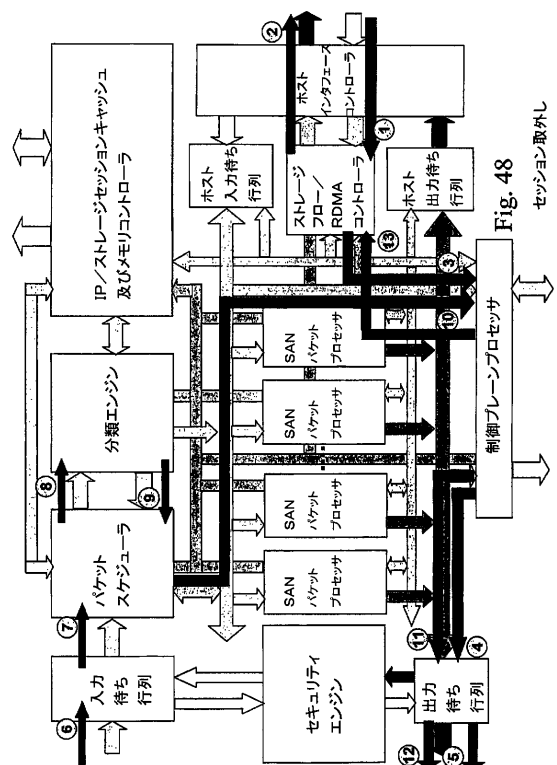
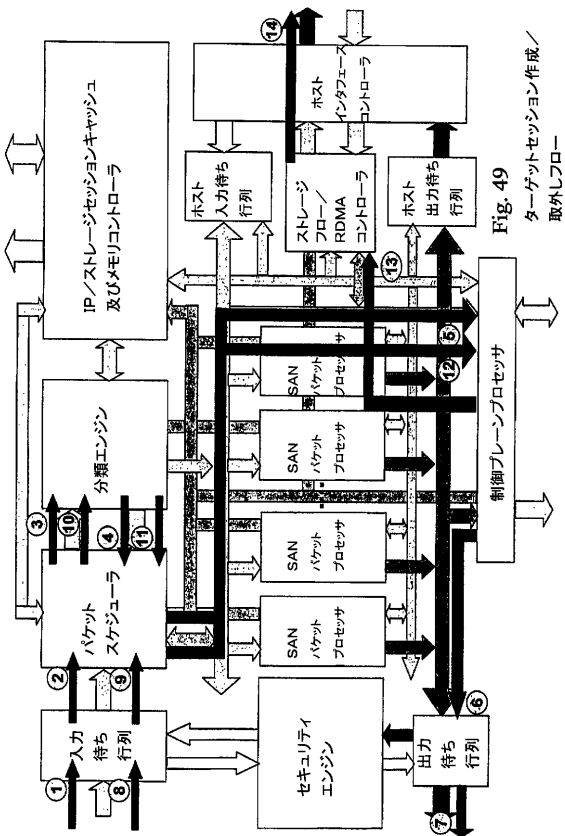
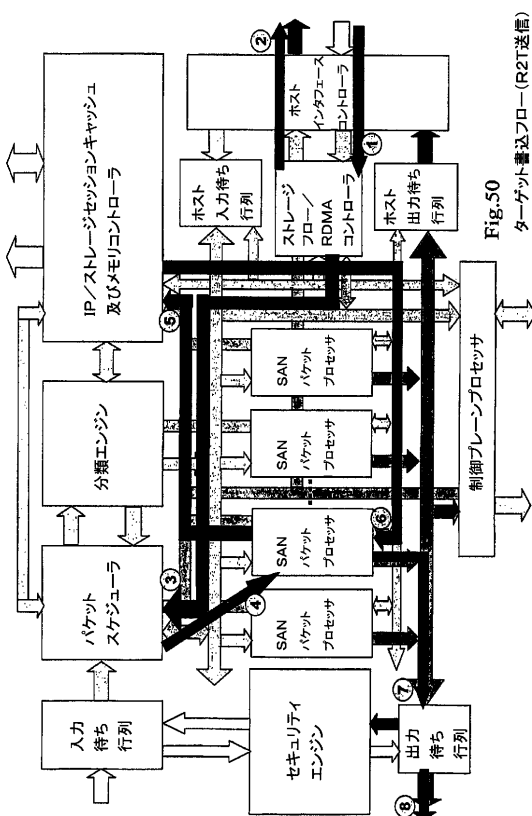


Fig. 48  
セッション取外し

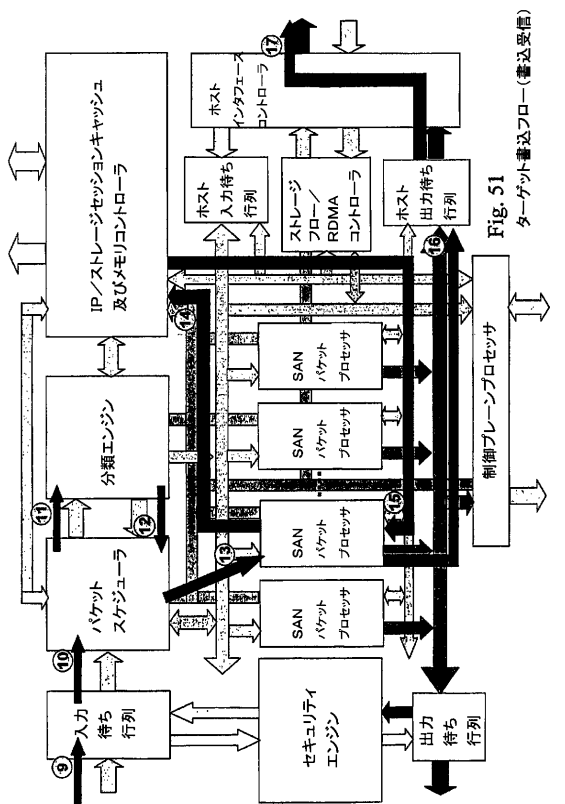
【 図 49 】



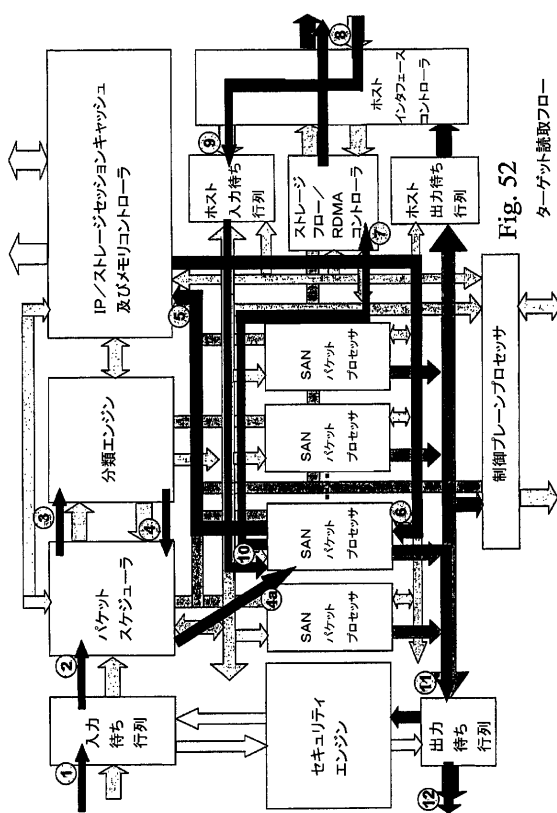
【 図 50 】



【 図 51 】



【 図 52 】





---

フロントページの続き

(72)発明者 パンドヤ アシシュ エイ  
アメリカ合衆国 カリフォルニア州 95762 エル ドラド ヒルズ ラファイエット ドラ  
イヴ 4318

審査官 横山 佳弘

(56)参考文献 特開2002-063060(JP,A)  
特開2001-268159(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 13/10

G06F 13/00

G06F 3/06