

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4066382号
(P4066382)

(45) 発行日 平成20年3月26日(2008.3.26)

(24) 登録日 平成20年1月18日(2008.1.18)

(51) Int.Cl.

F I

H04L 12/56 (2006.01)

H04L 12/56 200Z

請求項の数 3 (全 40 頁)

(21) 出願番号 特願2005-282931 (P2005-282931)
 (22) 出願日 平成17年9月28日(2005.9.28)
 (62) 分割の表示 特願2001-520649 (P2001-520649)
 の分割
 原出願日 平成12年8月24日(2000.8.24)
 (65) 公開番号 特開2006-20371 (P2006-20371A)
 (43) 公開日 平成18年1月19日(2006.1.19)
 審査請求日 平成17年9月28日(2005.9.28)
 (31) 優先権主張番号 09/384,692
 (32) 優先日 平成11年8月27日(1999.8.27)
 (33) 優先権主張国 米国(US)

(73) 特許権者 390009531
 インターナショナル・ビジネス・マシー
 ズ・コーポレーション
 INTERNATIONAL BUSIN
 ESS MASCHINES CORPO
 RATION
 アメリカ合衆国10504 ニューヨーク
 州 アーモンク ニュー オーチャード
 ロード
 (74) 代理人 100086243
 弁理士 坂口 博
 (74) 代理人 100091568
 弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 ネットワーク・スイッチ及びコンポーネント及び操作方法

(57) 【特許請求の範囲】

【請求項 1】

制御点プロセッサと、該制御点プロセッサに動作的に接続されるネットワーク・プロセッサとを含む通信装置において、前記ネットワーク・プロセッサの制御の下で、レジスタ及びメモリをアクセスすることを可能にする情報を有するガイド・フレームを経路指定する方法であって、

前記制御点プロセッサ内に配置される制御点機能を使用し、ガイド・フレームを生成するステップと、

前記制御点プロセッサ内のデバイス・ドライバを使用し、前記ガイド・フレームを、前記ネットワーク・プロセッサに関連付けられる複数のメディア・インタフェースの1つに送信するステップと、

前記メディア・インタフェース内の媒体アクセス制御ハードウェアを使用し、前記送信されたガイド・フレームを回復するステップと、

回復された前記ガイド・フレームをメモリに記憶するステップと、

記憶された前記ガイド・フレームを、前記ガイド・フレーム内で識別されるエンティティに経路指定するステップと

を含む方法。

【請求項 2】

前記ガイド・フレームにより伝搬される命令に従い、前記ガイド・フレームを前記エンティティにより処理するステップと、

10

20

前記ガイド・フレームにより伝搬される情報により要求される場合、処理された前記ガイド・フレームを前記制御点機能に戻すステップと

を含む、請求項 1 記載の方法。

【請求項 3】

前記ガイド・フレームをネットワーク・ルーティング情報によりカプセル化するステップを含む、請求項 1 記載の方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、様々なタイプ及び能力の情報処理システムまたはコンピュータをリンクするために使用される通信ネットワーク装置、及びこうした装置のコンポーネントに関する。特に、本発明は、こうした装置をアセンブルするのに有用なスケラブル・スイッチ装置及びコンポーネントに関する。本発明はまた、改善された多機能インタフェース装置、並びに、こうした装置とメディア速度ネットワーク・スイッチを提供する他の要素との組み合わせに関する。本発明は更に、ネットワーク・スイッチのデータ・フロー処理能力を改善する、こうした装置の操作方法に関する。

10

【背景技術】

【0002】

以下の説明は、ネットワーク・データ通信、並びにこうした通信ネットワークで使われるスイッチ及びルータに関する知識を前提とする。特に、ネットワーク動作をレイヤに分割するネットワーク・アーキテクチャの ISO モデルに精通していることを前提とする。ISO モデルにもとづく典型的なアーキテクチャは、物理パスまたはメディアであるレイヤ 1（ときに "L1" と示される）から始まり、信号は、上流のレイヤ 2、3、4 などを通じてレイヤ 7 に至る。レイヤ 7 は、ネットワークにリンクされるコンピュータ・システム上で実行されるアプリケーション・プログラミング・レイヤである。本明細書では、L1、L2 などは、ネットワーク・アーキテクチャの対応レイヤを指し示すものとする。本開示はまた、こうしたネットワーク通信においてパケット及びフレームとして知られる、ビット・ストリングの基本的理解を前提とする。

20

【0003】

今日のネットワーク化された世界では、帯域幅が重要な資源である。インターネットや他の新たなアプリケーションの登場による、ネットワーク・トラフィックの増加が、ネットワーク・インフラストラクチャの能力を圧迫している。こうした流れに追随するために、様々な組織や機関が、トラフィックの増加や音声とデータのコンバージェンスをサポートするための、より優れた技術や方法を検討している。

30

【0004】

ネットワーク・トラフィックの今日の劇的な増加は、インターネットの普及、情報へのリモート・アクセスの必要性の増加、及び新たなアプリケーションなどによる。インターネットは電子商取引の爆発的な成長により、ときにネットワーク・バックボーンにサポート不能な負荷を課する。データ・トラフィック量の増加の最も重要な 1 原因は、それが初めて音声トラフィックを超えることにある。電子メール、データベース・アクセス、及び

40

【0005】

音声とデータのコンバージェンスは、将来のネットワーク環境を定義する上で重要な役割を果たす。現在、インターネット・プロトコル (IP) ネットワークを介するデータの伝送は無料である。音声伝送は自然に最廉価の経路を通るので、音声は必然的にデータとコンバージされる。ボイス・オーバー IP (VoIP)、ボイス・オーバー ATM (VoATM)、ボイス・オーバー・フレーム・リレー (VoFR) などの技術は、この変化の速い市場においてコスト効率の高い技術と言える。しかしながら、これらの技術への移行を可能にするために、業界は音声のサービス品質 (QoS) を保証し、データ回線上での音声転

50

送のためにどのように課金すべきかを決定しなければならない。1996年のTelecommunication Deregulation Actが、更にこの環境を複雑化する。この法規は、選択音声プロトコルすなわちATMと、選択データ・プロトコルすなわちIPとの間の共生関係を強化する。

【0006】

新たな製品及び機能が使用可能になると、レガシー・システム（システム資産）を統合することが、組織にとって重大な問題となる。既存の設備及びソフトウェアへの投資を保護するために、組織は彼らの現在の運用を混乱させることなく、新たな技術への移行を可能にする解決策を要求する。

【0007】

サービス・プロバイダにとって、ネットワーク障害を取り除くことが先決である。ルータはしばしば、これらの障害の原因となる。しかしながら、ネットワーク輻輳は一般に、しばしば帯域幅問題と誤診され、より高い帯域幅を求めることにより取り組まれる。今日、メーカはこの難題を認識しつつある。彼らはネットワーク・プロセッサ技術に移行することにより、帯域幅資源をより効率的に管理し、ルータ及びネットワーク・アプリケーション・サーバで一般に見いだされる高度データ・サービスを有線速度で提供しつつある。これらのサービスにはロード・バランシング、QoS、ゲートウェイ、ファイアウォール、セキュリティ、及びウェブ（Web）キャッシングが含まれる。

【0008】

リモート・アクセス・アプリケーションでは、性能、要求時帯域幅、セキュリティ及び認証が最優先の部類に入る。QoS及びCoSの統合、統合音声処理、及びより複雑なセキュリティ・ソリューションが、将来のリモート・アクセス・ネットワーク・スイッチの設計を方向付けるであろう。更に、リモート・アクセスは、ISDN、T1、E1、OC-3乃至OC-48、ケーブル、及びxDSLモデムなど、より多くの物理メディアに対応しなければならない。

【0009】

業界コンサルタントはネットワーク・プロセッサ（ここでは"N P"とも称す）を、次に挙げる機能の1つ以上を実行できるプログラマブル通信集積回路として定義した。すなわち、

1) パケット分類：アドレスやプロトコルなどの既知の特性にもとづく、パケットの識別。

2) パケット変更：IP、ATM、または他のプロトコルに従う、パケットの変更（例えばIPのヘッダ内の存続時間（time-to-live）フィールドの更新）。

3) キュー/ポリシ管理：特定のアプリケーションにおけるパケットのパケット・キューイング、デキューイング、及びスケジューリングの設計方針の反映。

4) パケット転送：スイッチ・ファブリックを介するデータの伝送及び受信、及び適切なアドレスへのパケットの転送または経路指定。

【0010】

この定義は早期NPの基本フィーチャの正確な記述であるが、NPの完全な潜在的能力及び利点はまだ実現されていない。ネットワーク・プロセッサは、従来ソフトウェアで処理されたネットワーク・タスクをハードウェアで実行することにより、帯域幅を増加し、広範囲のアプリケーションにおける待ち時間問題を解決することができる。更に、NPは並列分散処理及びパイプライン処理設計などのアーキテクチャを通じて、速度の改善を提供できる。これらの機能は効率的な検索エンジンを可能にし、スループットを向上させ、複雑なタスクの迅速な実行を提供する。

【0011】

ネットワーク・プロセッサは、PCにおけるCPUのように、ネットワークの基本ネットワーク構築ブロックになるものと期待される。NPにより提供される一般的な機能は、リアルタイム処理、セキュリティ、蓄積交換、スイッチ・ファブリック、及びIPパケット処理及び学習機能である。NPはISOレイヤ2乃至5をターゲットとし、ネットワー

10

20

30

40

50

ク特定タスクを最適化するように設計される。

【 0 0 1 2 】

プロセッサ・モデル N P は、複数の汎用プロセッサ及び特殊論理を組み込む。提供者はこの設計に傾倒することにより、変更に適宜に且つ低廉に対応できる、スケーラブルで柔軟性のあるソリューションを提供しようとしている。プロセッサ・モデル N P は、低レベルの統合において分散処理を可能にし、より高度なスループット、柔軟性及び制御を提供する。プログラマビリティが、新たな A S I C 設計を要求することなく、新たなプロトコル及び技術への容易な移行を可能にする。プロセッサ・モデル N P により、N E V は払戻し不能なエンジニアリング・コストの低減、及び製品化までの時間の短縮の利益を甘受することができる。

10

【 発明の開示 】

【 発明が解決しようとする課題 】

【 0 0 1 3 】

本発明の第 1 の目的は、転送されるデータの処理速度を向上する一方で、サポート機能を各種の潜在的な要求にサイズ変更可能な、データ通信ネットワークで使用されるスケーラブル・スイッチ・アーキテクチャを提供することである。

【 0 0 1 4 】

本発明の別の目的は、単一基板上に集積化され、レイヤ 2、レイヤ 3、レイヤ 4 及びレイヤ 5 を含むフレームのメディア速度切替えを提供するように協働する、複数のサブアセンブリを含むインタフェース装置またはネットワーク・プロセッサ（これらの用語は互換性をもって使用される）を提供することである。

20

【 課題を解決するための手段 】

【 0 0 1 5 】

前記第 1 の目的は、関連処理ユニットの作業負荷から、従来に比較してより多くのデータ処理を削除するコンポーネント、及びコンポーネントの集合を提供することにより実現される。

【 0 0 1 6 】

インタフェース装置は、ワーク・グループ・スイッチとして、第 1 レベルの機能を提供する独立型のソリューションとして使用されるか、またはより高度な機能のワーク・グループ・スイッチを提供する相互接続ソリューションとして使用されるか、或いは、スイッチング・ファブリック装置との協働により、更に機能向上に向けてスケーリングされる。

30

【 発明を実施するための最良の形態 】

【 0 0 1 7 】

[数 1]

$$Y/\overline{X}$$

は、本明細書では Y / X バーと記載する。

【 0 0 1 8 】

本発明は、本発明の好適な実施例を示す添付の図面を参照して、以下で詳述されるが、説明の冒頭に当たり、当業者であれば、本発明の好適な結果を達成するように、ここで述べる発明を変更することができよう。従って、以下で述べる説明は本発明を制限するものではなく、当業者に対する広範囲の教示の開示として理解されるべきである。

40

【 0 0 1 9 】

ここで開示される装置はスケーラブルで、デスクトップまたはワークグループ・スイッチを相互接続し、こうしたスイッチをネットワーク・バックボーンに統合し、バックボーン・スイッチング・サービスを提供するように機能することができる。この装置はレイヤ 2、レイヤ 3、及びレイヤ 4 + 転送をハードウェアでサポートできる。特定の形態の装置は、デスクトップまたはワークグループ・スイッチ統合のために設計され、また他のもの

50

はコア・バックボーン・スイッチとして設計される。

【 0 0 2 0 】

装置のために使用されるアーキテクチャは、インタフェース装置またはネットワーク・プロセッサ・ハードウェア・システムと、制御点上で実行されるソフトウェア・ライブラリとにもとづき、これについては本明細書の中で詳述される。インタフェース装置またはネットワーク・プロセッサ・サブシステムは、L 2、L 3 及び L 4 + プロトコル・ヘッダの構文解析及び変換のために設計される高性能フレーム転送エンジンである。これはプロトコルがハードウェアの使用により、より高速にスイッチされることを可能にする。インタフェース装置またはネットワーク・プロセッサ・サブシステムは、ボックスを通じて高速パスを提供する一方、ソフトウェア・ライブラリ及び制御点プロセッサは、高速パスを維持するために必要な管理及びルート発見機能を提供する。制御点プロセッサ及びそれ上で動作するソフトウェア・ライブラリは、一緒にシステムの制御点 (C P) を定義する。C P では、透過型ブリッジング及び O S P F などの、実際のブリッジング及びルーティング・プロトコルが実行される。これはシステムの低速パスとも呼ばれる。

10

【 0 0 2 1 】

ここで開示される装置は、マルチレイヤ転送をハードウェアでサポートするが、装置は L 2 専用スイッチとしても動作でき、開示される最も単純な形態では、これがそのデフォルト動作モードである。各ポートは単一のドメインに配置され、任意の装置が任意の他の装置と通信することを可能にする。装置は L 2 において構成可能であり、システム管理者に次のようなフィーチャ、すなわち、ポートを別々のドメインまたはトランクにグループ化したり、仮想 L A N (V L A N) セグメントや、ブロードキャスト及びマルチキャスト・トラフィックを制御するフィルタを構成するなどの能力を提供する。

20

【 0 0 2 2 】

このスケーラブル装置は多くの利点を有する。第 1 に、これはシステム管理者に、L 2 で使用されるハードウェアと同一のものを使用して、同一の速度で、I P 及び I P X トラフィックの L 3 転送及びルーティングを構成する能力を提供する。第 2 に、構内建物を相互接続するために、外部ルータを使用する必要性を排除する一方で、同時に性能を向上させる。第 3 に、L 2 / L 3 サービスの管理を単純化または結合することにより、これらを単一の制御点に組み込む。最後に、スケーラブル装置は付加価値機能を L 4 + 機能と共に提供し、これはサーバ間での負荷平準化を目的として、主幹業務のアプリケーション及びネットワーク・ディスパッチャをサポートするために、システム管理者が異なるトラフィック分類を割当てて能力を提供する。

30

【 0 0 2 3 】

スケーラブル装置は、インタフェース装置またはネットワーク・プロセッサを使用するモジュラ・ユニットとして、または制御点 (C P) として、或いはその基本ビルディング・ブロックとしての、オプションのスイッチング・ファブリック装置として設計される。インタフェース装置は好適には、L 2 / L 3 / L 4 + 高速パス転送サービスを提供し、C P は高速パスを維持するために必要な管理及びルート発見機能を提供する。3 つ以上のインタフェース装置サブシステムが一緒に結合される場合には、オプションのスイッチング・ファブリック装置が使用される。オプションのスイッチング・ファブリック装置は、1 9 9 1 年 4 月 1 6 日発行の米国特許第 5 0 0 8 8 7 8 号 "High Speed Modular Switching Apparatus for Circuit and Packet Switched Traffic" で開示される。

40

【 0 0 2 4 】

装置は、プリント回路基板要素を用いて組み立てられるものと見込まれる。これはここでは "ブレード" とも称される。プリント回路基板要素は回路要素を実装され、装置ハウジング内に設けられるコネクタに収容される。同様の装置は "オプション・カード" としても知られる。装置は、適切なコネクタ及びバックプレーン電気接続が設けられる条件の下で、ブレードが様々なシャーシまたはハウジングの間で交換されることを考慮する。全てのブレード上で見いだされる基本コンポーネントは、キャリア・サブシステムである。キャリア・サブシステムから始まり、3 つのタイプのブレードが生成される。第 1 のタイプは

50

C P専用ブレードであり、これはキャリア・サブシステム及びC Pサブシステムから成る。C P専用ブレードは、冗長性が一番の関心事である製品に対して主に使用される。第2のタイプはC P+メディア・ブレードであり、これはキャリア・サブシステム、C Pサブシステム、及び1対3メディア・サブシステムから成る。C P+メディア・ブレードは、冗長性よりもポート密度が重要視される製品に対して、主に使用される。第3のタイプはメディア・ブレードであり、これはキャリア・サブシステムと1対4メディア・サブシステムとから成る。メディア・ブレードは任意のシャーシ内で使用され、使用されるメディア・サブシステムのタイプは構成可能である。

【0025】

ブレード管理は障害検出、電力管理、新たな装置の検出、初期化及び構成を含む。この管理は様々なレジスタ、入出力信号、及びC Pとキャリア・サブシステムとの間で通信するために使用されるガイド・セル・インタフェースを用いて行われる。しかしながら、シャーシと異なり、全てのブレード上には、プログラマブル装置及びメモリが存在する。プログラマビリティは、ブレードのタイプに依存する。C Pサブシステムがブレード上に存在する場合、C P及びキャリア・サブシステムの両方がプログラマブルである。メディア・サブシステムもプログラマブルであるが、キャリア・サブシステムを通じて間接的にプログラマブルなだけである。

【0026】

高機能製品では更に、スイッチング・ファブリック装置サブシステムを含むスイッチ・ブレードが存在する。このブレードの管理は、故障検出、電力管理、新たな装置の検出、及び初期化を含む。この管理は、C Pサブシステム内にマップされる様々なレジスタ及び入出力信号を用いて行われる。

【0027】

最も単純な形態では、本発明により考慮されるスイッチ装置は制御点プロセッサと、制御点プロセッサに動作可能なように接続されたインタフェース装置とを有する。好適には、ここで開示されるように、インタフェース装置（ネットワーク・プロセッサとして知られる）は、半導体基板と、基板上に形成される複数のインタフェース・プロセッサと、インタフェース・プロセッサによりアクセスされる命令を記憶する、前記基板上に形成される内部命令メモリと、インタフェース装置を通過し、インタフェース・プロセッサによりアクセスされるデータを記憶する、前記基板上に形成される内部データ・メモリと、複数の入出力ポートとを有する、単一の超大規模集積回路（VLSI）装置またはチップである。インタフェース・プロセッサはここではときに、ピコプロセッサまたは処理ユニットと称される。提供されるポートは、内部データ・メモリを外部データ・メモリに接続する少なくとも1つのポートと、インタフェース・プロセッサの指示の下で、インタフェース装置を通過するデータを外部ネットワークと交換する少なくとも2つの他のポートとを含む。制御点プロセッサはインタフェース装置と協働して、インタフェース・プロセッサにより実行される命令を命令メモリにロードする。そして、これらの命令が実行されて、入出力ポート間でのデータの交換、及びデータ・メモリを介するデータの流れを指示する。

【0028】

ここで開示されるネットワーク・プロセッサは、それが組み込まれるスイッチ・アセンブリとは切り離して、本発明に関わる。更に、ここで開示されるネットワーク・プロセッサは、ここで述べられるその要素内に、ここでは詳細に述べられない他の発明を有するものとみなされる。

【0029】

図1は、基板10と、基板上に集積化される複数のサブアセンブリとを含むインタフェース素子チップのブロック図である。サブアセンブリは、アップサイドすなわち上流側構成と、ダウンサイドすなわち下流側構成とに編成される。本明細書では、"アップサイド"は、ネットワークから本明細書で開示される装置に送られるデータ・フローを指し示し、"ダウンサイド"は、逆に装置から、この装置によりサービスされるネットワーク・サービスに送られるデータを指し示す。データ・フローは反復構成に従う。結果的に、アップサ

イド・データ・フロー及びダウンサイド・データ・フローが存在する。アップサイドのサブアセンブリには、エンキュー・デキュー・スケジューリング・アップ (E D S - U P) 論理 16、多重化 M A C アップ (P P M - U P) 14、スイッチ・データ・ムーバ・アップ (S D M - U P) 18、システム・インタフェース (S I F) 20、データ・アライン・シリアル・リンク A (D A S L A) 22、及びデータ・アライン・シリアル・リンク B (D A S L B) 24 が含まれる。データ・アライン・シリアル・リンクは、1999年6月11日出願の米国特許出願第09/330968号 "High Speed Parallel/Serial Link for Data Communication" で詳細に述べられている。ここで開示される本発明の装置の好適な形態は D A S L リンクを使用するが、本発明は、特にデータ・フローが V L S I 構造内にあるように制限される場合などのように、かなり高いデータ・フロー・レートを実現するために、他の形態のリンクが使用されることも考慮する。

10

【0030】

ダウンサイドのサブアセンブリには、D A S L A 26、D A S L B 28、S I F 30、S D M - D N 32、E D S - D N 34、及び P P M - D N 36 が含まれる。チップは更に複数の内部 S R A M、トラフィック管理スケジューラ 40、及び組み込みプロセッサ・コンプレックス (E P C) 12 を含む。インタフェース装置 38 は、イーサネット (R) フィジカル (E N E T P H Y) または A T M フレーマなどの、任意の好適な L 1 回路である。インタフェースのタイプは、部分的に、チップが接続されるネットワーク・メディアにより決定される。複数の外部 D R A M 及び S R A M が設けられ、チップにより使用される。

20

【0031】

ここでは特に、関連スイッチング及びルーティング装置の外部の一般データ・フローが、建物内に導入された配線やケーブルなどの電気導体を通過するネットワークが開示されるが、本発明は、ここで開示されるネットワーク・スイッチ及びコンポーネントが、無線環境において使用されることも同様に考慮する。例として、ここで述べられるメディア・アクセス制御 (M A C) 要素は、ここで述べられる要素を直接ワイヤレス・ネットワークにリンクする機能を有する、好適な無線周波数要素により置換されてもよく、例えば、既知のシリコン・ゲルマニウム技術などが使用される。こうした技術が適切に使用される場合には、無線周波数要素は当業者により、ここで開示される V L S I 構造内に統合される。或いは、無線周波数、または赤外線応答装置などの他のワイヤレス応答装置が、ワイヤレス・ネットワーク・システム用のスイッチ装置を実現する他の要素と一緒に、ブレード上に実装されてもよい。

30

【0032】

図中の矢印は、インタフェース装置内のデータの一般的な流れを示す。イーサネット (R) M A C から受信されるフレームは、E D S - U P により、内部データ・ストア・バッファに配置される。これらのフレームは、通常データ・フレームまたはシステム制御ガイド・フレーム (Guided Frame) として示され、E P C (図1) にエンキューされる。E P C は、最大 N フレーム ($N > 1$) を並列に処理できる N 個のプロトコル・プロセッサを含む。10 プロトコル・プロセッサ (図14) の実施例では、10 個のプロトコル・プロセッサの内の2個が特殊化される。すなわち、1つはガイド・フレームを処理するように (ガイド・セル・ハンドラ (G C H))、また他は制御メモリ内にルックアップ・データを作成するように (汎用ツリー・ハンドラ (G T H))、特殊化される。図13に示されるように、E P C は、新たなフレームをアイドル・プロセッサに突き合わせるディスパッチャと、フレーム・シーケンスを保持する完了 (completion) ユニットと、10 個の全プロセッサにより共用される共通命令メモリと、フレーム分類を決定する分類子ハードウェア補助機構及びフレームの開始命令アドレスを決定する支援をするコプロセッサと、フレーム・バッファの読み出し及び書き込み操作を制御する入口側 (以降イングレス (ingress) と称す) 及び出口側 (以降イーグレス (egress) と称す) データ・ストア・インタフェースと、10 個のプロセッサが制御メモリを共用することを可能にする制御メモリ・アービタと、ウェブ (Web) 制御と、内部インタフェース装置データ構造へのデバッグ・アクセ

40

50

スを可能にするアービタ及びインタフェースと、その他のハードウェア構造を含む。

【 0 0 3 3 】

ガイド・フレームは、G C Hプロセッサが使用可能になると、ディスパッチャにより G C Hプロセッサに送信される。レジスタ書込み、カウンタ読出し、イーサネット (R) M A C 構成変更など、ガイド・フレーム内でエンコードされる操作が実行される。M A C または I P エントリの追加など、ルックアップ・テーブル変更が、メモリ読出しや書込みなどの制御メモリ操作のために、ルックアップ・データ・プロセッサに渡される。M I B カウンタ読出しなどの幾つかのコマンドは、応答フレームの作成、及び適切なインタフェース装置上の適切なポートへの転送を要求する。ときとして、ガイド・フレームはインタフェース装置のイーグレス側に合わせてエンコードされる。これらのフレームは、照会されるインタフェース装置のイーグレス側に転送され、これが次にエンコード操作を実行し、適切な応答フレームを作成する。

10

【 0 0 3 4 】

データ・フレームは、次の使用可能なプロトコル・プロセッサにディスパッチされ、そこでフレーム参照が実行される。フレーム・データは、分類子ハードウェア補助機構 (C H A) エンジンからの結果と一緒に、プロトコル・プロセッサに渡される。C H A は I P または I P X を構文解析する。結果がツリー構造探索アルゴリズム及び開始共通命令アドレス (C I A) を決定する。サポートされるツリー構造探索アルゴリズムは、固定マッチ・ツリー (レイヤ 2 イーサネット (R) M A C テーブルなど、正確なマッチを要求する固定サイズ・パターン)、最長プレフィックス・マッチ・ツリー (サブネット I P 転送など、可変長マッチを要求する可変長パターン)、及びソフトウェア管理ツリー (フィルタ規則で使用されるような、範囲またはビット・マスク・セットのいずれかを定義する 2 つのパターン) を含む。

20

【 0 0 3 5 】

ルックアップ (探索) は、各プロトコル・プロセッサの一部である、ツリー構造探索エンジン (T S E) コプロセッサの支援により実行される。T S E コプロセッサは制御メモリ・アクセスを実行し、プロトコル・プロセッサを解放して実行を継続する。制御メモリは全てのテーブル、カウンタ、及びピココードにより必要とされる他のデータを記憶する。制御メモリ操作は、10個のプロセッサ・コンプレックスの間のメモリ・アクセスを調停する制御メモリ・アービタにより管理される。

30

【 0 0 3 6 】

フレーム・データはデータ・ストア・コプロセッサを通じてアクセスされる。データ・ストア・コプロセッサは、基本データ・バッファ (フレーム・データの最大 8 個の 16 バイト・セグメント) と、スクラッチ・パッド・データ・バッファ (フレーム・データの最大 8 個の 16 バイト・セグメント) と、データ・ストア操作のための幾つかの制御レジスタとを含む。一旦マッチが見いだされると、イングレス・フレーム変更が V L A N ヘッダ挿入またはオーバレイを含み得る。この変更は、インタフェース装置プロセッサ・コンプレックスにより実行されるのではなく、ハードウェア・フラグが導出され、他のイングレス・スイッチ・インタフェース・ハードウェアが変更を実行する。他のフレーム変更はピココード及びデータ・ストア・コプロセッサにより、イングレス・データ・ストア内に保持されるフレーム内容を変更することにより達成される。

40

【 0 0 3 7 】

他のデータは、フレームをスイッチ・ファブリック装置に送信する前に収集され、スイッチ・ヘッダ及びフレーム・ヘッダを作成するために使用される。制御データには、フレームの宛先ブレードなどのスイッチ情報や、イーグレス・インタフェース装置の情報が含まれ、宛先ポートのフレーム探索、マルチキャストまたはユニキャスト操作、及びイーグレス・フレーム変更を推進する。

【 0 0 3 8 】

完了時、エンキュー・コプロセッサが、フレームをスイッチ・ファブリックにエンキューするのに必要なフォーマットを作成し、それらを完了ユニットに送信する。完了ユニッ

50

トは、10個のプロトコル・プロセッサからスイッチ・ファブリック・キューへのフレーム順序を保証する。スイッチ・ファブリック・キューからのフレームは、64バイト・セルにセグメント化され、それらがPrisma-Eスイッチに伝送されるとき、フレーム・ヘッダ・バイト及びスイッチ・ヘッダ・バイトが挿入される。

【0039】

スイッチ・ファブリックから受信されたフレームは、イーグレスEDS34により、イーグレス・データ・ストア（イーグレスDS）バッファに配置され、EPCにエンキューされる。フレームの一部はディスパッチャにより、アイドル・プロトコル・プロセッサに送信され、フレーム探索が実行される。フレーム・データは、分類子ハードウェア補助機構からのデータと一緒に、プロトコル・プロセッサにディスパッチされる。分類子ハードウェア補助機構は、イーグレス・インタフェース装置により作成されるフレーム制御データを使用して、開始コード命令アドレス（CIA）の決定を支援する。

10

【0040】

イーグレス・ツリー探索は、イーグレス探索においてサポートされるのと同じのアルゴリズムをサポートする。ルックアップはTSEコプロセッサにより実行され、プロトコル・プロセッサを解放して実行を継続する。全ての制御メモリ操作は、制御メモリ・アービタにより管理され、アービタは10個のプロセッサ・コンプレックスの間で、メモリ・アクセスを割当てる。

【0041】

イーグレス・フレーム・データは、データ・ストア・コプロセッサを通じてアクセスされる。データ・ストア・コプロセッサは基本データ・バッファ（フレーム・データの最大8個の16バイト・セグメント）と、スクラッチ・パッド・データ・バッファ（フレーム・データの最大8個の16バイト・セグメント）と、データ・ストア操作のための幾つかの制御レジスタを含む。探索の成功結果は転送情報を含み、ときとしてフレーム変更情報を含む。フレーム変更はVLANヘッダ検出、存続時間増分（IPX）または減分（IP）、IPヘッダ・チェックサム再計算、イーサネット（R）フレームCRCオーバレイまたは挿入、及びMAC DA/SAオーバレイまたは挿入を含む。IPヘッダ・チェックサムはチェックサム・コプロセッサにより用意される。変更はインタフェース装置プロセッサ・コンプレックスにより実行されるのではなく、ハードウェア・フラグが作成され、PMMイーグレス・ハードウェアが変更を実行する。完了時に、エンキュー・コプロセッサが、フレームをEDSイーグレス・キューにエンキューし、それらを完了ユニットに送信するために必要なフォーマットを作成する。完了ユニットは、10個のプロトコル・プロセッサから、EDSイーグレス・キューへのフレーム順序を保証して、イーグレス・イーサネット（R）MAC36に供給する。

20

30

【0042】

完了フレームは最終的に、PMMイーグレス・ハードウェアによりイーサネット（R）MACに送信され、イーサネット（R）ポートから出力される。

【0043】

ウェブと呼ばれる内部バスは、内部レジスタ、カウンタ及びメモリへのアクセスを可能にする。ウェブ（Web）はまた、命令ステップや、デバッグ及び診断のための割り込み制御を制御する外部インタフェースを含む。

40

【0044】

ツリー構造探索エンジン・コプロセッサは、メモリ範囲チェック及び不当メモリ・アクセス通知を提供し、ツリー構造探索命令（メモリ読出し、書込み、またはread-add-writeなど）を、プロトコル・プロセッサ実行と並列に実行する。

【0045】

共通命令メモリは、1つの1024×128RAMと、2セットのデュアル512×128RAMとから構成される。デュアルRAMの各セットは、同一のピココードの2つのコピーを提供し、プロセッサによる同一アドレス範囲内の命令への独立したアクセスを可能にする。各128ビット・ワードは4つの32ビット命令を含み、合計8192命令を

50

提供する。

【 0 0 4 6 】

ディスパッチャは、10個のプロトコル・プロセッサへのフレームの受け渡しを制御し、割込み及びタイマを管理する。

【 0 0 4 7 】

完了ユニットは、プロセッサ・コンプレックスから、スイッチ・ファブリック及びターゲット・ポート・キューへのフレーム順序を保証する。豊富な命令セットには、条件付き実行、(入力ハッシュ・キーの)パッキング、条件付き分岐、符号付き及び符号無し演算、先行0のカウントなどが含まれる。

【 0 0 4 8 】

分類子ハードウェア補助機構エンジンは、各フレームのレイヤ2及びレイヤ3プロトコル・ヘッダを構文解析し、フレームがプロトコル・プロセッサにディスパッチされるとき、この情報をフレームに提供する。

【 0 0 4 9 】

制御メモリ・アービタは、内部及び外部メモリの両方へのプロセッサ・アクセスを制御する。

【 0 0 5 0 】

外部制御メモリ・オプションは、5個乃至7個のDDR DRAMサブシステムを含む。各サブシステムは、1対の2M×16ビット×4バンクDDR DRAM、または1対の4M×16ビット×4バンクDDR DRAMをサポートする。DDR DRAMインタフェースは、133MHzクロック・レート及び266MHzデータ・ストロークで動作し、構成可能なCAS待ち時間及びドライブ強度をサポートする。オプションの133MHz ZBT SRAMが、128K×36、2×256K×18、または2×512K×18構成のいずれかで追加される。

【 0 0 5 1 】

イーグレス・フレームは1つの外部データ・バッファ(例えばDS0)、または2つの外部データ・バッファ(DS0及びDS1)のいずれかに記憶される。各バッファは、1対の2M×16ビット×4バンクDDR DRAM(最大256Kの64バイト・フレームを記憶可能)、または1対の4M×16ビット×4バンクDDR DRAM(最大512Kの64バイト・フレームを記憶可能)から構成される。2.28Mbpsの単一の外部データ・バッファ(例えばDS0)を選択するか、第2のバッファ(例えばDS1)を追加し、4.57Mbpsのレイヤ2及びレイヤ3スイッチングをサポートする。第2のバッファの追加は性能を改善するが、フレーム容量を増加しない。外部データ・バッファ・インタフェースは、133MHzクロック・レート、及び266MHzのデータ・ストロークで動作し、構成可能なCAS待ち時間及びドライブ強度をサポートする。

【 0 0 5 2 】

内部制御メモリは、2つの512×128ビットRAMと、2つの1024×36ビットRAMと、1つの1024×64ビットRAMとを含む。

【 0 0 5 3 】

内部データ記憶は、イングレス方向(UP)に、最大2048個の64バイト・フレームをバッファリングすることができる。

【 0 0 5 4 】

固定フレーム変更は、イングレス方向のVLANタグ挿入と、VLANタグ削除と、持続時間増分/減分(IP、IPX)と、イーグレス(DOWN)方向のイーサネット(R)CRCオーバレイ/挿入及びMAC DA/SAオーバレイ/挿入とを含む。

【 0 0 5 5 】

ポート・ミラーリングは、プロトコル・プロセッサ資源を用いることなく、1つの受信ポート及び1つの送信ポートを、システム指定の観測ポートにコピーすることを可能にする。ミラーリングされたインタフェース装置ポートは、フレーム及びスイッチ制御データを追加するように構成される。別々のデータ・バスが、イングレス・スイッチ・インタフ

10

20

30

40

50

エースへの直接フレーム・エンキューイングを可能にする。

【 0 0 5 6 】

インタフェース装置は、4つのイーサネット(R)マクロを統合する。各マクロは、1ギガビットまたは10/100高速イーサネット(R)モードのいずれかで動作するように、個々に構成される。各イーサネット(R)マクロは、最大10個の10/100Mbps MACか、4つのマクロの各々に対して、1つの1000Mbps MACをサポートする。

【 0 0 5 7 】

図2は、MACコアのブロック図を示す。各マクロは、3つのイーサネット(R)コア設計、すなわちマルチポート10/100Mbps MACコア(FEnet)、1000Mbps MACコア(GEnet)、及び100Mbps物理コーディング・サブレイヤ・コア(PCs)を含む。

10

【 0 0 5 8 】

マルチポート・イーサネット(R)10/100MACフィーチャ:

1) 物理レイヤとの10個のシリアル・メディア独立インタフェースをサポートする。
2) 10Mbpsまたは100Mbpsメディア速度の10個のポートを、任意に混在させて処理できる。

3) 単一のMACが時分割多重インタフェースにより、10個の全てのポートをサービスする。

4) 全てのポート上で全2重/半2重動作をメディア速度でサポートする。

20

5) IEEE 802.3バイナリ指数バックオフをサポートする。

【 0 0 5 9 】

1000Mbpsイーサネット(R)MACコア・フィーチャ:

1) 物理PCSレイヤとの、または直接的に物理レイヤとのギガビット・メディア独立インタフェース(GMII)をサポートする。

2) PCSコアにより、完全なTBI(8b/10b)ソリューションをサポートする。

3) 全2重Point-to-Point接続をメディア速度でサポートする。

4) IBM PCSコア有効バイト信号方式をサポートする。

【 0 0 6 0 】

1000Mbpsイーサネット(R)物理コーディング・サブレイヤ・コア・フィーチャ:

30

1) 8b/10bエンコード及びデコードを実行する。

2) IEEE 802.3zで定義されるPMA(10ビット)サービス・インタフェースをサポートする。このインタフェースは、IEEE 802.3zに準拠する任意のPMAに接続する。

3) PMAから受信されるデータ(2フェーズ・クロック)を、MAC(1フェーズ)クロックに同期させる。

4) 次の2ページを含むオートネゴシエーションをサポートする。

5) 規格で定義された2フェーズ・クロック・システムを、1フェーズ・クロックに変換する。

40

6) 新たなデータを含むクロック・サイクルを示す信号を、MACに提供する。

7) 受信コード・グループ(10ビット)内のCOMMAをチェックし、ワード同期を確立する。

8) 8b/10b実行中ディスパリティを計算及びチェックする。

【 0 0 6 1 】

図3の(A)乃至(D)は、インタフェース素子チップの異なる構成を示す。これらの構成はDASLと、スイッチング・ファブリック装置への接続とによりサポートされる。各DASLは2つのチャネル、すなわち送信チャネルと受信チャネルとを含む。

【 0 0 6 2 】

図3の(A)は、単一インタフェース装置のラップ(wrap)構成を示す。この構成では

50

、送信チャネルが受信チャネルにラップ、すなわち折り返される。

【 0 0 6 3 】

図 3 の (B) は、2 つのインタフェース素子チップが接続される構成を示す。各インタフェース素子チップは、少なくとも 2 つの D A S L を提供される。この構成では、1 つのチップ上の 1 D A S L 上のチャネルが、他のチップ上のマッティング D A S L のチャネルに、動作可能なように接続される。各チップ上の他の D A S L は、ラップされる。

【 0 0 6 4 】

図 3 の (C) は、複数のインタフェース装置がスイッチ・ファブリックに接続される構成を示す。両頭矢印は両方向の伝送を示す。

【 0 0 6 5 】

図 3 の (D) は、メイン・スイッチ及びバックアップ・スイッチが複数のインタフェース装置に接続される構成を示す。メイン・スイッチが故障すると、バックアップ・スイッチが使用可能になる。

【 0 0 6 6 】

制御点 (C P) はシステム・プロセッサを含み、これは各構成に接続される。とりわけ、C P のシステム・プロセッサは、チップに対して、初期化及び構成サービスを提供する。C P は 3 つの位置、すなわちインタフェース素子チップ内か、チップが実装されるブレード上か、或いはブレードの外部のいずれかに配置される。ブレードの外部の場合、C P はリモートとなり、どこか別の場所に内蔵され、インタフェース装置及び C P が接続されるネットワークを介して通信する。C P の要素が図 2 0 に示され、メモリ素子 (キャッシュ、フラッシュ及び S D R A M)、メモリ制御装置、P C I バス、及びバックプレーン及び L 1 ネットワーク・メディア用のコネクタを含む。

【 0 0 6 7 】

図 2 1 は、単一チップ・ネットワーク・プロセッサと、E D S - U P、トラフィック管理 (M G T) スケジューラ、及び E D S - D O W N (D N) により提供される機能とを示す。U 字形アイコンはキューを表し、キュー内の内容を追跡する制御ブロック (C B) は、矩形アイコンにより表される。

【 0 0 6 8 】

各要素及びそれらの機能及び相互作用は、次の通りである。

【 0 0 6 9 】

P M M : これは M A C (F E n e t、P O S、G E n e t) を含み、外部 P H Y 装置に接続されるネットワーク・プロセッサの一部である。

【 0 0 7 0 】

U P - P M M : この論理は P H Y からバイトを受け取り、それを F I S H (1 6 バイト) にフォーマットし、U P - E D S に受け渡す。P M M 内には 4 つの D M U が存在し、各々は 1 G E n e t または 1 0 F E n e t 装置と連携することができる。

【 0 0 7 1 】

U P - E D S : この論理は U P - P M M から F I S H を受け取り、それらを U P - D A T A ストア (内部 R A M) に記憶する。これは 1 度に 4 0 フレームを処理することができ、適切な数のバイトが受信されると、フレームを E P C にエンキューする。E P C がフレーム処理を完了すると、U P - E D S がフレームを適切なターゲット・ポート・キューにエンキューし、U P - S D M にフレームの送信を開始する。U P - E D S は全てのバッファ及びフレーム管理の責任を負い、U P - S D M への転送が完了するとき、バッファ / フレームを空きプールに戻す。

【 0 0 7 2 】

E P C : この論理はピコプロセッサを含み、組み込み P o w e r P C を含み得る。この論理はフレーム・ヘッダを突き止め、フレームをどのように処理すべきかを決定する (転送、変更、フィルタリングなど)。E P C は幾つかのルックアップ・テーブルをアクセスでき、ピコプロセッサがネットワーク・プロセッサの高帯域幅要求に応じるためのハードウェア補助機構を有する。

10

20

30

40

50

【 0 0 7 3 】

UP - SDM : この論理はフレームを受け取り、それらをスイッチ・ファブリックへの伝送のために、Prismaセルにフォーマットする。この論理はVLANヘッダをフレームに挿入することができる。

【 0 0 7 4 】

UP - SIF : この論理はUP - DASLマクロを含み、外部スイッチ入出力 (I / O) に接続する。

【 0 0 7 5 】

DN - SIF : この論理はDN - DASLマクロを含み、外部スイッチ入出力 (I / O) からPrismaセルを受信する。

10

【 0 0 7 6 】

DN - SDM : この論理はPrismaセルを受信し、フレーム再組み立てのために、それらを事前処理する。

【 0 0 7 7 】

DN - EDS : この論理は各セルを受け取り、それらを再度フレームに組み立てる。セルは外部データ・ストアに記憶され、バッファと一緒にリンクされて、フレームを形成する。全フレームが受信されると、フレームはEPCにエンキューされる。EPCがフレーム処理を終えると、フレームはスケジューラ・キュー (但し存在する場合) またはターゲット・ポート・キューにエンキューされる。DN - EDSは次に、フレーム、変更情報、及び制御情報などをDN - PMMに送信することにより、フレームを適切なポートに送信する。

20

【 0 0 7 8 】

DN - PMM : DN - EDSから情報を受け取り、フレームをイーサネット (R) 、POSなどにフォーマットし、フレームを外部PHYに送信する。

【 0 0 7 9 】

SPM : この論理は、ネットワーク・プロセッサが外部装置 (PHY、LED、FLASHなど) とインタフェースするために使用されるが、3つの入出力 (I / O) を要求するだけである。ネットワーク・プロセッサはシリアル・インタフェースを用いて、SPMと通信し、次にSPMがこれらの外部装置を管理するために必要な機能を実行する。

30

【 0 0 8 0 】

アップサイド・フロー :

1) フレームがPHYに到来する。
2) バイトがUP - PMMにより受信される。
3) UP - PMMがFISHをUP - EDSに送信する (FISHはフレームの一部を意味する) 。

4) UP - EDSがFISHをUP - DSに記憶する。

5) UP - EDSがヘッダをEPC送信する。

6) EPCがヘッダを処理し、エンキュー情報をUP - EDSに返送する。

7) UP - EDSがフレームの残りをUP - PMMから受信し続ける。

8) スwitchへの適切なデータの送信準備が整うと、UP - EDSが情報をUP - SDMに送信する。

40

9) UP - SDMがフレーム・データを読み出し、それをPrismaセルにフォーマットする。

10) UP - SDMがセルをUP - SIFに送信する。

11) UP - SIFがDASLシリアル・リンクを介して、セルをPrismaに転送する。

12) 全てのデータが受け取られると、UP - EDSがバッファ/フレームを解放する。

【 0 0 8 1 】

ダウンサイド・フロー :

50

- 1) DN - S I F が P r i z m a セルを受信する。
- 2) DN - S D M がセルを記憶し、それらを再組み立て情報として事前処理する。
- 3) DN - E D S がセル・データ及び再組み立て情報を受信し、セルをダウンサイド側の新たなフレームにリンクする。
- 4) DN - E D S がセルを DN - D S に記憶する。
- 5) 全てのデータが受信されると、DN - E D S がフレームを E P C にエンキューする。
- 6) E P C がヘッダを処理し、エンキュー情報を DN - E D S に返送する。
- 7) DN - E D S がフレームをスケジューラ・キュー（但し存在する場合）またはターゲット・ポート・キューにエンキューする。
- 8) DN - E D S がキューをサービスし、フレーム情報を P C B に送信する。
- 9) DN - E D S が P C B を用いてフレームを"解体" (unravel) し、適切なデータを読出し、そのデータを DN - P M M に送信する。
- 10) DN - P M M がデータをフォーマットし（必要に応じて変更を加える）、フレームを外部 P H Y に送信する。
- 11) DN - P M M が DN - E D S に、バッファがもはや必要とされないことを通知し、DN - E D S がこれらの資源を解放する。

【 0 0 8 2 】

フレーム制御フロー：

- 1) ヘッダが U P - D S または DN - D S から E P C に送信される。
- 2) E P C がルックアップ・テーブルでヘッダ情報を調査し、フレーム・エンキュー情報を受信する。
- 3) E P C がエンキュー情報を E D S に返送し、フレームが適切なキューにエンキューされる。
- 4) セル・ヘッダ及びフレーム・ヘッダがフレーム・データと一緒に送信され、再組み立て及びフレーム転送を支援する。

【 0 0 8 3 】

C P 制御フロー：

- 1) 制御点がガイド・フレームをフォーマットし、それをネットワーク・プロセッサに送信する。
- 2) ネットワーク・プロセッサがガイド・フレームを、G C H ピコプロセッサにエンキューする。
- 3) G C H がガイド・フレームを処理し、Rainierの要求領域を読み書きする。
- 4) G C H がテーブル更新要求を G T H に渡す。
- 5) G T H が適切なテーブルを、ガイド・フレームからの情報により更新する。
- 6) 肯定応答ガイド・フレームが C P に返送される。

【 0 0 8 4 】

ネットワーク・プロセッサ制御フロー：

- 1) ピコプロセッサがガイド・フレームを作成し、情報を別のRainierまたは制御点に送信する。
- 2) ガイド・フレームが適切な位置に送信され処理される。

【 0 0 8 5 】

単一インタフェース装置は、最大 40 の高速イーサネット (R) ポート (図 3 の (A) 参照) のための、メディア速度切替えを提供する。I B M のデータ整合同期リンク (D A S L : Data Aligned Synchronous Link) 技術を用いて、2 つのインタフェース装置が相互接続される場合、80 個の高速イーサネット (R) ポートがサポートされる (図 3 の (B) 参照) 。各 D A S L 差動対は、440 M b p s のデータを伝搬する。従って、8 対の 2 つのセットが、3.5 G b p s 全 2 重接続 (各方向において 440 M b p s の 8 倍) を提供する。図 3 の (C) 及び (D) に示されるように、複数のインタフェース装置を、I B M の P r i z m a - E スイッチなどのスイッチに相互接続することにより、より大規模

10

20

30

40

50

なシステムが構成される。インタフェース装置は、2つの3.5 Gbps全2重DASL接続を提供し、これらの一方が基本接続で他が補助接続であり、後者は、ローカル・フレーム・トラフィックのためのラップ・バックパス(wrap-backpath)を提供するか(2つのインタフェース装置が直接接続される場合(図3の(B))、或いは、冗長スイッチ・ファブリックへの接続を提供する(図3の(D)、バックアップ・スイッチ)。以上から、単一ネットワーク・プロセッサ・チップは、1つのチップがローエンド・システム(比較的低いポート密度(例えば40)を有する)から、ハイエンド・システム(比較的高いポート密度(例えば80)を有する)まで提供するために使用されるという点でスケラブルである。

【0086】

10

システム内の1つのインタフェース装置が、最大10個の10/100Mbps高速イーサネット(R)ポートか、単一の1000Mbpsイーサネット(R)ポートを介して、システム・プロセッサに接続される。システム・プロセッサへのイーサネット(R)構成は、インタフェース装置に接続されるEEPROM内に記憶され、初期化の間にロードされる。システム・プロセッサは、イーサネット(R)フレームとしてカプセル化される特殊ガイド・フレームを作成することにより、システム内の全てのインタフェース装置と通信する(図3参照)。カプセル化されたガイド・フレームは、DASLリンクを介して他の装置に転送され、システム内の全てのインタフェース装置が、単一点から制御されることになる。

【0087】

20

ガイド・フレームは、制御点(CP)と組み込みプロセッサ・コンプレックスとの間で、及びインタフェース装置内で制御情報を伝達するために使用される。ここでの議論を明らかにするガイド・セルの従来の開示が、1998年3月3日発行の米国特許第5724348号"Efficient Hardware/Software Interface for a Data Switch"で述べられている。

【0088】

CPから発信されるガイド・フレーム・トラフィックでは、CPがそのローカル・メモリ内のデータ・バッファ内に、ガイド・フレームを作成する。CPのデバイス・ドライバがガイド・フレームを、ネットワーク・プロセッサのメディア・インタフェースの1つに送信する。メディア・アクセス制御(MAC)ハードウェアがガイド・フレームを回復し、それを内部データ・ストア(UDS)メモリに記憶する。ガイド・フレームは適切なブレードに経路指定され、処理され、必要に応じて再度CPに経路指定される。外部CPとインタフェース装置との間でやり取りされるガイド・フレームは、外部ネットワークのプロトコルに適應するようにカプセル化される。その結果、外部ネットワークがイーサネット(R)を含む場合、ガイド・フレームはイーサネット(R)フレームとしてカプセル化される。

30

【0089】

イーサネット(R)カプセル化は、CPとインタフェース装置との間のガイド・トラフィックのトランスポート手段を提供する。インタフェース装置のイーサネット(R)MAC(Enet MAC)は、フレームの受信時に、宛先アドレス(DA)またはソース・アドレス(SA)を分析しない。この分析はEPCピココードにより実行される。ガイド・トラフィックは、インタフェース装置が構成されておらず、DA及びSAがEPCピココードにより分析されないものと仮定する。従って、これらのフレームは元来、自己ルーティングである。しかしながら、イーサネット(R)MACは、イーサネット(R)タイプ・フィールドを分析し、ガイド・トラフィックとデータ・トラフィックとを区別する。ガイド・フレームのこのイーサネット(R)タイプ値の値は、E_Type_Cレジスタにロードされる値と合致しなければならない。このレジスタは、インタフェース装置のブート・ピココードにより、フラッシュ・メモリからロードされる。

40

【0090】

CPはそのローカル・メモリ内のデータ・バッファ内に、ガイド・フレームを構築する

50

。C P プロセッサ内の 3 2 ビット・レジスタの内容が、図 4 に示されるように、ビッグ・エンディアン形式でローカル・メモリ内に記憶される。ガイド・フレームが構築されると、C P のデバイス・ドライバがイーサネット (R) フレームを送信する。このフレームは、特定のガイド・セル・ハンドラ (G C H) の D A 、C P の大域 M A C アドレスまたは特定インタフェースの M A C アドレスに対応する S A 、ガイド・フレームを示す特殊イーサネット (R) タイプ・フィールド、及びガイド・フレーム・データを含む。ポートに到来する全てのイーサネット (R) フレームは、イーサネット (R) M A C により受信され、分析される。E_Type_C レジスタの内容に合致するイーサネット (R) タイプ値を有するフレームの場合、イーサネット (R) M A C が D A 、S A 及びイーサネット (R) タイプ・フィールドを剥ぎ取り、ガイド・フレーム・データを U _ D S メモリに記憶する。バイトが 1 つずつ、イーサネット (R) M A C により、F I S H と呼ばれる 1 6 バイトのブロック内に収集される。これらのバイトはビッグ・エンディアン形式で記憶され、ガイド・フレームの第 1 バイトが、F I S H の最上位バイト位置 (バイト 0) に記憶される。続くバイトは、F I S H 内の続くバイト位置 (バイト 1 , バイト 2 , . . . , バイト 1 5) に記憶される。これらの 1 6 バイトは次に、U _ D S 内のバッファに、F I S H 0 位置から記憶される。続く F I S H は、バッファ内の連続する F I S H 位置 (F I S H 1 、F I S H 2 、F I S H 3 など) に記憶される。ガイド・フレームの残りを記憶するために、必要に応じて、空きプールから追加のバッファが獲得される。

【 0 0 9 1 】

インタフェース装置 1 0 内のガイド・トラフィックのフローが、図 5 に示される。インタフェース装置のイーサネット (R) M A C 機能が、フレーム・ヘッダ情報を調査し、フレームがガイド・フレームであることを判断する。イーサネット (R) M A C はガイド・フレームからフレーム・ヘッダを除去し、その残りの内容をインタフェース装置の内部 U _ D S メモリにバッファリングする。イーサネット (R) M A C は、フレームが G C H により処理されるために、汎用制御 (G C) キューにエンキューされるべきことを示す。ガイド・フレームの終わりに達すると、エンキュー、デキュー及びスケジュール (E D S) 論理が、フレームを G C キューにエンキューする。

【 0 0 9 2 】

C P に局所的に接続されるブレード上の G C H ピココードは、フレーム制御情報 (図 7 参照) を調査し、ガイド・フレームがシステム内の他のブレードに向けられるか、及びガイド・フレームがインタフェース装置のダウンスайд側で実行されるべきかを判断する。フレームが、局所的に接続されるブレード以外のブレード、またはそれに加えたブレードに向けられる場合、G C H ピココードがフレーム制御ブロック (F C B) 内の T B 値を、ガイド・フレームのフレーム制御情報からの T B 値により更新し、E D S にフレームを、マルチキャスト・ターゲット・ブレード・フレーム開始 (T B _ S O F) キューにエンキューするように命令する。性能上の理由から、示された宛先ブレードの数とは無関係に、全てのガイド・トラフィックがマルチキャスト T B _ S O F キューにエンキューされる。

【 0 0 9 3 】

フレームが局所的に接続されたブレードだけに向けられる場合、G C H ピココードがフレーム制御情報の U P / D O W N パー・フィールド (以降、'パー'は'0' b で有効なことを意味する) (図 7 参照) を調査し、ガイド・フレームがインタフェース装置のアップサイドまたはダウンスайдのどちらで実行されるべきかを判断する。ガイド・フレームがインタフェース装置のダウンスайдで実行される場合、G C H ピココードが F C B 内の T B 値を、ガイド・フレームのフレーム制御情報からの T B 値により更新し、E D S にフレームをマルチキャスト T B _ S O F キューにエンキューするように命令する。フレーム制御情報が、ガイド・フレームがアップサイドで実行されるべきことを示す場合、G C H ピココードがガイド・フレームを分析し、そこに含まれるガイド・コマンドにより示される操作を実行する。

【 0 0 9 4 】

ガイド・コマンドを処理する前に、ピココードはフレーム制御情報の A C K / N O A C

Kバー・フィールドの値をチェックする、この値が'0'bの場合、ガイド・フレームが処理に続き廃棄される。但し、ガイド読出しコマンドは、この範疇ではない。

【0095】

ACK/NOACKバー・フィールドが'1'bで、EARLY/LATEバー・フィールドが'1'bの場合、ガイド・フレーム内の任意のガイド・コマンドを処理する前に、ピココードが早期ACKガイド・フレームを作成する。このとき、フレーム制御のTBフィールドの値は、早期ACKガイド・フレームの内容に等しく、フレーム制御のTBフィールドの値は、My__TBレジスタの内容に等しい。ピココードは、フレームのFCB内のTB値を、LAN制御点アドレス(LAN_CP_Addr)レジスタのTBフィールドに含まれる値により更新し、EDSにフレームをマルチキャストTB__SOFキューにエンキューするように命令することにより、早期ACKガイド・フレームをCPに返送する。ピココードは次に、ガイド・フレームのガイド・コマンドを処理し、ガイド・フレームを廃棄する。但し、ガイド読出しコマンドはこの範疇ではない。

10

【0096】

他方、ACK/NOACKバー・フィールドの値が'1'bで、EARLY/LATEバー・フィールドの値が'0'bの場合、ピココードはフレーム制御情報のRESP/REQバー・フィールドを'1'bに変更して、ガイド・フレーム応答を示し、TBフィールドをMy__TBレジスタの内容により更新し、ガイド・フレーム内の各ガイド・コマンドを処理する。ガイド・コマンドの処理の間、ピココードは次のガイド・コマンドの完了コード・フィールドを、現ガイド・コマンドの完了ステータス・コード値により更新する。ピココードは、FCB内のTB値をソース・ブレードに対応する値(CPのLAN_CP_Addrレジスタ値)により更新し、EDSにフレームをマルチキャストTB__SOFキューにエンキューするように命令することにより応答をソースに返送する。

20

【0097】

TB__SOFキュー内に存在するフレームは、EDSにより、転送をスケジュールされる。スイッチ・データ・ムーバ(SDM)が、FCBに含まれる情報から、スイッチング・ファブリック・セル・ヘッダ及びインタフェース装置フレーム・ヘッダを作成する。これらのセルはスイッチング・ファブリック装置を通過し、ターゲット・ブレードに到達し、そこでセルがD-Dメモリ内で、フレームに再組み立てされる。ダウンサイドのSDMは、フレームがガイド・フレームであることを認識し、EDSにそれをGCキューにエンキューするように伝える。

30

【0098】

GCキューまたはGTキューからの圧力により、ピココードはガイド・フレームをアクセスし分析する。ダウンサイドに到来する全てのガイド・フレームは、最初にGCキューにエンキューされる。これらのフレームにおけるフレーム制御情報のGTH/GCHバー値が、GCHピココードにより調査される。GTH/GCHバー値が'0'bの場合、ガイド・フレームがGTキューにエンキューされる。それ以外では、GCHピココードが、フレーム制御情報のRESP/REQバー・フィールドを調査し、ガイド・フレームが既に実行されたか否かを判断する。RESP/REQバー・フィールドが値'1'bを有する場合、ガイド・フレームは既に実行されており、CPに経路指定される。CP接続に対応するターゲット・ポート値は、EPCピココードにより保持される。これらのターゲット・ポート・キューからのフレームが、インタフェース装置からCPに返送される。

40

【0099】

RESP/REQバー・フィールドが値'0'bを有する場合、ブレードはCPに対してローカルまたはリモートである。これはLAN_CP_AddrレジスタのTBフィールドの値を、マイ・ターゲット・ブレード(My__TB)レジスタの内容と比較することにより解明される。これらが一致する場合、ブレードはCPにとってローカルであり、それ以外では、ブレードはCPにとってリモートである。いずれの場合にも、ピココードはフレーム制御情報のUP/DOWNバー値を調査する。UP/DOWNバー値が'1'bの場合、フレームはラップTPキューにエンキューされ、U__DSへ転送されて、アップサイドのGCH

50

により処理される。それ以外では、ピココード（GCHまたはGTH）が、ガイド・フレームに含まれるガイド・コマンドにより指示される操作を実行する。ガイド・コマンドの処理に先立ち、ピココードがフレーム制御情報のACK/NOACKバー・フィールドの値をチェックする。この値が'0'bの場合、ガイド・フレームが処理に続き廃棄される。但し、ガイド読出しコマンドはこの範疇ではない。

【0100】

ACK/NOACKバー・フィールドの値が'1'bで、EARLY/LATEバー・フィールドの値が'1'bの場合、ガイド・フレーム内の任意のガイド・コマンドを処理する前に、ピココードが早期ACKガイド・フレームを作成し、このとき、フレーム制御情報のTBフィールドの値は、My__TBレジスタの内容に等しい。ブレードがCPにとってリモートの場合、ピココードは早期ACKガイド・フレームをラップ・ポートに経路指定する。それ以外では、ブレードはCPにとってローカルであり、フレームはCPに対応するポート・キューに経路指定される。ピココードはガイド・コマンドを処理する一方、ラップ・ポートが早期ACKガイド・フレームをD__DSからU__DSに転送し、フレームをアップサイドのGCキューにエンキューするか、フレームがポート・キューからCPに返送される。U__DSにラップバック（wrap back）されるフレームに対して、GCHピココードは再度このフレームを調査するが、RESP/REQバー・フィールドは値'1'bを有する。GCHピココードは、FCB内のTBフィールドを、LAN_CP_AddrレジスタのTBフィールドに含まれる値により更新し、EDSにフレームをマルチキャストTB__SOFキューにエンキューするように命令することにより、フレームをCPに返送する。TB__SOFキュー内に存在するフレームは、EDSにより、転送をスケジュールされる。スイッチ・データ・ムバ（SDM）が、FCBに含まれる情報から、Prizmaセル・ヘッダ及びインタフェース装置フレーム・ヘッダを作成する。このフレームからのセルはPrizmaを通過し、CPのローカル・ブレード上で、フレームに再組み立てされる。ダウンサイドのSDMは、フレームがガイド・フレームであることを認識し、EDSにそれをGCキューにエンキューするように伝える。GCHピココードがフレームを分析するとき、RESP/REQバー・フィールドは値'1'bを有する。このことは、このブレードがCPに局所的に接続されることを意味し、従ってガイド・フレームは、CPに対応するポート・キューに経路指定される。このキュー内のフレームは、インタフェース装置からCPに返送される。

【0101】

他方、ACK/NOACKバー・フィールドの値が'1'bで、EARLY/LATEバー・フィールドの値が'0'bの場合、ピココードはRESP/REQバー・フィールドを'1'bに変更し、ガイド・フレーム応答であることを示し、TBフィールドをMy__TBレジスタの内容により置換し、次にガイド・フレーム内の各ガイド・コマンドを処理する。ガイド・コマンドの処理の間、ピココードは次のガイド・コマンドの完了コード・フィールドを、現ガイド・コマンドの完了ステータス・コード値により更新する。ブレードがCPにとってリモートの場合、ピココードはガイド・フレームをラップ・ポートに経路指定する。それ以外では、ブレードはCPにとってローカルであり、フレームはCPに対応するポート・キューに経路指定される。ラップ・ポートがガイド・フレームをD__DSからU__DSに転送し、フレームをアップサイドのGCキューにエンキューするか、フレームがポート・キューからCPに返送される。U__DSにラップバックされるフレームに対して、GCHピココードは再度このフレームを調査するが、RESP/REQバー・フィールドは値'1'bを有する。GCHピココードは、FCB内のTBフィールドを、LAN_CP_AddrレジスタのTBフィールドに含まれる値により更新し、EDSにフレームをマルチキャストTB__SOFキューにエンキューするように命令することにより、フレームをCPに返送する。TB__SOFキュー内に存在するフレームは、EDSにより、転送をスケジュールされる。SDMが、FCBに含まれる情報から、Prizmaセル・ヘッダ及びインタフェース装置フレーム・ヘッダを作成する。このフレームからのセルはPrizmaを通過し、CPのローカル・ブレード上のダウンサイドでフレームに再組み立てされる

10

20

30

40

50

。ダウンサイドのSDMは、フレームがガイド・フレームであることを認識し、EDSにそれをGCキューにエンキューするように伝える。GCHピココードがD_DSからのフレームを分析するとき、RESP/REQバー・フィールドは値'1'bを有する。このことは、このブレードがCPに局所的に接続されることを意味し、ガイド・フレームが、CPに対応するポート・キューに経路指定される。このキュー内のフレームは、インタフェース装置からCPに返送される。

【0102】

何らかの理由により、GCHピココードが、フレーム制御情報のTBフィールドが'0000'hに等しいガイド・フレームに遭遇する場合、GCHピココードはこのフレームを、このブレードだけに向けられるものと解釈し、それにもとづき処理を行う。このアクションは、My_TBレジスタの値が全てのブレードに対して'0000'hである初期化の間に要求される。CPは、フレーム制御情報が'0000'hのTB値を有するガイド・フレーム内の、書込みガイド・コマンドを送信することにより、局所的に接続されるブレードのMy_TBレジスタを初期化する。

10

【0103】

EPC内の任意のピコプロセッサがガイド・フレームを生成できる。このフレームは非送信請求ガイド・フレームか、他の形式のガイド・フレームである。このタイプの内部的に生成されるフレームは、肯定応答を許可しないように生成される（すなわちACK/NOACKバー='0'b）。これらのフレームは、同一のEPC内の2つのピコプロセッサの一方（GCHまたはGTH）に送信されるか、他のブレードのGCHまたはGTHに送信される。

20

【0104】

非送信請求ガイド・フレームはCPにも送信され得る。同一のEPCに向けられるガイド・フレームは、D_DS内のデータ・バッファを用いて構成される。これらのフレームは次に処理のために、GCまたはGTキューにエンキューされる。これらのフレームは次に通常通り処理され、廃棄される。局所的に接続されるCPに向けられる非送信請求ガイド・フレームは、D_DS内のデータ・バッファを用いて構成される。これらのフレームは、それらがEPCにより実行されたことを示すように構成される（すなわち、RESP/REQバー='1'b、TB=My_TB）。これらのフレームは、CPに対応するポート・キューにエンキューされる。このキューからフレームは、CPに返送される。

30

【0105】

別のブレードに向けられるガイド・フレームは、D_DSまたはU_DS内のデータ・バッファを用いて構成される。CPに向けられる非送信請求ガイド・フレームは、それらがEPCにより実行されたことを示すように構成される（すなわち、RESP/REQバー='1'b、TB=My_TB）。D_DS内のデータ・バッファを用いて構成されたフレームは、ラップ・ポートにエンキューされる。これらのフレームはU_DSに転送され、アップサイドのGCキューにエンキューされる。'1'bのRESP/REQバー値を有する非送信請求ガイド・フレームは、LAN_CP_Addrレジスタ内のTB値を用いてCPに経路指定される。それ以外では、GCHピココードが、ガイド・フレームのフレーム制御情報のTB値を用いて、これらのフレームを経路指定する。受信ブレードにおいて、フレームがダウンサイドのGCキューにエンキューされる。このブレードのGCHは、フレームを実行し廃棄するか（RESP/REQバー='0'b、GTH/GCHバー='1'）、フレームをGTキューにエンキューするか（RESP/REQバー='0'b、GTH/GCHバー='0'）、フレームをCPに対応するポート・キューにエンキューする（RESP/REQバー='1'b）。U_DS内のデータ・バッファを用いて構成されたフレームは、アップサイドのGCキューに直接エンキューされる。この点から以降、これらのフレームは同一の経路に従い転送され、D_DSデータ・バッファを用いて構成されたフレームと同様に処理される。図6は、ガイド・フレームの汎用フォーマットを示す。

40

【0106】

図示のフォーマットは、左側に最上位バイトを有し、右側に最下位バイトを有する論理

50

表現である。4 バイト・ワードが先頭のワード 0 から開始し、ページの終わりに向けて増加する。

【 0 1 0 7 】

インタフェース装置が C P により構成される前に、ガイド・フレームは経路指定され、処理されなければならないので、これらのフレームは自己ルーティングでなければならない。探索及び分類により通常獲得される結果が、ガイド・フレームのこのフレーム制御情報フィールド内に含まれ、チップが探索操作を実行することなく、F C B をこの情報により更新することを可能にする。ガイド・フレームに含まれるターゲット・ブレード情報は、ガイド・フレーム・ハンドラにより、F C B のリーフ・ページ・フィールドを用意するために使用される。C P がターゲット・ブレード情報を提供する一方、G C H ピココードが F C B 内の他のフィールドを記入する。この F C B 情報は S D M により、セル・ヘッダ及びフレーム・ヘッダを用意するために使用される。ガイド・フレームのフレーム制御情報フィールドのフォーマットが図 7 に示される。

10

【 0 1 0 8 】

次に、図 7 の各ビット位置の略語について説明する。

【 0 1 0 9 】

R E S P / R E Q バー：応答及び非要求バー標識値。このフィールドは、要求（未処理）ガイド・フレームと応答ガイド・フレームとを区別するために使用される。

0：要求

1：応答

20

【 0 1 1 0 】

A C K / N O A C K バー：肯定応答または無肯定応答制御値。このフィールドは、G C H ピココードがガイド・フレームを肯定応答するか否かを制御するために使用される（肯定応答する場合 A C K、そうでない場合 N O A C K）。ガイド・フレームが肯定応答されない場合、読出しを実行するいずれの形式のガイド・コマンドも含まれない。

0：無肯定応答

1：肯定応答

【 0 1 1 1 】

E A R L Y / L A T E バー：早期及び遅延肯定応答制御値。このフィールドは、要求される肯定応答（A C K / N O A C K バー = ' 1 ' b）が、ガイド・フレームが処理される前に発生するか（E A R L Y）、または後に発生するか（L A T E）を制御するために使用される。A C K / N O A C K バー = ' 0 ' の場合、このフィールドは無視される。

30

0：ガイド・フレーム処理後の肯定応答

1：ガイド・フレーム処理前の肯定応答

【 0 1 1 2 】

N E G / A L L バー：否定応答または全肯定応答制御値。このフィールドは、ガイド・コマンドが成功裡に完了しない場合を除き、A C K / N O A C K バー・フィールドが値 ' 0 ' b を有するとき無視される。

0：A C K / N O A C K バー・フィールドが ' 1 ' b の場合、全てのガイド・フレームを肯定応答する。早期または遅延肯定応答は、E A R L Y / L A T E バーの値により決定される。

40

1：成功裡に完了しないガイド・フレームだけを肯定応答する。この肯定応答は、A C K / N O A C K バー及び E A R L Y / L A T E バーの値に関係なく発生し、もちろん遅延肯定応答である。

【 0 1 1 3 】

U P / D O W N バー：アップまたはダウン制御値。この値は、フレームがアップサイドまたはダウンサイドのどちらで処理されるかを制御するために使用される。このフィールドは、R E S P / R E Q バーが ' 1 ' b のとき無視される。全てのマルチキャスト・ガイド・フレームは、' 0 ' b の U P / D O W N バー値を有する。更に、G T H ハードウェア補助機構命令の使用を要求するガイド・コマンドは、' 0 ' b の U P / D O W N バー値を有する

50

。

- 0 : ダウンサイド処理
- 1 : アップサイド処理

【 0 1 1 4 】

G T H / G C H バー : 汎用ツリー・ハンドラまたはガイド・セル・ハンドラ制御値。この値は、ガイド・フレームを適切なピコプロセッサに転送するために使用される。

- 0 : G C H ピコプロセッサ
- 1 : G T H ピコプロセッサ

【 0 1 1 5 】

T B : ターゲット・ブレード値。R E S P / R E Q バーが ' 0 ' b のとき、このフィールドは P r i z m a により使用されるルーティング情報を含む。各ビット位置はターゲット・ブレードに対応する。この値が ' 0 0 0 0 ' h のとき、ガイド・フレームがこのブレードに当てはまるとみなされ、従って実行される。T B フィールドの 1 つ以上のビット位置の ' 1 ' b の値は、セルが対応するターゲット・ブレードに経路指定されることを示す。R E S P / R E Q バーが ' 1 ' b のとき、このフィールドは応答ブレードの M y _ T B 値を含む。

10

【 0 1 1 6 】

ガイド・フレームのワード 1 は、相関関係子の値を含む (図 8)。この値は C P ソフトウェアにより割当てられ、ガイド・フレーム応答をそれらの要求に相関付ける。相関関係子は、機能を割当てられた複数のビットを含む。

20

【 0 1 1 7 】

あらゆるガイド・コマンドは、コマンド制御情報フィールドで開始する。このコマンド制御は、G C H ピココードがガイド・フレームを処理するのを支援する情報を含む。この情報のフォーマットが図 9 に示される。

【 0 1 1 8 】

レンジス値 : この値は、制御情報に含まれる 3 2 ビット・ワードの総数 (コマンド・ワード 0)、アドレス情報 (コマンド・ワード 1)、及びガイド・フレームのオペランド (コマンド・ワード 2 +) 部分を含む。

【 0 1 1 9 】

完了コード値 : このフィールドは C P により初期化され、ガイド・コマンドを処理するとき、G C H ピココードにより変更される。G C H ピココードはこのフィールドを、コマンド・リスト内の先行ガイド・コマンドの完了ステータスとして使用する。全てのガイド・コマンド・リストは、終了区切り文字ガイド・コマンドで終了するので、最後のコマンドの完了ステータスが、終了区切り文字の完了コード・フィールド内に含まれる。

30

【 0 1 2 0 】

ガイド・コマンド (G C) タイプ値 (シンボル名) :

【表 1】

シンボル名	タイプ値	タイプ記述	
End_Delimiter	0000	ガイド・フレーム・シーケンスの終わりをマーク	
Build_TSE_Free_List	0001	フリー・リストを作成	
Software_Action	0010	ソフトウェア・アクションを実行	
Unsolicited	0011	EPCピココードにより開始されるフレーム	10
Block_Write	0100	データ・ブロックを連続アドレスに書込む	
Duplicate_Write	0101	重複データをレジスタまたはメモリに書込む	
Read register	0110	レジスタまたはメモリ・データの読出しの要求または応答	
	0111	予約済み	
Insert_Leaf	1000	リーフを探索ツリーに挿入	
Update_Leaf	1001	探索ツリーのリーフを更新	20
Read_Leaf	1010	リーフの読出しの要求及び応答	
Leaf	1011	予約済み	
Delete_Leaf	1100	探索ツリーのリーフを削除	
	1101-1111	予約済み	

【0 1 2 1】

ガイド・フレームに含まれるアドレス情報は、ネットワーク・プロセッサのアドレッシング機構内の要素を識別する。アドレス情報フィールドの汎用形式が、図 1 0 に示される。

30

【0 1 2 2】

インタフェース装置は 3 2 ビット・アドレッシング方式を採用する。このアドレッシング方式は、アドレス値を、インタフェース装置のあらゆるアクセス可能な構造に割当てて。これらの構造はプロセッサの内部に存在するか、プロセッサの制御に従い、インタフェースに接続される。これらの構造のあるものは、組み込みプロセッサ・コンプレックス (EPC) により、ウェブ・インタフェースと呼ばれる内部インタフェースを介してアクセスされる。残りの構造は、メモリ制御装置インタフェースを介してアクセスされる。全ての場

40

【0 1 2 3】

ネットワーク制御装置は、主要チップ・アイランドに細分化される。各アイランドは固有のアイランド ID 値を与えられる。この 5 ビット・アイランド ID 値は、そのチップ・アイランドにより制御される構造のアドレスの、最上位 5 ビットを形成する。エンコードされたアイランド ID 値とチップ・アイランド名との対応が、図 1 2 に示される。ウェブ・アドレスの第 2 の部分は、次の最上位 2 3 ビットを含む。このアドレス・フィールドは、構造アドレス部分と、エレメント・アドレス部分とに区分化される。各セグメントに対して使用されるビット数は、アイランド毎に異なる。一部のアイランドは数個の大きな構造だけを含むのに対して、他のアイランドは多くの小さな構造を含む。そうした理由から、これらのアドレス・セグメントには固定サイズが存在しない。構造アドレス部分は、

50

アイランド内のアレイをアドレス指定するために使用されるのに対して、エレメント・アドレス部分は、アレイ内の要素をアドレス指定するために使用される。アドレスの残りの部分は、ウェブ・インタフェースの32ビット・データ・バス制限を調整する。この4ビット・ワード・アドレスは、アドレス指定されるエレメントの32ビット・セグメントを選択するために使用される。これは32ビットより広い構造エレメントを、ネットワーク制御装置のウェブ・データ・バスを介して、移動するために必要である。ワード・アドレス値'0'hは、構造エレメントの最上位32ビットを指し示すのに対して、順次ワード・アドレス値は、構造エレメントの連続的な下位セグメントに対応する。アドレスのワード・アドレス部分は、ウェブ・インタフェースを介してアクセスされない構造に対しては、要求されない。この理由から、アップ・データ・ストア、制御メモリ、及びダウン・データ・ストアは、アドレスの最下位27ビット全てを使用し、構造エレメントをアクセスする。このフォーマットの別の例外は、SPMインタフェースのアドレスである。この場合、アドレスの27ビット全てが使用され、32ビットよりも大きな幅のエレメントは存在しない。

10

【0124】

組み込みプロセッサ・コンプレックス（EPC）は、インタフェース素子チップのプログラマビリティを提供し制御する。これは次のコンポーネントを含む（図13参照）。

【0125】

N個の処理ユニット（G×Hと呼ばれる）：G×Hは、共通命令メモリに記憶されるピココードを同時に実行する。各G×Hは処理ユニット・コア（CLPと呼ばれる）を含み、これは3ステージ・パイプライン、16GPR、及びALUを含む。各G×Hはまた、例えばツリー構造探索エンジンのように、幾つかのコプロセッサを含む。

20

【0126】

命令メモリ：初期化の間にロードされ、フレームを転送し、システムを管理するピココードを含む。

【0127】

ディスパッチャ：アップ及びダウン・ディスパッチャ・キューからフレーム・アドレスをデキューする。デキューの後、ディスパッチャがアップまたはダウン・データ・ストア（DS）から、フレーム・ヘッダの一部をプリフェッチし、これを内部メモリに記憶する。G×Hがアイドルになると、ディスパッチャが直ちに、コード命令アドレス（CIA）のような適切な制御情報を有するフレーム・ヘッダを、G×Hに渡す。ディスパッチャはまた、タイマ及び割込みを処理する。

30

【0128】

ツリー構造探索メモリ（TSM）アービタ：各G×Hが使用可能な多数の共用内部メモリ位置及び外部メモリ位置が存在する。このメモリは共用されるので、メモリへのアクセスを制御するために、アービタが使用される。TSMはピココードにより直接アクセスされ、例えばエージング・テーブルを記憶するために使用される。TSMはまた、ツリー構造探索の間に、TSEによりアクセスされる。

【0129】

完了ユニット（CU）：完了ユニットは2つの機能を実行する。第1に、これはN個の処理ユニットを、アップ及びダウンEDS（エンキュー、デキュー、及びスケジュール・アイランド）とインタフェースする。EDSはエンキュー・アクションを実行する。すなわち、フレーム・アドレスが、FCBページと呼ばれる適切なパラメータと一緒に、伝送キュー、廃棄キュー、またはディスパッチャ・キューのいずれかに、待ち行列化される。第2に、完了ユニットはフレーム・シーケンスを保証する。複数のG×Hが、同一のフローに属するフレームを処理している可能性があるため、これらのフレームがアップまたはダウン伝送キューに正しい順序でエンキューされるように、事前注意が払われなければならない。完了ユニットは、フレーム・ディスパッチ時に分類子ハードウェア補助機構により生成されるラベルを使用する。

40

【0130】

50

分類子ハードウェア補助機構：アップ・フレームでは、分類子ハードウェア補助機構は、フレーム・フォーマットの周知のケースを分類する。分類結果はフレーム・ディスパッチの間に、C I A 及び 1 つ以上のレジスタの内容に関して、G x H に渡される。ダウン・フレームでは、分類子ハードウェア補助機構は、フレーム・ヘッダに応じて、C I A を決定する。アップ及びダウン・フレーム・ディスパッチの両方のために、分類子ハードウェア補助機構は、フレーム・シーケンスを維持するために、完了ユニットにより使用されるラベルを生成する。

【 0 1 3 1 】

アップ/ダウン・データ・ストア・インタフェース及びアービタ：各 G x H はアップ及びダウン・データ・ストアをアクセスできる。すなわち、"より多くの F I S H" を読出すとき、読出しアクセスが提供され、F I S H プールの内容をデータ・ストアに書戻すとき、書込みアクセスが提供される。N 個の処理ユニットが存在し、1 度にそれらの 1 つだけがアップ・データ・ストアをアクセスでき、1 度に 1 つだけがダウン・データ・ストアをアクセスできるので、各データ・ストアに対して、1 つのアービタが必要とされる。

10

【 0 1 3 2 】

ウェブ・アービタ及びウェブウォッチ・インタフェース：ウェブ・アービタは G x H の間で、ウェブへのアクセスを調停する。全ての G x H はウェブをアクセスでき、このことは、インタフェース装置内の全てのメモリ及びレジスタ機能をアクセスすることを可能にする。これにより、G x H は全ての構成領域を変更または読出すことができる。ウェブはインタフェース装置のメモリ・マップとみなされる。ウェブウォッチ・インタフェースは、3 つのチップ入出力ピンを用いて、チップの外部からウェブ全体をアクセスできるようにする。

20

【 0 1 3 3 】

デバッグ、割込み及び単一ステップ制御：ウェブは G C H またはウェブウォッチャが、必要に応じて、チップ上の各 G x H を制御することを可能にする。例えば、ウェブは G C H またはウェブウォッチャにより、G x H 上で命令を単一ステップ操作するために使用される。

【 0 1 3 4 】

P o w e r P C などの組み込み汎用プロセッサ：

4 つのタイプの G x H が存在する（図 1 4 参照）。

30

【 0 1 3 5 】

G D H（汎用データ・ハンドラ）：8 個の G D H が存在する。各 G D H は完全な C L P と、5 個のコプロセッサ（次のセクションで述べられる）を有する。G D H は主に、フレームを転送するために使用される。

【 0 1 3 6 】

G C H（ガイド・セル・ハンドラ）：G C H は、G D H と正に同一のハードウェアを有する。しかしながら、ガイド・フレームだけが G C H により処理される。G C H がデータ・フレームについても処理できるか否かは（この場合、G C H は G D H の役割をする）、ウェブ上でプログラマブルである（C L P _ E n a レジスタ）。G C H は G D H に比較して、追加のハードウェア、すなわち、ツリー挿入及び削除を実行するハードウェア補助機構を有する。G C H は、ガイド・セル関連ピココードを実行したり、エージングなどの、チップ及びツリー管理関連ピココードを実行したり、或いは制御情報を C P や別の G C H と交換するために使用される。実行すべきこうしたタスクが存在しない場合、G C H はフレーム転送関連ピココードを実行し、この場合、丁度 G D H のように動作することになる。

40

【 0 1 3 7 】

G T H（汎用ツリー・ハンドラ）：G T H はツリー挿入、ツリー削除、及びローブ管理を実行する追加のハードウェア補助機構を有する。G P Q 内に（ツリー管理コマンドを含む）フレームが存在しない場合、G T H はデータ・フレームを処理する。

【 0 1 3 8 】

50

G P H (汎用 P o w e r P C ハンドラ) : G P H は G D H 及び G T H に比較して、追加のハードウェアを有する。G P H はメールボックス・インタフェース (I / F) を介して、汎用プロセッサとインタフェースする。

【 0 1 3 9 】

G x H の数 (1 0 個) は " 最善推量 " (best-guess) である。性能評価により、実際に要求される G x H の個数が決定される。アーキテクチャ及び構造は、より多くの G x H に向けて完全にスケーラブルであり、唯一の制限はシリコン面積である (より大きなアービタ及び命令メモリを含むはずである) 。

【 0 1 4 0 】

各 G x H は、図 1 5 に示されるように構造化される。汎用レジスタ (G R P) 及び演算論理ユニット (A L U) を有する C L P に加え、各 G x H は次の 5 つのコプロセッサを含む。

【 0 1 4 1 】

(D S) コプロセッサ・インタフェース : ディスパッチャ、並びにアップ及びダウン・データ・ストアへの読出し及び書込みアクセスを提供するサブアイランドとインタフェースする。D S インタフェースは、いわゆる F I S H プールを含む。

【 0 1 4 2 】

ツリー構造探索エンジン・コプロセッサ (T S E) : T S E はツリー内の探索を実行し、ツリー構造探索メモリ (T S M) とインタフェースする。

【 0 1 4 3 】

エンキュー・コプロセッサ : 完了ユニット・インタフェースとインタフェースし、F C B ページを含む。このコプロセッサは、ピココードがエンキュー・パラメータを含む F C B ページを作成するために使用する追加のハードウェア補助機構と共に、2 5 6 ビット・レジスタを含む。一旦 F C B ページが作成されると、ピコプロセッサがエンキュー命令を実行し、F C B ページが完了ユニットに転送される。

【 0 1 4 4 】

ウェブ・インタフェース・コプロセッサ : このコプロセッサはウェブ・アービタとのインタフェースを提供し、インタフェース装置への書込み及び読出しを可能にする。

【 0 1 4 5 】

チェックサム・コプロセッサ : F I S H プール (後述) に記憶されるフレーム上に、チェックサムを生成する。

【 0 1 4 6 】

処理ユニットは、インGRES 処理とイーGRES 処理との間で共用される。インGRES 処理対イーGRES 処理において、どれだけの帯域幅が確保されるかは、プログラマブルである。現インプリメンテーションでは、2 つのモデルが存在し、一方は 5 0 / 5 0 (すなわち、インGRES 及びイーGRES が同じ帯域幅を獲得する) で、他は 6 6 対 3 4 (すなわち、インGRES がイーGRES の 2 倍の帯域幅を獲得する) である。

【 0 1 4 7 】

処理ユニットの動作はイベント・ドリブン方式である。すなわち、フレームの到来がイベントとして、またタイマまたは割込みのポッピングとして扱われる。ディスパッチャは異なるイベントを同様に扱う。但し、優先順位は存在する (第 1 が割込みで、第 2 がタイマ・イベントで、第 3 がフレーム到来イベント) 。イベントが処理ユニットに渡されると、適切な情報が処理ユニットに与えられる。フレーム到来イベントでは、これはフレーム・ヘッダの一部と、ハードウェア分類子からの情報とを含む。タイマ及び割込みでは、これはコード・エントリ・ポイントと、イベントに関連する他の情報とを含む。

【 0 1 4 8 】

フレームがインGRES 側に到来し、このフレームの受信バイト数がプログラマブルしきい値を超える場合、フレーム制御ブロックのアドレスが G Q に書込まれる。

【 0 1 4 9 】

完全なフレームがイーGRES 側で再組み立てされた場合、フレーム・アドレスが G Q に

10

20

30

40

50

書込まれる。4つのタイプのGQが存在する（そして、各タイプに対して、図14に示されるように、イングレス・バージョン及びイーグレス・バージョンが存在する）。すなわち、

G C Q : G C Hにより処理されなければならないフレームを含む。

G T Q : G T Hにより処理されなければならないフレームを含む。

G P Q : G P Hにより処理されなければならないフレームを含む。

G D Q : 任意のG D H（またはG C H / G T Hがデータ・フレームを処理できる場合には、G C H / G T H）により処理され得るフレームを含む。G D Qについては、複数の優先順位が存在し、高い優先順位でエンキューされたフレームは、低い優先順位でエンキューされたフレームより先に処理される。

10

【0150】

一部の処理ユニットは特殊化され得る。現インプリメンテーションでは、4つのタイプの処理ユニット（G x H）が存在する（図14参照）。すなわち、

G D H（汎用データ・ハンドラ）：G D Hは主に、フレームを転送するために使用される。

【0151】

G C H（ガイド・セル・ハンドラ）：G C Hは、G D Hと正に同一のハードウェアを有する。しかしながら、ガイド・フレームだけがG C Hにより処理される。G C Hがデータ・フレームについても処理できるか否かは（この場合、G C HはG D Hの役割をする）、ウェブ上でプログラマブルである（C L P _ E n aレジスタ）。

20

【0152】

G T H（汎用ツリー・ハンドラ）：G T HはG D H及びG C Hに比較して、追加のハードウェア、すなわち、ツリー挿入、ツリー削除、及びロープ管理を実行するハードウェア補助機構を有する。G P Q内に（ツリー管理コマンドを含む）フレームが存在しない場合、G T Hはデータ・フレームを処理する。

【0153】

G P H（汎用PowerPCハンドラ）：G P HはG D H及びG T Hに比較して、追加のハードウェアを有する。G P Hはメールボックス・インタフェースを介して、組み込みPowerPCとインタフェースする。

【0154】

30

実際のインプリメンテーションでは、G C H、G T H及びG P Hの役割は、単一の処理ユニット上で実現される。例えば、1インプリメンテーションはG C H及びG P Hに対して、1つの処理ユニットを有する。類似のコメントがG C Q、G T O及びG P Qについても適用できる。

【0155】

データ・ストア・コプロセッサの目的は、メディアから受信されたフレームを含むアップ・データ・ストアと、及びP r i z m a A t l a n t i cから受信される再組み立て済みフレームを含むダウン・データ・ストアとインタフェースすることである。

【0156】

データ・ストア・コプロセッサはまた、タイマ・イベントまたは割込みのディスパッチの間に、構成情報を受信する。

40

【0157】

データ・ストア・コプロセッサは、フレーム上のチェックサムを計算できる。

【0158】

データ・ストア・コプロセッサはF I S Hプール（8 F I S Hを保持できる）と、スクラッチ・メモリ（8 F I S Hを保持できる）と、アップまたはダウン・データベースからF I S Hプール内容を読み書きする幾つかの制御レジスタとを含む。F I S Hプールは、データ・ストアのためのある種の作業領域とみなされる。すなわち、データ・ストアを直接読み書きする代わりに、大量のフレーム・データがデータ・ストアからF I S Hプールに読出されるか、大量のデータがF I S Hプールからデータ・ストアに書込まれる。転送

50

の単位は F I S H すなわち 1 6 バイトである。

【 0 1 5 9 】

F I S H プールは、8 F I S H、すなわち各々が 1 2 8 ビットの 8 ワード分を含むことができるメモリとみなされる。C L P プロセッサ・アーキテクチャでは、F I S H プールは 1 2 8 バイトのレジスタ・アレイである。F I S H プール内の各バイトは、7 ビット・バイト・アドレス (0 乃至 1 2 7) を有し、アクセスは 1 6 ビットまたは 3 2 ビット・ベースで行われる。全てのレジスタ・アレイ同様、F I S H プールは循環アドレス方式を有する。すなわち、F I S H プール内のロケーション 1 2 6 で開始するワード (すなわち 4 バイト) のアドレス指定は、バイト 1 2 6、1 2 7、0 及び 1 を返却する。更に、データ・ストア・コプロセッサの観点から、F I S H プール内の F I S H ロケーションは、3 ビットの F I S H アドレスを有する。

10

【 0 1 6 0 】

フレーム・ディスパッチに際して、フレームの最初の N 個の F I S H が、ディスパッチャにより自動的に F I S H プールにコピーされる。N の値は、ポート構成メモリ (PortConfigMemory) 内でプログラマブルである。一般に、アップ・フレーム・ディスパッチでは、N は 4 に等しく、ダウン・ユニキャスト・フレーム・ディスパッチでは 2 で、ダウン・マルチキャスト・フレーム・ディスパッチでは 4 で、割込み及びタイマでは 0 である。

【 0 1 6 1 】

ピココードはフレームからより多くのバイトを読出すことができ、この場合、データ・ストア・コプロセッサが自動的に、フレーム・データを F I S H プール内の次の F I S H アドレスに読出し、F I S H プールの境界に達すると、自動的に 0 に循環する。また、ピココードはアップ / ダウン・データ・ストアを絶対アドレスにおいて読み書きすることもできる。

20

【 0 1 6 2 】

ウェブ・コプロセッサは E P C ウェブ・アービタとインタフェースする。E P C ウェブ・アービタは 1 0 個の G x H とウェブウォッチャとの間で、インタフェース装置のウェブ・インタフェース上でマスタになるものを調停する。これは全ての G x H がウェブ上で読み書きすることを可能にする。

【 0 1 6 3 】

インタフェース装置のメモリ・コンプレックスは、図 1 3 に示される組み込みプロセッサ・コンプレックス (E P C) の記憶機構を提供する。メモリ・コンプレックスはツリー構造探索メモリ (T S M) アービタと、複数のオンチップ及びオフチップ・メモリを含む。メモリはツリー構造、カウンタ、及びピココードによりメモリ・アクセスを要求される他のものを記憶する。更に、メモリはフリー・リストやキュー制御ブロックなど、ハードウェアにより使用されるデータ構造を記憶するために使用される。ツリーのために割当てられない、またはハードウェアにより使用されないメモリ・ロケーションは、デフォルトでは、ピココードがカウンタやエージング・テーブルとして使用することができる。

30

【 0 1 6 4 】

図 1 6 は、メモリ・コンプレックスの詳細ブロック図を示す。ツリー構造探索メモリ (T S M) アービタが、組み込みプロセッサ (G x H) とメモリとの間の通信リンクを提供する。メモリは 5 個のオンチップ S R A M と、1 個のオフチップ S R A M と、7 個のオフチップ D R A M とを含む。T S M アービタは、1 0 個の要求制御ユニット (各々が組み込みプロセッサ G x H の 1 つに接続される) と、1 3 個のメモリ・アービタ・ユニット (各々が各メモリに対応する) とを含む。各制御ユニット及びその接続される G x H が、全てのメモリをアクセスできるように、バス構造が要求制御ユニット及びアービタ・ユニットを相互接続する。

40

【 0 1 6 5 】

制御ユニットは、組み込みプロセッサ (G x H) とアービタとの間で、データを方向付けするために必要なハードウェアを含む。

【 0 1 6 6 】

50

S R A Mアービタ・ユニットはとりわけ、組み込みプロセッサ G x Hとオンチップ及びオフチップ S R A Mとの間の、データのフローを管理する。

【 0 1 6 7 】

D R A Mアービタ・ユニットはとりわけ、組み込みプロセッサ G x Hとオフチップ D R A M素子との間のデータのフローを管理する。

【 0 1 6 8 】

各メモリ・アービタは"バックドア" (back-door) アクセスを含み、これは一般にチップの他の部分により使用され、最も高いアクセス優先順位を有する。

【 0 1 6 9 】

D R A Mメモリは2つの動作モードで実行できる。すなわち、

10

1) T D Mモード: D R A M内の4つのバンクへのメモリ・アクセスが、読出しウィンドウ及び書込みウィンドウを交互することにより実行される。読出しウィンドウでは、4つのバンクのいずれかへのアクセスが読出し専用であり、書込みウィンドウでは、4つのバンクのいずれかへのアクセスが書込み専用である。複数のD R A Mに対してT D Mモードを使用することにより、D R A M間で幾つかの制御信号を共用することが可能になり、希少資源である幾つかのチップ入出力を節約できる。

【 0 1 7 0 】

2) 非T D Mモード: D R A M内の4つのバンクへのメモリ・アクセスが、読出しと書込みの組み合わせとなり得、これは特定の規則に従わねばならない。例えば、アクセス・ウィンドウ内で、バンクAでは読出しを、バンクCでは書込みを実行できる。

20

【 0 1 7 1 】

T S Mアービタは、N個のリクエストが同時にM個のメモリをアクセスすることを可能にする。複数のリクエストが同一のメモリをアクセスしたい場合、ラウンドロビン・アービトレーションが実行される。

【 0 1 7 2 】

M個のメモリは異なるプロパティを有することができる。現インプリメンテーションでは、3つのメモリ・タイプ、すなわち、内部S R A M、外部S R A M、及び外部D R A Mが存在する。

【 0 1 7 3 】

M個のメモリ及びNリクエストは同種 (homogeneous) であり、任意のリクエストが任意のメモリをアクセスできる。

30

【 0 1 7 4 】

一部のメモリは複数のサブメモリ (D R A M内の4つのバンクなど) に論理的に分割され、これらは論理的に同時にアクセスできる。

【 0 1 7 5 】

M個のメモリの一部は、内部的に使用されるデータ構造を含む制御メモリとして使用され、これらはピコプロセッサと比較して、高い優先アクセスを有する。このことはまたチップのでバッグを可能にする。なぜなら、ピコプロセッサは制御メモリの内容を読出すことができるからである。

【 0 1 7 6 】

40

アービタは読出しアクセス、書込みアクセス、及びread-add-writeをサポートする。それにより、Nビット整数がアトミック演算において、メモリの内容に追加される。

【 0 1 7 7 】

また、メモリ内のオブジェクトの物理ロケーションが透過的となるように、汎用アドレス方式が、M個のメモリをアクセスするために使用される。

【 0 1 7 8 】

ツリーの概念は、ツリー構造探索エンジンにより、情報を記憶及び検索するために使用される。検索、すなわちツリー探索及び挿入並びに削除は、キーにもとづき実行される。ここでキーは、例えばM A Cソース・アドレスのようなビット・パターンであるか、I Pソース・アドレスとI P宛先アドレスの連結である。少なくともキーを含む情報が、リー

50

フと呼ばれる制御ブロック内に記憶される（後述のように、記憶ビット・パターンは実際にはハッシュ・キーである）。リーフは更に、エージング情報などの追加の情報や、ターゲット・ブレード及びターゲット・ポート番号などの情報を転送するユーザ情報を含み得る。

【0179】

3つのタイプ（FM、LPM、SMT）（すなわち、固定マッチ、最長プレフィックス・マッチ及びソフトウェア管理ツリー）、及び関連ツリー・タイプ探索が存在する。ツリー探索の間にリーフをチェックするオプションの追加の基準は、ベクトルマスクである。ローピング、エージング及びラッチは、探索性能を向上させるために使用される。

【0180】

FMツリーの探索アルゴリズムが、図17に示される。探索アルゴリズムはキーを含む入力パラメータに作用し、キーにハッシュを実行し、ダイレクト・テーブル（DT）をアクセスし、パターン探索制御ブロック（PSCB）を通じてツリーを探索し、リーフに行き着く（図17）。3つのタイプのツリーが存在し、各々は異なる規則に従いツリー探索を発生させる独自の探索アルゴリズムを有する。例えば、固定マッチ（FM）ツリーでは、データ構造がパトリシア・ツリーである。リーフが見いだされるとき、このリーフは入力キーに一致する唯一の可能な候補である。ソフトウェア管理ツリーでは、リンク・リスト内でチェーニングされる複数のリーフが存在し得る。この場合、一致が見いだされるか、チェーンが尽きるまで、チェーン内の全てのリーフが入力キーと符合される。入力キーをリーフに記憶されるパターンと比較するいわゆる"最終比較"（compareat the end）操作が、リーフが真に入力キーに一致するか否かを確認する。リーフが見いだされ、一致が発生するとき、探索結果はOKであり、他の全ての場合にはKOである。

【0181】

探索操作への入力は、次のパラメータを含む。すなわち、

キー（128ビット）：キーは探索（または挿入／削除）の前に、特殊なピココード命令を用いて作成される。1つのキー・レジスタだけが存在する。しかしながら、ツリー構造探索が開始されると、キー・レジスタはTSEが探索を実行するのと並行して、ピココードにより、次の探索のキーを作成するために使用される。これはTSEがキーをバッシュ（bash）し、結果を内部ハッシュドキー・レジスタに記憶することによる（従って、実際には2つのキー・レジスタが存在する）。

【0182】

キー・レングス（7ビット）：このレジスタはキーの長さを含み、ビットで表す。これはキーの作成の間に、ハードウェアにより自動的に更新される。

【0183】

LU定義指標（LUDefindex）（8ビット）：これは探索が発生するツリーの完全な定義を含むLU定義テーブル（LUDefTable）への指標である。LU定義テーブルは以下で詳述される。

【0184】

TSRNR（1ビット）：探索結果はツリー探索結果領域0（TSR0）またはTSR1に記憶される。これはTSRNRにより指定される。TSEが探索を行っている間、ピココードは他のTSRをアクセスし、前の探索の結果を分析することができる。

【0185】

ベクトル指標（6ビット）：ベクトルマスクをイネーブルされるツリーでは（LU定義テーブル内で指定される）、ベクトル指標がベクトルマスク内のビットを示す。探索の終わりに、このビットの値が返却され、ピココードにより使用される。

【0186】

図17に示されるように、入力キーはハッシュドキーにハッシュされる。使用可能な6つの固定ハッシュ・アルゴリズムが存在する（1つの"アルゴリズム"はハッシュ関数を実行しない）。どのアルゴリズムが使用されるかについては、LU定義テーブル内で指定される。プログラマブル・ハッシュ関数が、柔軟性を追加するために使用され得る。

10

20

30

40

50

【 0 1 8 7 】

ハッシュ関数の出力は常に 1 2 8 ビット数であり、オリジナル入力キーとハッシュ関数の出力との間には、1 対 1 対応が存在する。後述のように、この特性はダイレクト・テーブルの後で開始するツリーの深さを最小化する。

【 0 1 8 8 】

図 1 7 の場合のように、カラーがツリーに対して許可される場合、1 6 ビット・カラー・レジスタが 1 2 8 ビット・ハッシュ関数出力に挿入される。挿入はダイレクト・テーブルの直後に発生する。すなわち、図示のように、ダイレクト・テーブルが 2^N 個のエントリを含む場合、1 6 ビット・カラー値がビット位置 N に挿入される。ハッシュ関数の出力は、挿入されるカラー値（イネーブルされる場合）と一緒に、ハッシュドキー・レジスタ

10

【 0 1 8 9 】

ハッシュ関数は、その出力内のほとんどのエントロピが最上位ビット側に存在するように定義される。ハッシュドキー・レジスタの最上位 N ビットは、ダイレクト・テーブル（ DT ）への指標を計算するために使用される。

【 0 1 9 0 】

探索はダイレクト・テーブルへのアクセスにより開始する。すなわち、ダイレクト・テーブル（ DT ）エントリがダイレクト・テーブルから読出される。 DT エントリを読出するために使用されるアドレスは、ハッシュドキーの最上位 N ビットから計算され、 LU 定義テーブル内で定義されるツリー属性に関しても同様である。これについては以下で詳述される。 DT エントリはツリーのルートとみなされる。使用される特定のツリー・データ構造は、ツリー・タイプに依存する。この時点では、パトリシア・ツリー・データ構造が FM ツリーとして、及び LP M 及び SM T ツリーのためのパトリシア・ツリーの拡張として使用されることを述べれば十分であろう。

20

【 0 1 9 1 】

8 個のエントリ・ダイレクト・テーブル（ DT ）の使用例が、図 1 8 に示される。 DT を使用することにより、探索時間（すなわち、アクセスされなければならない $PSCB$ の数）が低減される。従って、 DT サイズを増加することにより、メモリ使用と探索性能との間でトレードオフが生じる。

【 0 1 9 2 】

図 1 8 から明らかなように、 DT エントリは次の情報を含む。すなわち、

- 1) エンプティ：この DT エントリに接続されるリーフは存在しない。
- 2) リーフを指し示すポインタ：この DT エントリに接続される 1 つのリーフが存在する。
- 3) $PSCB$ を指し示すポインタ：この DT エントリに接続される 2 つ以上のリーフが存在する。 DT エントリはツリーのルートを定義する。

30

【 0 1 9 3 】

ソフトウェア管理ツリーの探索アルゴリズム、及びツリーを生成するアルゴリズムが、米国特許出願第 0 9 / 3 1 2 1 4 8 号で述べられている。

【 0 1 9 4 】

"選択ビット・アルゴリズム"と称されるアルゴリズムは、特定のメトリックを使用し、規則のセットまたは領域内の"ルール"と称されるアイテムから選択されるビットにもとづき、二分探索ツリーを作成する。ここで述べる全ての例は、インターネット・プロトコル（ IP ）ヘッダに関連して議論されるが、任意のタイプの固定フォーマット・ヘッダが代わりに使用され得る。

40

【 0 1 9 5 】

IP では、各ルールは次のサブセクション、すなわちソース・アドレス（ SA ）、宛先アドレス（ DA ）、ソース・ポート（ SP ）、宛先ポート（ DP ）、及びプロトコル（ P ）で作成される、特定のキーに関する。これらのデータはそれぞれ 3 2、3 2、1 6、1 6 及び 8 ビット長であり、従ってテストされるキーは 1 0 4 ビットから成る。選択ビット

50

・アルゴリズムは、104ビットの内の特に有用な幾つかのビットを見いだす。実際に2、3のビットをテストすることは、1つのルールを除く全てを、または2、3のルールを除く全てを、可能なアプリケーションから削除する。一部のルールでは、単純な比較操作による不等テストもまたふさわしい。ビット・テスト及び比較は、二分木において論理的に編成される。ツリーは、ビットを高速にテスト可能なハードウェア構造にマップされる。こうしたテストの結果、キーが適合する1つのルールまたは少数のルール（リーフ・チェーンと呼ばれる）が生成される。前者の場合、キーがルールにより詳細にテストされる。後者の場合、比較及び完全ルール・テストにより、キーがテストの枠内でテストされる。

【0196】

10

ルール・セット内の各ルールは、規則がキーに適合する最高優先順位の規則の場合に取られるアクションに関連付けられる。ルールは交差する（すなわち、1つのキーが2つまたはそれ以上のルールに適合する）。その場合、ルールは優先順位番号1、2、3、・・・を与えられ、任意の2つの交差するルールが異なる優先順位を有する（キーが2つ以上のルールに適合する場合、管理者はどのルールが上位になるかを宣言しなければならない）。従って、ビット・テスト及び比較後に、依然2つ以上のルールがテストされる場合、ルールは優先順位に従いテストされる。低い優先順位番号が、高い優先順位のルールを指定する。

【0197】

全く適合が見いだされない場合、デフォルト規定が指定される。

20

【0198】

最長プレフィックス・マッチング法の探索アルゴリズムが、米国特許第5787430号で述べられている。この方法は、前記データベースのノード（ルートノード）において入力するステップと、次の（子）ノードを識別するのに必要なエントリ部分だけを含む探索索引数のセグメントと、第2のリンク情報とを、前記セグメントが消費されるか、前記第2のリンク情報を欠く（リーフ）ノードに達するまで連続的に処理することにより、ツリー状データベースを通じてあるノードから別のノードへの探索パスを決定するステップと、前記探索索引数を、前記探索パスが終了するノードに記憶されるエントリと比較するステップと、前記現ノードにおいて、前記探索索引数と前記エントリとの間に、少なくとも部分的な一致が見いだされない場合、前記現ノードの第1のリンク情報を処理することにより、前記探索パスを後戻りするステップと、前記少なくとも部分的な一致が見いだされるか、または前記ルート・ノードに達するまで、前の2つのステップを繰り返すステップとを含む。

30

【0199】

図19は、メイン・スイッチング・ファブリック装置の実施例を示す。好適には、各インタフェース素子チップは、少なくとも2つの集積化並列直列変換ポートを有し、これらは並列データを受信して、それを高速直列データ・ストリームに変換する。そして、直列データがシリアル・リンクを介して、スイッチング・ファブリック装置に転送される。高速シリアル・リンクを介してスイッチング・ファブリック装置から受信されるデータは、別のDASLにより、並列データに変換される。データ・アライン・シリアル・リンク（DASL）と称されるシリアライザ/デシリアライザの実施例が、ここでは述べられる。

40

【0200】

少なくとも1つのDASLが、スイッチング・ファブリック装置をシリアル・リンクにインタフェースする。シリアル・リンクからのデータは、並列データに変換され、これがスイッチング・ファブリック装置に転送する。同様に、スイッチング・ファブリック装置からの並列データは直列データに変換され、シリアル・リンクに転送される。シリアル・リンクはスループットを向上するために集約される。

【0201】

図19を更に参照すると、スイッチング・システムはスイッチ・ファブリック11と、スイッチ・ファブリック入力ポート15（15-1，・・・，15-k）に接続される入力

50

スイッチ・アダプタ 13 (13 - 1 , . . . , 13 - k) と、スイッチ・ファブリック出力ポート 19 (19 - 1 , . . . , 19 - p) に接続される出力スイッチ・アダプタ 17 (17 - 1 , . . . , 17 - p) とを含む。

【 0 2 0 2 】

着信及び発信伝送リンク 21 (21 - 1 , . . . , 21 - g) 及び 23 (23 - 1 , . . . , 23 - r) は、それぞれ回線 (リンク) アダプタ 25 (25 - 1 , 25 - g) 及び 27 (27 - 1 , . . . , 27 - r) により、スイッチ・システムに接続される。伝送リンクは、ワークステーションや電話器などの接続ユニット (リンク指定 W S) に、或いはローカル・エリア・ネットワーク (リンク指定 L A N) に、更にサービス統合デジタル網 (I S D N) (リンク指定 I S D N) や他の通信システムに、回線交換またはパケット交換トラフィックを伝搬する。更に、プロセッサはスイッチ・アダプタ 13 及び 17 に直接接続される。回線アダプタ (L A) 及びスイッチ・アダプタ (S A) は、共通のインタフェースを有する。

10

【 0 2 0 3 】

入力スイッチ・アダプタでは、パケット交換及び回線交換インタフェースから様々なサービスが収集され、一様なミニパケット (幾つかの可能な固定長の 1 つを有する) に変換される。ルーティング情報を含むヘッダは、スイッチの要求出力ポート (及び出力リンク) を指定する。入力スイッチ・アダプタ内でのミニパケット・フォーマット及びミニパケット生成に関する詳細、及び出力スイッチ・アダプタ内でのパケット解除に関する詳細について、次に述べることにする。

20

【 0 2 0 4 】

スイッチ・ファブリックはミニパケットを、高速自己ルーティング相互接続ネットワークを介して、任意の入力ポートから任意の出力ポートに経路指定する。自己ルーティング・ネットワークの構造は、ミニパケットが競合無しに、内部的に同時に経路指定されることである。

【 0 2 0 5 】

スイッチング・システムの心臓部は、スイッチ・ファブリックである。2つの異なるインプリメンテーションが考慮され、これらについて次に述べる。1インプリメンテーションでは、スイッチ・ファブリックは、各入力ポートに対して、それぞれの入力ポートを全ての出力ポートに接続する自己ルーティング二分木を含み、(k 個の入力ポートが提供される場合、) k 個のこうした二分木が組み合わされて提供される。他のインプリメンテーションでは、出力 R A M を有するバス構造が、各出力構造に対してスライスとして提供され、全ての入力ポートをそれぞれの出力ポートに接続する。(p 個の出力ポートが提供される場合、) スwitch・ファブリックは、p 個のこうしたスライスを組み合わせて提供される。

30

【 0 2 0 6 】

D A S L については、1999 年 6 月 11 日出願の米国特許出願第 09 / 330968 号で述べられている。D A S L インタフェースは、C M O S A S I C などのパラレル・インタフェースからデータを受信し、パラレル・インタフェースからのビットを、少数の並列ビット・ストリームに区分化する。少数の並列ビット・ストリームが次に高速直列ストリームに変換され、これが伝送メディアを介して、他のモジュールの受信機に移送される。制御インピーダンスを有する差動ドライバが、データの直列ビット・ストリームを伝送メディアに駆動する。

40

【 0 2 0 7 】

D A S L は、N ビット並列データとして表されるデータ・ストリームを、各々が n ビット (但し n は N の分数) を有する複数の部分に構文解析し、データ・ストリームの各 n ビット部分を直列化し、直列化された各部分を複数の並列チャネルの対応するチャネルを介して転送し、データ・ストリームの各転送部分を非直列化し、データ・ストリームを N ビット並列データとして復元する。

【 0 2 0 8 】

50

以上、図面を参照しながら、本発明の好適な実施例について述べてきた。説明の中で使用された特定の用語は総称的な意味を成すもので、限定的な意味で使用されるものではない。

【図面の簡単な説明】

【0209】

【図1】本発明に従うインタフェース装置のブロック図である。

【図2】MACのブロック図である。

【図3】異なるシステム構成内の他のコンポーネントと相互接続されるインタフェース装置を示す図である。

【図4】カプセル化されたガイド・フレームのフロー及び処理を示す図である。

10

【図5】内部ガイド・フレームのフロー及び処理を示す図である。

【図6】ガイド・セルの汎用フォーマットを示す図である。

【図7】フレーム制御情報のフォーマットを示す図である。

【図8】相関関係子のフォーマットを示す図である。

【図9】コマンド制御情報のフォーマットを示す図である。

【図10】アドレス情報のフォーマットを示す図である。

【図11】構造アドレッシングの汎用形式を示す図である。

【図12】アドレッシング、アイランド・エンコードを示す表である。

【図13】組み込みプロセッサ・コンプレックスのブロック図である。

【図14】組み込みプロセッサの概略図である。

20

【図15】G×Hプロセッサの構造を示す図である。

【図16】メモリ・コンプレックスのブロック図である。

【図17】固定マッチ(FM)検索アルゴリズムのフローチャートである。

【図18】ダイレクト・テーブルを使用する場合と、使用しない場合のデータ構造を示すフローである。

【図19】Prizmaなどのスイッチング・システムのブロック図である。

【図20】CPのブロック図である。

【図21】EDS-UP、EDS-DOWN及びEPCにおける、シングルチップ・ネットワーク・プロセッサの強調表示機能のブロック図である。

【符号の説明】

30

【0210】

10、38 インタフェース装置

12 組み込みプロセッサ・コンプレックス(EPC)

14 多重化MACアップ(PPM-UP)、多重化MAC

16 エンキュー・デキュー・スケジューリング・アップ(EDS-UP)

18 スイッチ・データ・ムーバ・アップ(SDM-UP)

20、30 システム・インタフェース(SIF)

22 データ・アライン・シリアル・リンクA(DASLA)

24 データ・アライン・シリアル・リンクB(DASLB)

32 SDM-DN

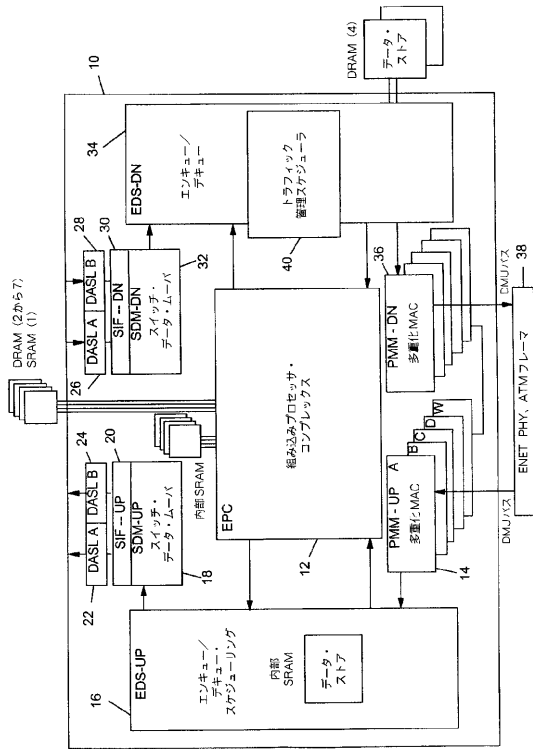
40

34 EDS-DN

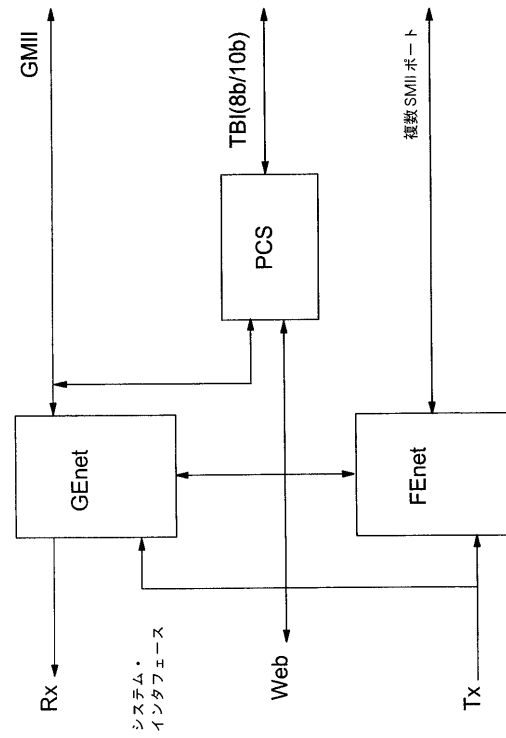
36 イーグレス・イーサネット(R)MAC

40 トラフィック管理スケジューラ

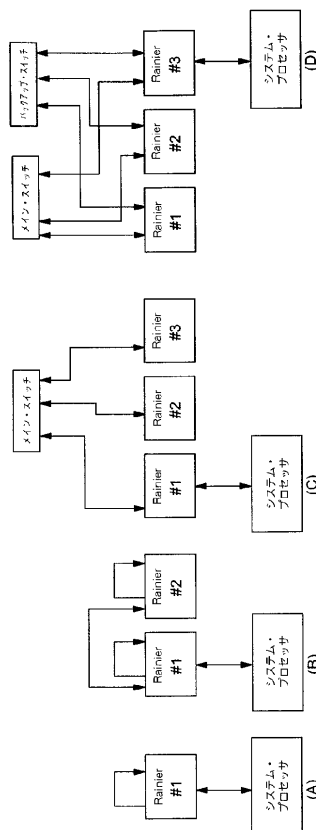
【 図 1 】



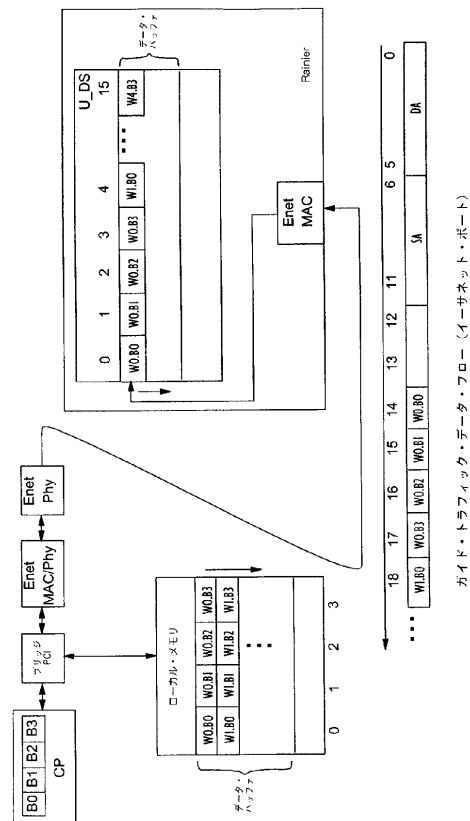
【 図 2 】



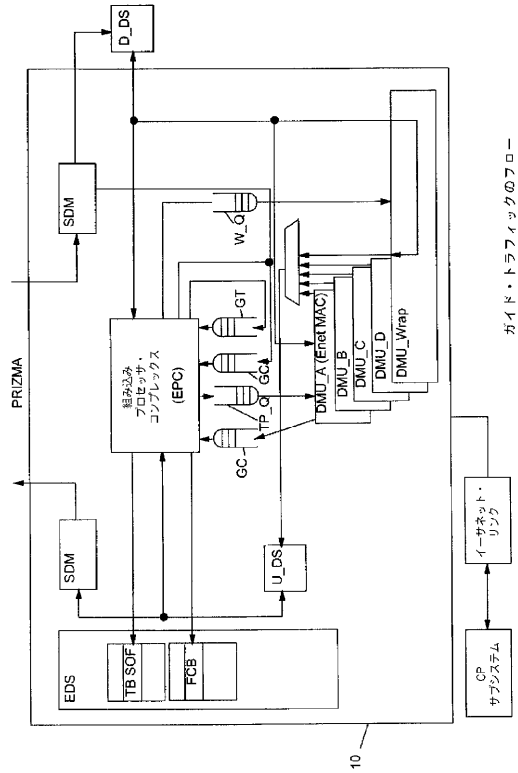
【 図 3 】



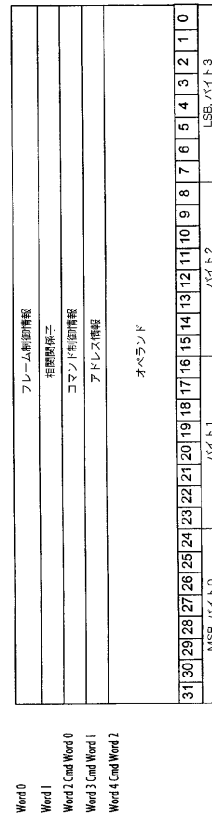
【 図 4 】



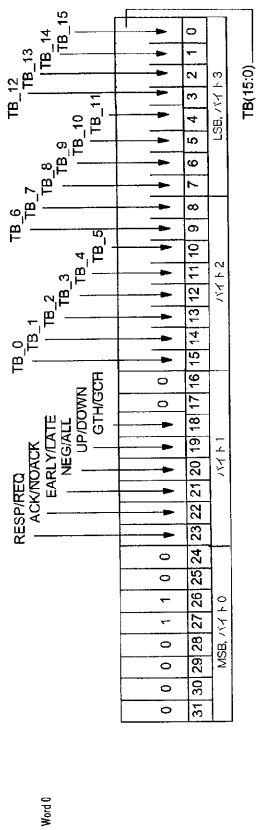
【 図 5 】



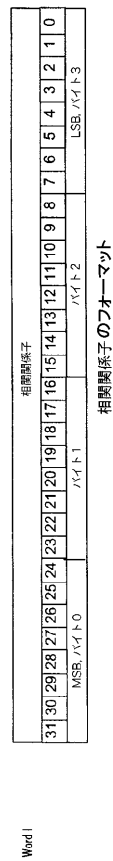
【 図 6 】



【 図 7 】



【 図 8 】



【図 1 1 1】

アイランドID	構造アドレス	エレメント・アドレス	ワード・アドレス
5	23		4
	32		

構造アドレスリングの汎用形式

【図 9】

L																変コード																res																GCタイプ															
1	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0																															
MSB, バイト 0																バイト 1																バイト 2																LSB, バイト 3															

コマンド制御情報のフォーマット

【図 1 1 2】

アイランドID	アイランド名
'00000'b	アップ・データ・ストア*
'00001'b	アップPMM
'00010'b	アップEDS
'00011'b	アップSDM
'00100'b	組み込みプロセッサ・コンプレックス
'00101'b	SPM
'00110-00111'b	予約済み
'01000'b	制御メモリ*
'01001-01111'b	予約済み
'10000'b	ダウン・データ・ストア*
'10001'b	ダウンPMM
'10010'b	ダウンEDS
'10011'b	ダウンSDM
'10100'b	構成レジスタ
'10101'b	DASL
'10110-11111'b	予約済み
*これらのアイランドはウェア・インタフェースを介してアクセスできない	

アドレスリング、アイランド・エンコーディング

【図 1 0】

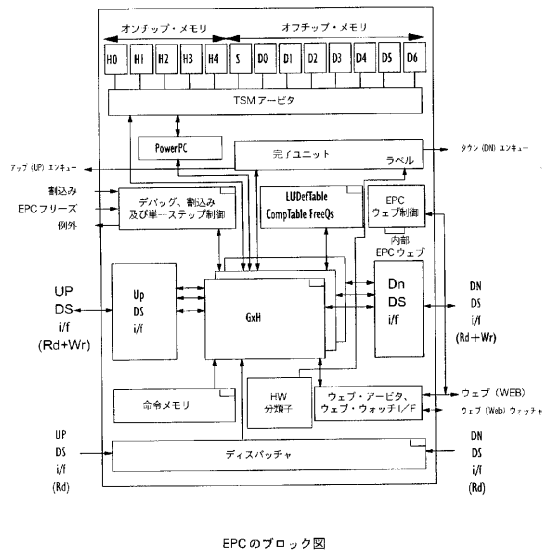
アドレス																																							
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0								
MSB, バイト 0																バイト 1								バイト 2								LSB, バイト 3							

アドレス情報のフォーマット

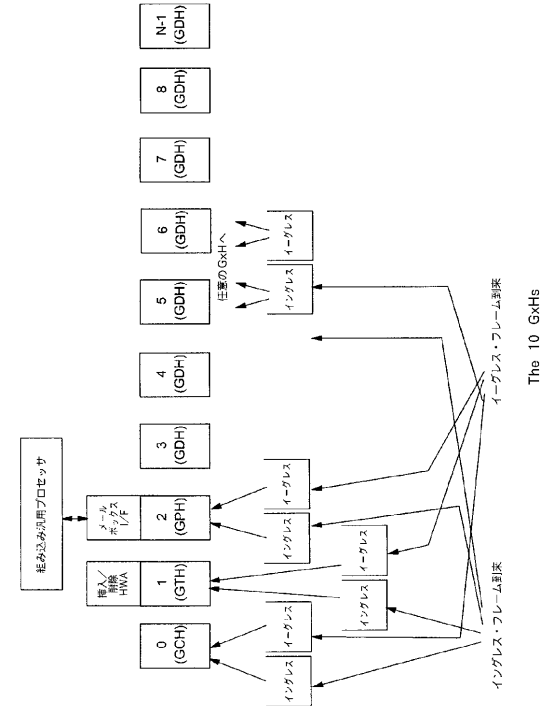
Word 1 End Word 1

Word 1 End Word 0

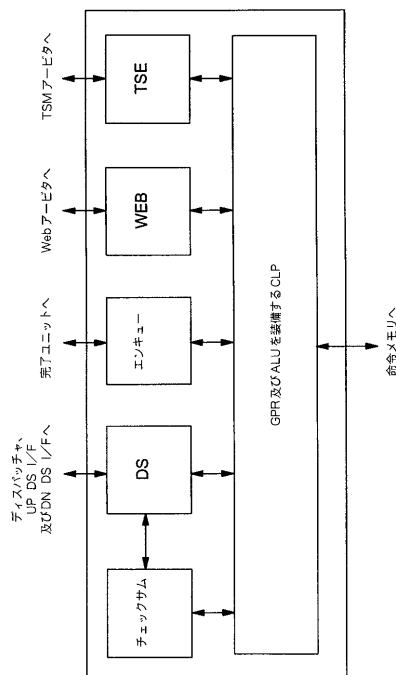
【図 13】



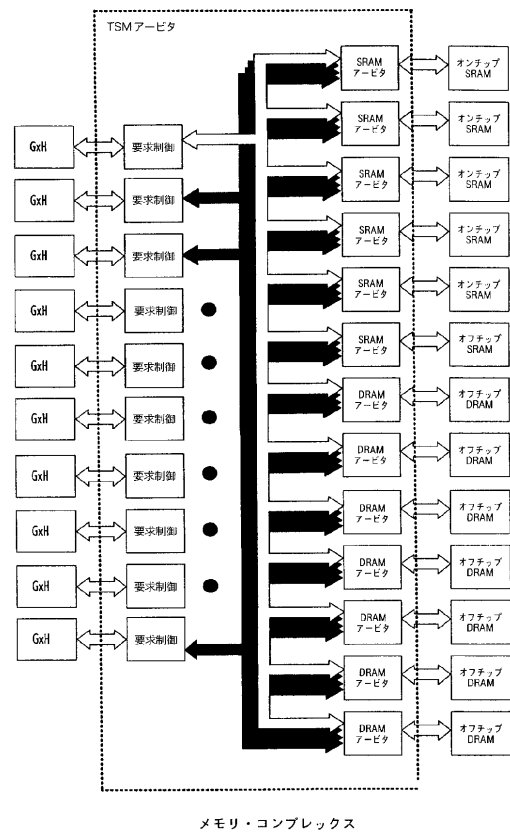
【図 14】



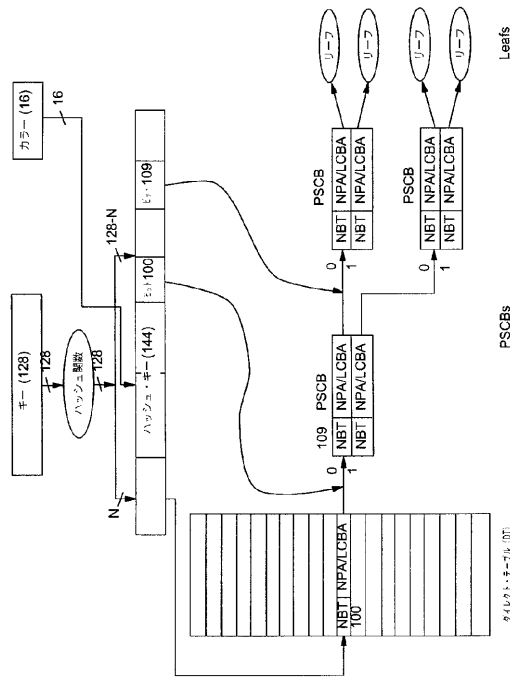
【図 15】



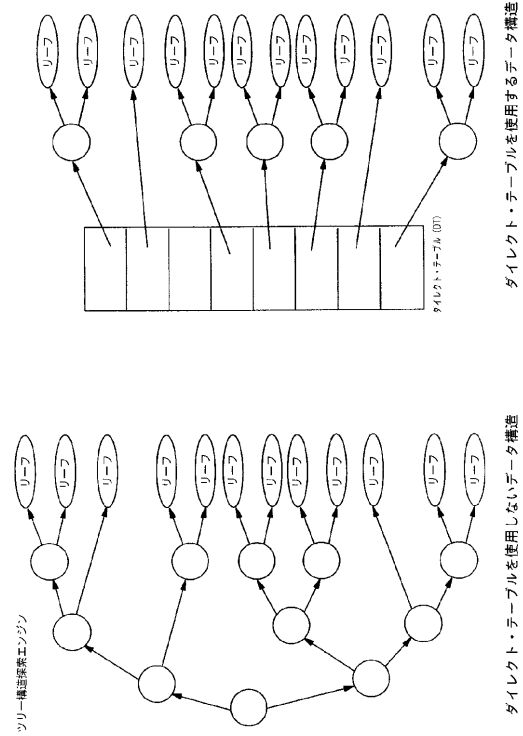
【図 16】



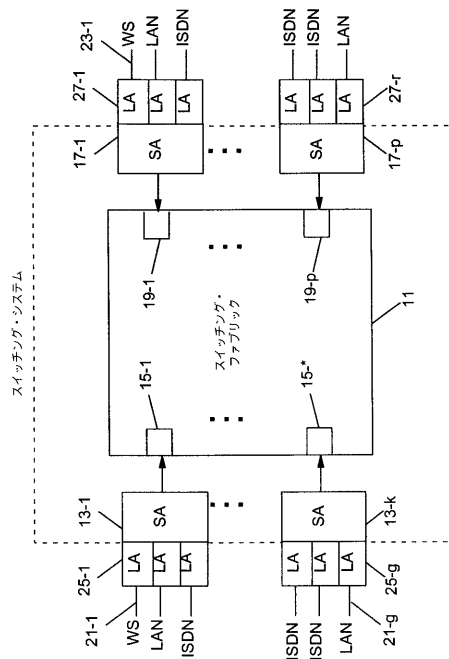
【 図 1 7 】



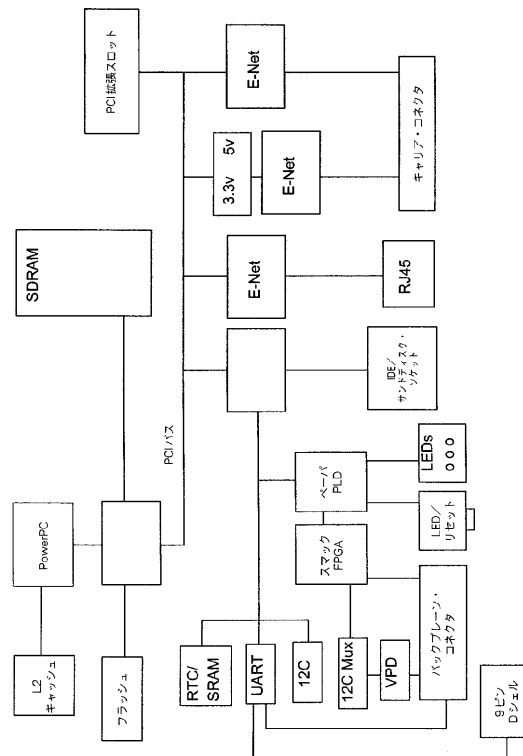
【 図 1 8 】



【 図 1 9 】



【 図 2 0 】



The diagram illustrates a system architecture with two main memory sections at the top: **DRAM (2から7)** and **SRAM (1)**. Below these, the system is divided into several functional blocks:

- Left Side (Input/Output and Control):** Includes **EDS-UP** (with TB リング, OCB, BCB, FCB, データ・ストア, and G-DCB), **EDS-DN** (with RCB and G-DCB), and a central control area with **SMT リープ・キャッチ**, **SMT-PCSB**, **キャッシュ・制御**, **アービタ**, **EPC**, and **ディバイス**.
- Right Side (Processing and Storage):** Includes **PMN-UP M-MACS**, **PMN-DN M-MACS**, **Power PC**, **GDH0**, **GDH1**, **GDH8**, **GDH9**, **トラフィック 管理スワッチャー**, **TD**, **QC8**, **PCB**, and **DS0**, **DS1**.
- Interconnections:** Various lines connect these blocks, including a **フラッシュ** (Flash) memory block and a **ネットワーク** (Network) interface at the bottom right.

フロントページの続き

- (72)発明者 バス、ブライアン、ミッチェル
アメリカ合衆国 2 7 5 0 2、ノース・カロライナ州アベックス、オールド・スターブリッジ・ドライブ 4 0 2 1
- (72)発明者 カルビンナック、ジーン、ルイス
アメリカ合衆国 2 7 5 1 1、ノース・カロライナ州カーリー、スプリング・ホロー・レーン 1 1 2
- (72)発明者 ガロー、アンソニー、マッテオ
アメリカ合衆国 2 7 5 0 2、ノース・カロライナ州アベックス、コーシャム・ドライブ 3 3 0 8
- (72)発明者 ヘデス、マルコ、シイ
アメリカ合衆国 2 7 6 1 2、ノース・カロライナ州ラーレー、グランド・メナー・コート 4 1 0 9、ナンバー 3 0 8
- (72)発明者 ラオ、スリダール
アメリカ合衆国 2 7 6 1 2、ノース・カロライナ州ラーレー、ビーバーブルック・ロード 5 0 2 0、アパートメント 2 0 4
- (72)発明者 シーゲル、マイケル、スティーブン
アメリカ合衆国 2 7 6 1 3、ノース・カロライナ州ラーレー、ロワリー・ドライブ 1 0 6 2 5
- (72)発明者 ヤングマン、ブライアン、アラン
アメリカ合衆国 2 7 5 1 1、ノース・カロライナ州カーリー、ニューポート・サークル 7 0 2
- (72)発明者 ヴァーブランケン、ファブリス、ジーン
フランス 0 6 6 1 0、ラ・ゴード、ルート・ド・カーニュ 9 1 5 2

審査官 清水 稔

(56)参考文献 特開平 1 1 - 0 8 8 3 4 5 (J P , A)

(58)調査した分野(Int.Cl. , D B 名)
H 0 4 L 1 2 / 5 6