

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5220974号
(P5220974)

(45) 発行日 平成25年6月26日 (2013. 6. 26)

(24) 登録日 平成25年3月15日 (2013. 3. 15)

(51) Int. Cl.

F I

G 0 6 F 13/00 (2006. 01)

G 0 6 F 13/00 3 5 4 A

請求項の数 1 (全 31 頁)

(21) 出願番号	特願2001-530282 (P2001-530282)	(73) 特許権者	505160898
(86) (22) 出願日	平成12年10月12日 (2000. 10. 12)		ブルアーク ユーケー リミテッド
(65) 公表番号	特表2003-511777 (P2003-511777A)		イギリス国 アールジー 1 2 1 アールビ
(43) 公表日	平成15年3月25日 (2003. 3. 25)		ー パークシャー, ブラックネル, ク
(86) 国際出願番号	PCT/EP2000/010277		ックハム ロード, クイーンズゲート
(87) 国際公開番号	W02001/028179		ハウス
(87) 国際公開日	平成13年4月19日 (2001. 4. 19)	(74) 代理人	110000279
審査請求日	平成19年10月12日 (2007. 10. 12)		特許業務法人ウィルフォート国際特許事務
審判番号	不服2011-22151 (P2011-22151/J1)		所
審判請求日	平成23年10月13日 (2011. 10. 13)	(72) 発明者	バレル、ジェフリー・エス
(31) 優先権主張番号	09/418, 558		英国、パークシャー・アールジー 4 6・4
(32) 優先日	平成11年10月14日 (1999. 10. 14)		エイチワイ、ロウアー・アーリー、ボウチ
(33) 優先権主張国	米国 (US)		ーフ・クロース 9

最終頁に続く

(54) 【発明の名称】 ハードウェア実行又はオペレーティングシステム機能の加速のための装置及び方法

(57) 【特許請求の範囲】

【請求項 1】

ネットワークノードからネットワークを介してネットワークサービス要求を受信し、前記ネットワークサービス要求を扱って、記憶装置に記憶装置アクセス要求を送信する要求処理装置であって、

ソフトウェアプログラムの制御下にあるシステムバスに接続されたプロセッサと

前記ネットワークサービス要求を受信するネットワークサブシステムと、

前記ネットワークサブシステムに接続され前記記憶装置アクセス要求を前記記憶装置に送信するサービスサブシステムと

を備え、

前記ネットワークサブシステム及びサービスサブシステムは、前記システムバスから独立した受信高速バスインターフェース及び送信高速バスインターフェースにより接続され、

前記ネットワークサブシステム及び前記サービスサブシステムは、それぞれ、前記ソフトウェアプログラムの直接の制御外で駆動する専用ハードウェアであり、

前記ネットワークサブシステムは、

前記システムバスと前記受信高速バスインターフェースとに接続されており、前記ネットワークサービス要求を前記サービスサブシステム或いは前記プロセッサに予定されているのかを決定し、前記サービスサブシステムに予定されているカプセルから取り出したネットワークサービス要求を、前記受信高速バスインターフェースを介して前記サービスサ

ブシステムに送信する受信モジュールと、

前記システムバスと前記送信高速バスインターフェースとに接続されており、前記ネットワークサービス要求のネットワークサービス応答を、前記サービスサブシステムから前記送信高速バスインターフェースを介して受信し、前記ネットワークサービス応答を前記ネットワークノードに送信する送信モジュールと

を備え、

前記サービスサブシステムは、

前記ネットワークサービス要求を受信し前記記憶装置アクセス要求を発するサービス/ファイルシステムモジュールと、

前記サービス/ファイルシステムモジュールから前記記憶装置アクセス要求を受信し、
前記記憶装置アクセス要求を記憶装置に送信する記憶装置モジュールと

10

を備え、

前記サービス/ファイルシステムモジュールが、

前記システムバスと前記受信高速バスインターフェースとに接続されており、前記受信高速バスインターフェースを介して前記ネットワークサービス要求を受信し、前記記憶装置アクセス要求を発する受信器と、

前記システムバスと前記送信高速バスインターフェースとに接続されており、前記記憶装置アクセス要求の記憶装置アクセス応答を前記記憶装置モジュールから受信し、前記記憶装置アクセス応答に基づくネットワークサービス応答を生成し、前記ネットワークサービス応答を、前記送信高速バスインターフェースを介して前記ネットワークサブシステムの
前記送信モジュールに送信する送信器と、

20

前記システムバス、前記受信高速バスインターフェース及び前記送信高速バスインターフェースから独立しており前記受信器と前記送信器とを結ぶ専用バスと、

前記システムバス、前記受信高速バスインターフェース及び前記送信高速バスインターフェースのいずれにも接続されておらず前記専用バスに接続されており前記受信器及び前記送信器の処理を制御する専用プロセッサと

を備え、

前記受信器が、

前記ネットワークサービス要求で要求されたデータの物理的位置の要求を出し、その要求の応答を受けて、第1の待ち行列に入れ、第1の待ち行列内の応答を基に、アクセス要求を発する受信制御エンジンと、

30

前記物理的位置の要求を受けて、その要求を第2の待ち行列に入れ、前記第2の待ち行列内の要求を処理することで、前記物理的位置を表す情報を発し、その要求の応答を前記受信制御エンジンに送る物理位置制御エンジンと、

前記受信制御エンジンからの前記アクセス要求と前記物理位置制御エンジンからの前記物理的位置を表す情報とを基に前記記憶装置アクセス要求を生成する生成エンジンと、

前記生成エンジンにより生成された前記記憶装置アクセス要求を発する要求送信インターフェースと

を備え、

前記記憶装置モジュールが、

40

前記記憶装置アクセス要求を受信し、前記記憶装置アクセス要求を前記記憶装置に適したフォーマットに変換する要求インターフェースと、

前記記憶装置からの応答を取得し、その応答を前記サービス/ファイルシステムモジュールに適したフォーマットに変換する承認インターフェースと、

データの部分への高速読込アクセスを可能にするために前記記憶装置に含まれるデータの一部のローカルコピーを保持するキャッシュ制御部と

を備え、

前記ネットワークノードから前記ネットワークサービス要求を受信したか否かの判断を含んだ第1のループと、記憶装置アクセス要求が完了したか否かを含んだ第2のループとで構成された単一スレッドが実行されるようになっており、

50

前記ネットワークサブシステムにおいて前記第１のループが実行され、前記サービスサブシステムにおいて前記第２のループが実行される、ことを特徴とする要求処理装置。

【発明の詳細な説明】

【０００１】

発明の技術分野

本発明は、オペレーティングシステム機能及びハードウェア実行、あるいは、そのような機能の加速に関する。

【０００２】

発明の背景技術

コンピュータのオペレーティングシステムは、コンピュータが外部供給源と通信することを可能にする。オペレーティングシステムは、典型的に、キーボード、ディスプレイ、ディスク記憶装置、ネットワーク設備、プリンタ、モデムなどを含むコンピュータ使用法に結び付けられたアイテムの直接制御を扱う。コンピュータのオペレーティングシステムは、典型的に、中央演算処理装置（ＣＰＵ）にローカル及びネットワークファイルシステム、メモリ、周辺装置ドライバ、及びアプリケーション処理を含む処理を管理することを含むタスクを実行させるように設計される。これらのすべての機能に対する責任をＣＰＵ上に置くことは、特に、オペレーティングシステムが例えば、Windows NT（登録商標）（ワシントン州レッドモンドのマイクロソフト社から入手可能）、Unix（登録商標）（カリフォルニア州サンタクルーズのSCO Softwareから、及びマサチューセッツ州ケンブリッジの Red Hat Softwareの「Linux」と呼ばれるバージョンを含む多くの供給源から利用可能）、及びNetWare（ユタ州プロボのNovellから利用可能）のように高性能であるとき、重大な処理負担を強要する。負担がアプリケーション結び付けられた以外の処理を実行するＣＰＵ上に置かれると、アプリケーションのパフォーマンスが低下し得る結果となり、アプリケーションを実行するのに利用可能なＣＰＵ時間が益々小さくなる。それに加えて、ＣＰＵの外部装置のスループットは、オペレーティングシステムがこれらの装置を管理する責任をＣＰＵ上に置くとき、ＣＰＵによって強要される制限を受けやすい。さらに、ＣＰＵを含み、オペレーティングシステムを実行し、装置に結び付けられるソフトウェア－ハードウェアシステム全体の信頼性は、とりわけ、オペレーティングシステムに依存する。オペレーティングシステムに内在する複雑さに起因して、ソフトウェア－ハードウェアシステム全体の安定性を害する不測の状態が発生し得る。

【０００３】

発明の概要

この概要で列挙される本発明のある態様は、これより同日に提出された他の出願の目的である。本発明の一態様では、ネットワークを介してサービス要求を扱うための装置が提供され、ネットワークはプロトコルを利用する。この態様では、該装置は、

a．ネットワークプロトコルを用いてネットワークサービス要求を送受信するためのネットワークサブシステムと、

b．前記ネットワークサービス要求を満たすための、前記ネットワークサブシステムに接続されるサービスサブシステムとを含む。

【０００４】

同じく、この態様では、ネットワークサブシステムとサービスサブシステムの少なくとも一つは、ハードウェアで実行され、ネットワークサブシステムとサービスサブシステムの他方は、任意的にハードウェアで加速されてもよい。その代わりに、又はそれに加えて、サービスサブシステムは、ハードウェアで加速されてもよい。

【０００５】

関連する実施の形態では、サービス要求は、長期の電子記憶装置へのデータの読み込み及び書き込みの一つを含み、追加的に、ネットワークサブシステムはハードウェアで加速される。同じく、追加的に、長期の記憶装置は、ネットワークを介してコンピュータにアクセス可能なネットワークディスク記憶装置である。その代わりに、長期の記憶装置は、口

10

20

30

40

50

ーカルコンピュータにアクセス可能であるが、ネットワークを介して他のあらゆるコンピュータにはアクセス不可能なローカルディスク記憶装置である。同じく追加的に、長期の記憶装置は、ネットワークを介して電子メールの供給に関連し、あるいは、それはネットワーク上のウェブページへのアクセスを供給してもよい。

【 0 0 0 6 】

同様に、サービス要求は、記憶装置システムのデータのアクセスを含んでもよく、サービスサブシステムもまた、記憶装置システムのデータの記憶を管理するためのハードウェアで実行されたモジュールを含んでもよい。したがって、一実施の形態では、そのような装置がファイルサーバーであって、記憶装置システムのデータがファイルに配列され、サービス要求は、記憶装置システムのファイルへの要求を含み、サービスサブシステムもまた、記憶装置システムに関連したファイルシステムを管理するためのハードウェアで実行されたモジュールを含む。

10

【 0 0 0 7 】

もう一つの関連する態様では、プロトコルは、ファイルシステムプロトコルを含み、ファイルシステムプロトコルは、ファイル読み込み及びファイル書込みを含む操作を定義する。装置がウェブサーバーであってもよく、記憶装置システムのデータがウェブページを含んでもよく、サービス要求は、記憶装置システム内のウェブページのための要求を含んでもよい。同様に、プロトコルがIPを含んでもよい。さらなる関連する態様では、記憶装置システムが記憶装置プロトコルを有し、サービスサブシステムは、記憶装置システムと相互接続するためのハードウェアで実行されたモジュールを含む。

20

【 0 0 0 8 】

もう一つの態様では、ネットワークを介してデータを送受信するためのサブシステムであって、ネットワークは層3及び4の少なくとも一つを有するプロトコルを用い、サブシステムは、

ネットワークからカプセル収納データを受信し、プロトコルに従ってそのようなデータをカプセルから出す受信器と、

プロトコルに従ってデータをカプセルに収納し、ネットワークを介してカプセル収納データを送信する送信器とを含む。

【 0 0 0 9 】

受信器及び送信器の少なくとも一つはハードウェアで実行され、その代わりに、又はそれに加えて、受信器及び送信器の少なくとも一つはハードウェアで加速される。さらなる実施の形態では、ネットワークは、TCP/IPプロトコルを用いる。関連する実施の形態では、データは、ネットワークを介してパケットで受信され、各パケットはプロトコルヘッダを有し、サブシステムは、また、受信器によって受信される各パケットのプロトコルヘッダ内に含まれる情報から独特な接続を決定する接続識別子を含む。もう一つの関連する実施の形態では、カプセル収納データは、ネットワーク接続に関連し、サブシステムは、接続の状態を格納する、ネットワーク接続に関連したメモリ領域をさらに含む。

30

【 0 0 1 0 】

もう一つの関連する態様では、記憶装置アクセス要求を生成し得るネットワークと記憶装置配置を相互接続するためのサービスサブシステムが提供される。この態様のサービスサブシステムは、

40

a. ネットワークサービス要求を受信し、そのようなサービス要求を満たし、そうして、データ記憶装置アクセス要求を発することができるサービスモジュールと、

b. 前記サービスモジュールからデータ記憶装置アクセス要求を受信し、そのような記憶装置アクセス要求を満たし、そうして、記憶装置配置アクセス要求を発することができる、該サービスモジュールに接続されたファイルシステムモジュールと、

c. 前記ファイルシステムモジュールから記憶装置配置アクセス要求を受信し、そのような記憶装置配置アクセス要求を満たすために、記憶装置配置を制御する、該ファイルシステムモジュールに接続された記憶装置モジュールとを含む。

【 0 0 1 1 】

50

モジュールの少なくとも一つはハードウェアで実行され、その代わりに、又はそれに加えて、モジュールの少なくとも一つは、ハードウェアで加速される。関連する実施の形態では、サービスモジュールは、

i. ネットワークサービス要求を受信し、そのような要求が適切な否か決定し、もし適切ならば、情報が利用可能か否か応答し、さもなければ、データ記憶装置アクセス要求を発する、受信制御エンジンと、

ii. 前記受信制御エンジンからの命令に基づいて、ネットワークサービス応答を生成し、前記データ記憶装置アクセス要求へのデータ記憶装置アクセス応答がある場合、該データ記憶装置アクセス応答を処理する、送信制御エンジンとを含む。

【0012】

10

エンジンの少なくとも一つはハードウェアで実行され、その代わりに、又はそれに加えて、エンジンの少なくとも一つはハードウェアで加速される。他の関連する実施の形態では、サービスサブシステムは、コンピュータのマザーボードに直接統合され、又は、コンピュータに接続され得るアダプタカードに統合される。

【0013】

もう一つの態様では、ネットワークサービス要求を受信し、そのようなサービス要求を満たすサービスモジュールが提供される。そのサービスモジュールは、

a. ネットワークサービスを受信し、そのような要求が適切であるか否かを決定し、もし適切ならば、情報が利用可能であるか否かを応答し、さもなければ、データ記憶装置アクセス要求を発する、受信制御エンジンと、

20

b. 前記受信制御エンジンからの命令に基づいて、ネットワークサービス応答を受信し、前記データ記憶装置アクセス要求へのデータ記憶装置アクセス応答がある場合、該データ記憶装置アクセス応答を処理する、送信制御エンジンとを含む。

【0014】

エンジンの少なくとも一つはハードウェアで実行され、その代わりに、又はそれに加えて、エンジンの少なくとも一つはハードウェアで加速される。関連する実施の形態では、ネットワークサービス要求は、CIFSプロトコル、SMBプロトコル、HTTPプロトコル、NFSプロトコル、FCPプロトコル、又はSMTPプロトコルである。さらに関連する実施の形態では、サービスモジュールは、受信器によって受信されたネットワーク要求が該要求を発する権限を有する情報源から発せられたか否かを決定する認証エンジンを含む。またさらなる実施の形態では、認証エンジンは、受信器によって受信されたネットワーク要求が要求された動作を実行する権限を有する情報源から発せられたか否かを決定する。同じく、サービスモジュールは、コンピュータのマザーボードに直接統合され、又は、コンピュータに接続され得るアダプタカードに統合される。

30

【0015】

もう一つの態様では、データ記憶装置アクセス要求を受信し、そのようなデータ記憶装置アクセス要求を満たすファイルシステムモジュールが提供される。そのファイルシステムモジュールは、

そのようなデータ記憶装置アクセス要求を受信し、それを翻訳し、そうして、記憶装置アクセス要求を発し得る、受信器と、

40

データ記憶装置アクセス応答を構築し、それを発する、前記受信器に接続される送信器であって、前記記憶装置アクセス要求への応答に基づいて適切なとき、そのような応答が情報を含む、前記送信器とを含む。

【0016】

受信器及び送信器の少なくとも一つは、ハードウェアで実行され、その代わりに、又はそれに加えて、受信器及び送信器の少なくとも一つは、ハードウェアで加速される。さらなる実施の形態では、記憶装置アクセス要求は、前記モジュールが接続され得る記憶装置によって用いられるプロトコルに矛盾がない。またさらなる実施の形態では、プロトコルは、NTFS、HPFS、FAT、FAT16、又はFAT32である。もう一つの関連する実施の形態では、ファイルシステムモジュールは、また、モジュールが接続され得る

50

記憶装置におけるファイルの物理的位置を画定するテーブルを格納する、前記受信器に接続されるファイルテーブルキャッシュを含む。種々の実施の形態では、プロトコルは、記憶装置内の連続的な物理定位置に置かれるべきファイルを要求しない。他の実施の形態では、ファイルシステムモジュールは、コンピュータのマザーボードに直接統合され、又はコンピュータに接続され得るアダプタカードに統合される。

【0017】

もう一つの態様では、要求情報源から記憶装置アクセス要求を受信し、そのような記憶装置アクセス要求を満たすために、記憶装置制御部と通信する記憶装置モジュールが提供される。その記憶装置モジュールは、

a. そのような記憶装置アクセス要求を受信し、それらを前記記憶装置制御部に適したフォーマットに変換する記憶装置要求インターフェースと、

b. 前記記憶装置制御部からの応答を取得し、そのような応答を前記要求情報源に適したフォーマットに変換する記憶装置承認インターフェースとを含む。

【0018】

記憶装置要求インターフェースと記憶装置承認インターフェースの少なくとも一つはハードウェアで実行され、その代わりに、又はそれに加えて、記憶装置要求インターフェースと記憶装置承認インターフェースの少なくとも一つはハードウェアで加速される。さらなる実施の形態では、記憶装置モジュールは、また、データの部分への高速読込アクセスを可能にするために前記記憶装置に含まれるデータの一部のローカルコピーを保持するキャッシュ制御部を含む。他の関連する実施の形態では、記憶装置要求インターフェースと記憶装置承認インターフェースは、ポートに接続され、そのポートは、光ファイバチャネルを介して、あるいはSCSI関連のプロトコルを利用する前記記憶装置制御部との通信を可能にする。さらなる実施の形態では、記憶装置モジュールは、コンピュータのマザーボードに直接統合され、又は、コンピュータに接続され得るアダプタカードに統合される。

【0019】

もう一つの態様では、記憶装置アクセス要求を置かれ得るラインと記憶装置配置を相互接続するためのシステムが提供される。この態様のシステムは、

a. 前記記憶装置アクセス要求を処理し、前記記憶装置配置へのアクセスに必要なところを生成し、応答の生成をさせる、前記記憶装置配置に接続されるサービス受信ブロックと、

b. 前記記憶装置配置内のファイルの物理的位置を画定するテーブルを格納する、前記受信ブロックに接続されるファイルテーブルキャッシュと、

c. 前記応答を送信するための、前記サービス受信ブロックに接続されるサービス送信ブロックとを含む。

【0020】

サービス受信ブロックとサービス送信ブロックの少なくとも一つは、ハードウェアで実行され、その代わりに、又はそれに加えて、サービス受信ブロックとサービス送信ブロックの少なくとも一つは、ハードウェアで加速される。さらなる実施の形態では、システムは、また、サービス受信ブロックとサービス送信ブロックのそれぞれに接続される応答情報メモリであって、該メモリは前記要求に関連したヘッダに存在する情報を格納し、該情報は前記応答を構築する前記サービス送信ブロックによって使用される、応答情報メモリを含む。もう一つの関連する実施の形態では、記憶装置アクセス要求は、ネットワーク要求である。まだもう一つの実施の形態では、記憶装置アクセス要求は、前記ラインが接続されるローカルプロセッサによって生成される。

【0021】

もう一つの態様では、多数のクライアントからの記憶装置アクセス要求を扱うための処理が提供される。その処理は、

前記クライアントのいずれかからの記憶装置アクセス要求の受信をテストするステップと、

10

20

30

40

50

あらゆる未決定の要求に従って記憶装置へのアクセスの完了をテストするステップとを含む。

【 0 0 2 2 】

この実施の形態では、記憶装置アクセス要求の受信テストと記憶装置へのアクセスの完了テストは、クライアントの数から独立して、多数のスレッドで実行される。さらなる実施の形態では、その処理は、また、要求の受信テストからの肯定的な決定を条件付けられ、該肯定的な決定を生じさせる要求を処理し、そのような要求に従って記憶装置アクセスを始めるステップを含む。関連する実施の形態では、前記処理は、また、未決定の要求に従って記憶装置へのアクセスの完了テストから肯定的な決定を条件付けられ、そのような未決定の要求を発するクライアントに応答を送るステップを含む。まだもう一つの関連する実施の形態では、スレッドの数が3よりも小さい。実際には、全処理は、単一のスレッドで具体化されてもよい。

10

【 0 0 2 3 】

もう一つの態様では、ネットワークを介してサービス要求を扱うための拡大縮小可能な装置であって、ネットワークはプロトコルを利用する拡大縮小可能な装置が提供される。この態様の装置は、

前記ネットワークプロトコルを用いるネットワークサービス要求を送受信するための第1の複数のネットワークサブシステムと、

前記ネットワークサービス要求を満たすための第2の複数のサービスサブシステムとを含む。

20

【 0 0 2 4 】

該ネットワークサブシステムと該サービスサブシステムのそれぞれ一つがハードウェアで実行され、あるいはハードウェアで加速されるものである。それに加えて、装置は、第1の複数のネットワークサブシステムのそれぞれを第2の複数のサービスサブシステムのそれぞれに接続するインタコネクトを含む。関連する実施の形態では、インタコネクトはスイッチであり、あるいはインタコネクトはバスである。

【 0 0 2 5 】

関連する態様では、記憶装置アクセス要求を生成され得るネットワークと記憶装置配置を相互接続するための拡大縮小可能なサービスサブシステムが提供される。サービスサブシステムは、

30

ネットワークサービス要求を受信し、そのようなサービス要求を満たし、そして、データ記憶装置アクセス要求を発する第1の複数のサービスモジュールと、

データ記憶装置アクセス要求を受信し、そのような記憶装置アクセス要求を満たし、そして、記憶装置配置アクセス要求を発する第2の複数のファイルシステムモジュールとを含む。

【 0 0 2 6 】

該サービスモジュールと該ファイルシステムモジュールのそれぞれ一つがハードウェアで実行され、あるいはハードウェアで加速されるものである。サービスサブシステムは、また、第1の複数のサービスモジュールのそれぞれを第2の複数のファイルシステムモジュールのそれぞれに接続するインタコネクトを含む。関連する実施の形態では、インタコネクトはスイッチであり、あるいはインタコネクトはバスである。さらなる実施の形態では、拡大縮小可能サービスサブシステムは、記憶装置配置アクセス要求を受信し、そのような記憶装置配置アクセス要求を満たすために該記憶装置配置を制御する第3の複数の記憶装置モジュールを含み、該記憶装置モジュールのそれぞれ一つはハードウェアで実行され、あるいはハードウェアで加速されるものである。同じく、サービスサブシステムは、前記ファイルシステムモジュールのそれぞれを前記記憶装置モジュールのそれぞれと接続する第2のインタコネクトを含む。同様に、関連する実施の形態では、インタコネクト及び第2のインタコネクトのそれぞれは、スイッチであり、あるいはインタコネクト及び第2のインタコネクトのそれぞれは、バスである。

40

【 0 0 2 7 】

50

好ましい実施の形態の詳細な記述

本発明の前述の特徴は、添付図面を参照するとともに、次の詳細な説明を参照することによっていっそう容易に理解されるだろう。

【 0 0 2 8 】

本記述と添付の特許請求の範囲の目的のために、次の用語は、文脈が別の方法で必要としないならば、示された意味を有する。

【 0 0 2 9 】

「ハードウェア実行 (hardware-implemented)」サブシステムは、主なサブシステム機能がソフトウェアプログラムの直接の制御外で駆動する専用ハードウェアで実行されるサブシステムを意味する。そのようなサブシステムがソフトウェアの制御下にあるプロセッサと相互作用するが、サブシステム自体がソフトウェアによって直接的に制御されていないことに注意されたい。「主な」機能は、最も頻繁に使用されるものである。

【 0 0 3 0 】

「ハードウェア加速 (hardware-accelerated)」サブシステムは、主なサブシステム機能が専用プロセッサ又は専用メモリを用いて実行され、それに加えて、(あるいはその代わりに) 特別な用途ハードウェア、すなわち、専用プロセッサ及びメモリがCPUに組み込まれるあらゆる中央演算処理装置 (CPU) 及びメモリと性質が異なるものを意味する。

【 0 0 3 1 】

「TCP/IP」は、他の場所間でwww.ietf.orgにおけるインターネット管理委員会のウェブサイト上に定義されるプロトコルである。それは、参照によってここに組み込まれる。「IP」は、同じ場所で定義されるインターネットプロトコルである。

【 0 0 3 2 】

「ファイル (file)」は、データの論理的結合である。

【 0 0 3 3 】

プロトコル「ヘッダ (header)」は、プロトコルのユーザーと結び付けられたデータの輸送のためにプロトコルによって指定されたフォーマットにおける情報である。

【 0 0 3 4 】

「SCSI 関連 (SCSI-related)」プロトコルは、SCSI、SCSI-2、SCSI-3、Wide SCSI、Fast SCSI、Fast Wide SCSI、Ultra SCSI、Ultra2 SCSI、Wide Ultra2 SCSI、又はあらゆる類似の若しくは後継のプロトコルを含む。SCSIは、「小型コンピュータシステムインターフェース (Small Computer System Interface)」であり、www.ansi.orgにウェブURLアドレスを持つ米国規格協会 (ANSI) に従うコンピュータ周辺機器の平行接続の標準である。

【 0 0 3 5 】

「層 3 及び 4 (layers 3 and 4)」への言及は、ISO 標準である開放型システム間相互接続 (OSI) 7 層モデルの層 3 及び 4 を意味する。ISO (国際標準化機構) は、www.iso.chにウェブURLアドレスを有する。

【 0 0 3 6 】

図 1 は、ネットワークを介してサービス要求を処理するように配置された本発明の一実施の形態の概略表示である。したがって、この実施の形態は、ファイルサーバー又はウェブサーバーが提供される配列を含む。本発明の実施の形態 11 は、ネットワークインターフェース 13 を介してネットワーク 10 に接続される。ネットワーク 10 は、例えば、複数のワークステーションへの通信リンクを含んでもよい。ここで、実施の形態 11 は、また、記憶装置内部接続 14 を介して複数の記憶装置に接続される。実施の形態 11 は、ハードウェアで実行され、あるいは、ハードウェアで加速され (又はハードウェア実行及びハードウェア加速を組み合わせる) てもよい。

【 0 0 3 7 】

図 2 は、図 1 に示された実施の形態のブロック図である。ネットワークサブシステム 21 は、ネットワークサービス要求及び応答を送受信する。ネットワークサブシステム 21

10

20

30

40

50

は、ネットワークサービス要求に応ずるサービスサブシステム 22 に接続される。ネットワークサブシステム 21、サービスサブシステム 22、あるいは両サブシステムは、ハードウェアで実行されるかハードウェアで加速されるかのいずれかであってよい。

【0038】

図3は、より詳細にファイルサーバーとして配置された図1の実施の形態のブロック図である。ネットワークサブシステム31は、ネットワークサービス要求及び応答を送受信する。ネットワークサブシステム31は、サービスサブシステム32に接続される。サービスサブシステムは、3つのモジュール、すなわち、サービスモジュール33、ファイルシステムモジュール34、及び記憶モジュール35を含む。サービスモジュール33は、サービスサブシステム32を通過したネットワークサービス要求を解析し、適切であるとき、対応する記憶アクセス要求を発する。ネットワークサービス要求は、CIFS、SMB、NFS、又はFCPのような種々のプロトコルのいずれかで伝達されてもよい。サービスモジュール33は、ファイルシステムモジュール34に接続される。もし、ネットワークサービス要求が記憶アクセス要求を含むならば、ファイルシステムモジュール34は、その要求を記憶媒体によって利用されるファイル記憶プロトコル（例えば、HTFS、NTFS、FAT、FAT16、又はFAT32）と一致するフォーマットに変換することによって、記憶装置へのアクセスの要求を変換する。記憶モジュール35は、サービスサブシステムが接続され得る記憶媒体に直接アクセスするためのバス要求と一致する（SCSTのような）フォーマットにファイルシステムモジュール34の出力を変換する。

【0039】

図4は、図3に類似しており、ウェブサーバーとして配置される図1の実施の形態のブロック図である。ネットワークサーバー41は、ネットワークサービス要求及び応答を送受信する。ネットワークサブシステム41は、サービスサブシステム42に接続される。サービスサブシステムは、3つのモジュール、すなわち、サービスモジュール43、ファイルシステムモジュール44、及び記憶モジュール45を含む。サービスモジュール43は、サービスサブシステム42に通過したネットワークサービス要求を解析し、適切であるとき、対応する記憶アクセス要求を発する。ここで、ネットワークサービス要求は、典型的に、HTTPプロトコルである。サービスモジュールは、記憶モジュール45に接続されるファイルシステムモジュール44に接続される。ファイルシステムモジュール44及び記憶モジュール45は、図3に関連して上述された、対応するモジュール34及び35に類似する方法で駆動する。

【0040】

図5は、図2～4の実施の形態のネットワークサブシステム及びサービスサブシステムである。ネットワークサブシステム51は、ネットワーク受信インターフェース54からカプセル収納データを受信し、TCP/IP又は他のプロトコルバス53に従ってそのデータをカプセルから出す。ネットワークサブシステム51は、また、ネットワークを介してデータにアクセスする（同じく、PCIバスに接続される）ローカルプロセッサに供給するために、PCIバス53に接続される。ネットワークサブシステム51は、また、サービスサブシステム52にデータを送信し、送信されたデータは、PCIバス53を介してネットワーク受信インターフェース54又なローカルプロセッサから来てもよい。サービスサブシステム52は、また、それぞれ図2、3、及び4のサービスサブシステム22、32、及び42に類似の方法で駆動する。

【0041】

図6は、図5のネットワークサブシステム51の詳細なブロック図である。図6のネットワークサブシステムは、（受信器601、受信バッファメモリ603、及び受信制御メモリ604を含む）受信器モジュール614、並びに、（送信器602、送信バッファメモリ605、及び送信制御メモリ606を含む）送信器モジュール613を含む。プロセッサ611は、受信器モジュール614と送信器モジュール613の両方に使用される。受信器601は、ネットワーク受信インターフェース607からカプセル収納データを受信し、それを翻訳処理する。受信器601は、受信制御メモリ604及び送信制御メモリ

10

20

30

40

50

606に含まれる制御情報を用いてそのデータをカプセルから出し、受信バッファメモリ605にそのカプセルから出されたデータを格納する。それは、受信バッファメモリ605からP C Iバス613を介してプロセッサ611によって取り出されるか、受信高速パス(fast path)インターフェース605に出力される。メモリ612は、データと命令の貯蔵のためにプロセッサ611に用いられる。送信器602は、送信後続パスインターフェース610かあるいはP C Iバス613を介してプロセッサ611から送信要求を受ける。

【0042】

送信器602は、送信バッファメモリ605にそのデータを格納する。送信器602は、送信制御メモリ606に含まれる制御情報を用いて送信データをカプセル収納し、ネットワーク送信インターフェース609を介してネットワーク上のカプセル収納されたデータを送信する。

10

【0043】

図7は、図6のネットワークサブシステムの受信モジュール614のブロック図である。パケットは、ネットワーク受信インターフェース607から受信エンジン701によって受信される。受信エンジン701は、パケットを解析し、そのパケットがエラーを含むか否か、TCP/IPパケットか、あるいは、TCP/IPパケットでないかを決定する。パケットは、パケット内に含まれるネットワークプロトコルヘッダの検査によってTCP/IPパケットかそうでないかを決定される。もしパケットがエラーを含むならば、それは落とされる。

【0044】

20

もし、パケットがTCP/IPパケットでないならば、そのパケットは、受信バッファメモリアービタ709を介して受信バッファメモリ603に格納される。パケットが受信されたという表示は、プロセッサイベント待ち行列702に書き込まれる。プロセッサ715は、P C Iバス704と受信P C Iインターフェースブロック703を用いて受信バッファメモリ603からパケットを受信することができる。

【0045】

もし、パケットがTCP/IPパケットならば、受信エンジン701は、パケットのプロトコルヘッダ内に含まれるネットワークアドレスとポートナンバーをこのパケットが属する接続、すなわち接続同定を特有の方法で識別する数に変換しようと試みるために、受信制御メモリ604に含まれるハッシュテーブルを用いる。もし、これが新しい接続道程ならば、パケットは、受信バッファメモリアービタ708を介して受信バッファメモリ603に格納される。パケットが受信されたという表示は、プロセッサイベント待ち行列702に書き込まれる。プロセッサ713は、P C Iバス704と受信P C Iインターフェースブロック703を用いて受信バッファメモリ603からパケットを受け取る。プロセッサは、TCP/IPプロトコルで指定されるように要求されるならば、新しい接続を確立することができ、あるいは他の適切なアクションをとり得る。

30

【0046】

もし、接続同定が既に存在するならば、受信エンジン701は、各接続状態についての情報を含むデータのテーブルへのインデックスとしてこの接続同定を用いる。この情報は、「TCP制御ブロック」(「TCB」)と呼ばれる。各接続のためのTCBは、送信制御メモリ606に格納される。受信エンジン701は、受信器TCBアクセスインターフェース710を介してこの接続のためのTCBにアクセスする。それは、TCP/IPプロトコルに従ってこのパケットを処理し、受信バッファメモリ605のこの接続のための受信されたバイトストリームに結果として生じるバイトを加える。もし、この接続におけるデータがプロセッサ713のために予定されているならば、幾らかのバイトが受信されたという表示は、プロセッサイベント待ち行列702に書き込まれる。プロセッサは、P C Iバス704と受信P C Iインターフェースブロック703を用いて受信バッファメモリ603からバイトを受け取る。もし、この接続におけるデータが高速パスインターフェース608のために予定されているならば、幾らかのバイトが受信されたという表示は、高速パスイベント待ち行列705に書き込まれる。受信DMAエンジン706は、受信バッファ

40

50

メモリ 603 からバイトを受け取り、高速バスインターフェース 608 にそれらを出力する。

【0047】

受信エンジン 701 によって受信されたいくつかのパケットは、IP パケットのフラグメントであり得る。もし、これがその場合ならば、フラグメントは、最初に受信バッファメモリ 603 で新たに組み立てられる。完全な IP パケットが新たに組み立てられたとき、通常のパケット処理が上述のように適用される。

【0048】

TCP プロトコルによれば、接続は、SYN_SENT、SYN_RECEIVED、及び ESTABLISHED を含む多くの異なる状態で存在し得る。ネットワークノードがネットワークサブシステムとの接続を確立することを望むとき、それは、最初に SYN フラグセットを持つ TCP/IP パケットを送信する。このパケットは、新しい接続同定を有するので、プロセッサ 713 によって検索される。プロセッサ 713 は、SYN_RECEIVED へのこの接続のために TCB 内に接続状態を設定することを含むすべての要求される初期化を実行する。SYN_RECEIVED から ESTABLISHED への移行は、TCP/IP プロトコルに従って受信エンジン 701 によって実行される。プロセッサ 713 がネットワークサブシステムを介してネットワークノードへの接続を確立することを望むとき、最初に SYN_SENT へのこの接続のために TCB 内に接続状態を設定することを含むすべての要求される初期化を実行する。それは、それから、SYN フラグセットを持つ TCP/IP パケットを送信する。SYN_SENT から ESTABLISHED への移行は、TCP/IP プロトコルに従って受信エンジン 701 によって実行される。

【0049】

もし、プロトコルヘッダ内に SYN フラグ又は FIN フラグあるいは RST フラグセットを有するパケットが受信され、そして、これがプロセッサ 713 によるアクションを要求するならば、受信エンジン 701 は、プロセッサ イベント待ち行列 702 へのエントリを書き込むことによって、このイベントをそのプロセッサに通知する。プロセッサ 713 は、TCP/IP プロトコルによって要求されるように適切なアクションをとり得る。

【0050】

受信されたパケットに TCP/IP プロトコルを適用する結果として、1 以上のパケットが今この接続で送信されるべきことが可能になる。例えば、受信データの承認が送信される必要があってもよく、あるいは、受信パケットは、そのようなデータが送信のために適用可能であるならば、より多くのデータがこの接続で送信されることを可能にする増加されたウィンドウサイズを示してもよい。受信エンジン 701 は、それに応じて TCB を変更し、それから、受信器送信待ち行列要求インターフェース 711 を介して図 8 の送信待ち行列 802 に接続同定を書き込むことによって送信試みを要求することによって、このことを達成する。

【0051】

受信データは、受信バッファメモリ 603 内の別々のユニット（バッファ）に格納される。バッファ内のすべてのデータがプロセッサ 713 によって検索されるか、高速バスインターフェース 605 に出力されるとすぐに、バッファは、解放、すなわち、新しいデータを格納するために再利用され得る。類似のシステムは、送信バッファメモリ 605 のために駆動するが、しかしながら、送信の場合、そのバッファは、その中のすべてのデータが TCP/IP プロトコルを用いて、送信データを受信するネットワークノードによって完全に承認されるときのみ解放される。送信されたデータが承認されたことをパケットのプロトコルヘッダが示すとき、受信エンジン 701 は、受信器フリー送信バッファ要求インターフェース 712 を介して図 8 のフリー送信バッファブロック 805 にこのことを示す。

【0052】

さらに、受信エンジン 701 が TCP/IP 自身と同様に TCP/IP をうまく処理する上部層プロトコル（ULP）を処理することを可能にする。この場合、完全な ULP プロトコルデータユニット（PDU）が受信されたときのみ、イベント待ち行列エントリは、処理又はイベント待ち行列 702 と高速バスイベント待ち行列 705 に書き込まれる。完全な ULP

10

20

30

40

50

のPDUのみがプロセッサ713によって受信され、高速バスインターフェース608に出力される。ULPの一例はNetBIOSである。ULP処理を可能にすることは、接続毎に基づいてなされてもよい、すなわち、いくつかの接続は、ULP処理を可能にし、その他のものはそうでなくてもよい。

【0053】

図8は、図6のネットワークサブシステムの送信モジュール613のブロック図である。TCP/IPを用いてネットワーク上で送信されるべきデータは、送信DMAエンジン807に入力される。このデータは、送信高速バスインターフェース610か、PCIバス704と送信PCIインターフェース808を介してプロセッサ713からのいずれかからの入力である。それぞれの場合、TCP/IP接続がデータを送信するために用いられるべきか否かを決定する接続同定は、同じく入力される。上述のように、各接続は、接続状態についての情報を含む、組み合わされたTCBを有する。

10

【0054】

送信DMAエンジンは、この接続のために格納されたバイトストリームに入力されたバイトを加えて、送信バッファメモリ605内にそのデータを格納する。入力の終わりに、それは、それに応じて接続のためにTCBを変更し、同じく、接続同定を送信待ち行列802に書き込む。

【0055】

送信行列802は、3つの情報源、すなわち、受信器送信待ち行列要求インターフェース711、タイマー機能ブロック806、及び送信DMAエンジン807から接続同定の形式で送信要求を受け取る。要求が受信されるので、それらは待ち行列に置かれる。待ち行列がからであるときはいつも、待ち行列の前における接続同定のための送信要求は、送信エンジン801へ通過される。送信エンジン801が送信要求の処理を完了すると、この接続同定は、待ち行列の前から取り除かれ、その処理が繰り返す。

20

【0056】

送信エンジン801は、送信待ち行列802から送信行列を受け取る。各要求のために、送信エンジン801は、その接続と要求される送信パケットにTCP/IPプロトコルを適用する。これをするために、それは、送信制御メモリアービタ805を介して送信制御メモリ606内の接続のためにTCBにアクセスし、送信バッファメモリアービタ804を介して送信バッファメモリ605から接続のための格納されたバイトストリームを検索する。

30

【0057】

接続のための格納されたバイトストリームは、送信バッファメモリ605内の別々のユニット(バッファ)に格納される。上述のように、各バッファは、TCP/IPプロトコルを用いて、送信データを受信するネットワークノードによってその中のすべてのデータが完全に承認されたときのみ解放され得る。送信されたデータが承認されたことをパケットのプロトコルヘッダが示すとき、受信エンジン701は、受信器フリー送信バッファ要求インターフェース712を介してフリー送信バッファブロック805にこのことを示す。フリー送信バッファブロック805は、完全に承認されたすべてのバッファを解放し、これらのバッファは、新しいデータを格納するために再利用され得る。

40

【0058】

TCP/IPは、ある条件が満たされるならば一定間隔で実行されるべきある動作を要求する多くのタイマー機能を有する。これらの機能は、タイマー機能ブロック806によって実行される。一定間隔で、タイマー機能ブロック806は、送信制御メモリアービタ803を介して各接続のためのTCBにアクセスする。もし、あらゆる操作が特定の接続のために実行される必要があるならば、その接続のためのTCBは、それに応じて変更され、接続同定は、送信待ち行列802に書き込まれる。

【0059】

そのうえ、送信DMAエンジン807がTCP/IPをうまく処理する上部層プロトコルを処理することができる。この場合、完全なULPプロトコルデータユニットのみは、プロセ

50

ッサ 713 か、送信高速バスインターフェース 610 のいずれかから送信 DMA エンジン 807 に入力される。送信 DMA エンジン 807 は、それから、PDU の前に ULP ヘッダを添付し、接続のための格納されたバイトストリームに「前に添付された」ULP ヘッダと入力されたバイトを加える。図 2 に関連して上述されるように、ULP の一例は NetBIOS である。ULP 処理を可能にすることは、接続毎に基づいてなされてもよい。すなわち、いくつかの接続が ULP 処理可能であり、そのたのものがそうでなくてもよい。

【0060】

もし、プロセッサ 713 が生のパケットを送信すること、すなわち、TCP/IP を用いてデータのハードウェアの自動送信することなく、データを送信することを望むならば、プロセッサ 713 が送信 DMA エンジン 807 にデータを入力するとき、それは、特定の接続同定を用いる。この特定の接続同定は、送信エンジン 801 に、プロセッサ 713 によって送信 DMA エンジン 807 に入力として正確に生のパケットを送信させる。

【0061】

図 9 は、ワークステーション又はサーバーのようなネットワークノードで使用するネットワークインターフェースアダプタとして図 5 のネットワークサブシステムの使用を示すブロック図である。この実施の形態では、ネットワークサブシステム 901 は、コンピュータに接続されるアダプタカード 900 に統合される。アダプタカード 900 は、ネットワークインターフェース 904 を介してネットワークに接続される。アダプタカード 900 はまた、PCI バス 907 と PCI ブリッジ 912 を介してコンピュータのマイクロプロセッサに接続される。PCI バス 907 はまた、ビデオシステム 913 のような周辺装置にアクセスするためにコンピュータによって用いられてもよい。受信モジュール 902 と送信モジュール 903 は、図 6 の受信モジュール 614 と送信モジュール 613 と類似の方法で動作する。その代わりに又はそれに加えて、アダプタカード 900 は、ネットワーク上のリモートノード又はマイクロプロセッサ 10 によって記憶装置配列への高速アクセスを提供するために、単一プロトコル高速受信パイプ 906 と単一プロトコル高速送信パイプ 905 を介して、図 2、3、4、又は 5 のアイテム 22、32、43、又は 52 のそれぞれのいずれかに相当するサービスモジュールに接続されてもよい。

【0062】

図 10 は、図 3 に示されるような一実施の形態で使用する図 3 の SMB サービスモジュール 33 とファイルシステムモジュール 34 のハードウェアで実行される組み合わせのブロック図である。図 10 の実施の形態では、SMB 要求は、サービス受信ブロック 101 への入力 105 で受信される。最終的に、この実施の形態による処理は、出力 106 上の対応する SMB 応答の送信を結果としてもたらす。この応答の一部はヘッダを含む。出力ヘッダを作るために、入力ヘッダは、SMB 応答情報メモリ 103 に格納される。ブロック 101 は、SMB 要求を処理し、応答を生成する。要求の性質に依存して、ブロック 101 は、ファイルテーブルキャッシュ 104 にアクセスし、ディスクアクセス要求を発してもよい。さもなければ、応答が送信ブロック 102 を直接中継される。サービス送信ブロック 102 は、出力 106 上にブロック 101 によって生成される応答を送信する。ディスクアクセス要求がブロック 101 によって発せられ、ディスク応答のライン 108 上で受信された場合、送信ブロック 102 は、ライン 106 上に適切な SMB 応答を発する。受信及び送信の両モジュール 101 及び 102 は、PCI バス 109 を介してホストシステムと選択的に通信する。供給されるとき、そのような通信は、伝統的なオペレーティングシステムの範囲外に高速で、ハードウェアで実行されるファイルシステムアクセスをホストシステムに与えるように、ホストシステムがネットワークを介する代わりに実施の形態と直接通信することを可能にする。

【0063】

図 11 は、図 3 に示されるような一実施の形態において使用する図 3 の SMB サービスモジュール 33 とファイルシステムモジュール 34 のハードウェアで加速される組み合わせのブロック図である。動作は、同じく番号を付されたブロック及びライン 105、107、108、及び 106 に関して図 10 に関連して上述されたものに類似している。しか

10

20

30

40

50

しながら、専用のファイルシステムプロセッサ 110 は、専用バス 112 を介して動作する専用メモリ 111 に関連して、ブロック 101 及び 102 の処理を制御する。それに加えて、これらのアイテムは、ソフトウェアで変更され得るので、そのような処理の取扱いに柔軟性を提供する。

【0064】

図 12A は、それぞれ図 3 又は図 4 のアイテム 33 又は 43 のようなハードウェアで実行されるサービスモジュールのブロック図である。サービスモジュール 1200 は、ネットワークサービス要求を受信し、そのようなサービス要求を実現し、データ記憶アクセス要求を発してもよい。サービスモジュール 1200 は、送信器 1202 に接続された受信器 1201 と、受信器 1201 及び送信器 1202 の両方に接続されたデータ記憶アクセスインターフェース 1203 とを含む。受信器 1201 は、ネットワークサービス要求を受信し、それを解釈する。サービス要求を受信し次第、受信器 1201 は、データ記憶アクセスインターフェース 1203 にその要求を送るか、あるいは、ネットワークサービス要求を満たす情報を送信器 1202 に送る。もし、その要求がデータ記憶アクセスインターフェース 1203 に送られるならば、データ記憶アクセスインターフェース 1203 は、データ記憶アクセス要求を構成し、それを発する。データ記憶アクセスインターフェースは、また、データ記憶アクセス要求への応答を受信し、オリジナルのネットワークサービス要求を満たすために必要とされる情報を抽出する。その情報は、それから送信器 1202 に送られる。送信器 1202 は、受信器 1202 又はデータ記憶アクセスインターフェース 1203 からそれに送られる情報を処理し、ネットワークサービス応答を構築し、それを発する。

【0065】

図 12B は、それぞれ図 3 又は図 4 のアイテム 34 又は 44 のようなハードウェアで実行されるファイルモジュールのブロック図である。ファイルシステムモジュール 1210 は、データ記憶アクセス要求を受信し、そのようなデータサービスアクセス要求を実現し、記憶装置アクセス要求を発してもよい。ファイルシステムモジュール 1210 は、送信器 1212 に接続された受信器 1211 と、受信器 1211 及び送信器 1212 の両方に接続されたデータ記憶装置アクセスインターフェース 1213 とを含む。受信器 1211 は、データ記憶アクセス要求を受信して解釈し、データ記憶装置アクセスインターフェース 1213 へその要求を送るか、あるいは、データ記憶アクセス要求を満たす情報を送信器 1212 に送る。もし、その要求がデータ記憶装置アクセスインターフェース 1213 に送られるならば、データ記憶装置アクセスインターフェース 1213 は、データ記憶装置アクセス要求を構築し、それを発する。データ記憶装置アクセスインターフェース 1213 は、また、データ記憶装置アクセス要求への応答を受信し、オリジナルのデータ記憶アクセス要求を満たすために必要とされる情報を抽出する。その情報は、それから送信器 1212 に送られる。送信器 1212 は、受信器 1211 又はデータ記憶装置アクセスインターフェースモジュール 1213 からそれに送られた情報を処理し、データ記憶アクセス応答を構築し、それを発する。

【0066】

図 12C は、結合されたサービスモジュール及びファイルモジュールを提供する、図 10 のハードウェアで実行されるサービスサブシステムの詳細なブロック図である。図 12C 内の点線 129 は、この実行の機能間の分割を示す。ライン 129 の左側にはサービスモジュール部があり、ライン 129 の右側にはファイルシステムモジュール部がある。(しかしながら、SMB 受信制御エンジン 121 と SMB 送信制御エンジン 122 を結ぶ二頭の矢がサービスモジュール部とファイルシステムモジュール部のそれぞれのためのエンジン 121 及び 122 間の 2 方向通信を適切に供給することを理解されるだろう。)

【0067】

図 12C では、SMB フレームは、ネットワーク受信インターフェース 121f を介してネットワークサブシステムから受信され、SMB フレーム翻訳エンジン 121b に送られる。ここで、フレームは解析され、多くのタスクが実行される。ヘッダの最初のセクシ

ョンは、SMB 応答の接続基礎毎に適切な情報をメモリ 103 に格納する、SMB 応答情報制御 123 にコピーされる。完全なフレームは、受信バッファメモリ 121c 内のバッファに書き込まれ、受信制御メモリ 121d は更新される。SMB フレームヘッダの適切な部分は、SMB 受信制御エンジン 121 に送られる。

【0068】

図 12C の SMB 受信制御エンジン 121 は、ヘッダからその情報を解剖し、適切な場合には、批准エンジン 124 からファイルアクセス許可を要求する。ファイルアクセスが要求される SMB フレームのために、SMB 受信制御エンジン 121 は、SMB フレームヘッダからファイルパス情報かファイル同定のいずれかを抽出し、MFT 制御エンジン 125 に要求されたファイルデータの物理的位置を要求する。

10

【0069】

MFT 制御エンジン 125 は、SMB 受信制御エンジン 121 からの要求を待ち行列に入れることができ、同様に、SMB 受信制御エンジン 121 は、MFT 制御エンジン 125 からの要求を待ち行列に入れることができる。これは、2つのエンジンが互いに非同期に動作することを可能にし、従って MFT 要求が顕著な間入ってくる SMB フレームが処理されることを可能にする。

【0070】

MFT 制御エンジン 125 は、SMB 受信制御エンジン 121 からの要求を処理する。典型的に、SMB OPEN コマンドのために、要求は、必要な物理的ファイル位置情報を得るためのディスクアクセスが必要である。これが必要な場合、MFT 制御エンジン 125 は、必要な圧縮 SCSI 要求を生成する圧縮 SCSI フレーム生成エンジン 121a に要求を送る。圧縮 SCSI プロトコル(「CSP」)は、SCSI コマンドが図 17A と他の図に関連して以下に記述される方法で生成され得るデータフォーマットに関する。圧縮 SCSI データが SCSI から得られないが、むしろ SCSI データが得られ得る情報源であるので、我々は、時々圧縮された SCSI データを「プロト-SCSI」データとして言及する。適切なプロト-SCSI 応答は、それが処理される MFT 制御エンジン 125 に送り返され、MFT キャッシュ 104 は更新され、物理的ファイル情報は、SMB 受信制御エンジン 121 に送り返さる。

20

【0071】

典型的に、最近アクセスされた小さなファイルに関して SMB READ 又は WRITE コマンドのために、ファイル情報は、MFT キャッシュ 104 内に存在する。したがって、ディスクアクセスは要求されない。

30

【0072】

SMB 受信制御エンジン 121 が MFT 要求からその応答を受信し、ファイルデータのためのディスクアクセスが要求されるとき、典型的な READ 又は WRITE コマンドに必要なように、1 以上のプロト-SCSI 要求は、プロト-SCSI フレーム生成エンジン 121a に送られる。

【0073】

プロト-SCSI フレーム生成エンジン 121a は、プロト-SCSI ヘッダを構築し、必要な場合、例えば、WRITE コマンドのために、受信バッファメモリ 121c からファイルデータを引くために、ファイルデータ DMA エンジン 121e をプログラムする。プロト-SCSI フレームは、プロト-SCSI 送信インターフェース 121g を介してプロト-SCSI モジュールに送られる。ディスクアクセスが要求されない場合、SMB 応答要求は、直接 SMB 送信制御エンジン 122 に送られる。

40

【0074】

プロト-SCSI フレームは、プロト-SCSI モジュールから受信され、プロト-SCSI 受信インターフェース 122f を介してプロト-SCSI フレーム翻訳エンジン 122b に送られる。ここで、そのフレームは解析され、多くのタスクが実行される。MFT 応答は、MFT 制御エンジン 125 に送り返される。他のすべてのフレームは、受信バッファメモリ 121c 内のバッファに書き込まれ、受信制御メモリ 121d は更新される

50

。プロト - S C S I フレームヘッダの適切な部分は、S M B 送信制御エンジン 1 2 2 に送られる。

【 0 0 7 5 】

各 S M B 接続は、独特な同定を以前に割り当てられた。すべてのプロト - S C S I フレームは、この同定を含み、S M B 送信制御エンジン 1 2 2 は、S M B 受信制御エンジン 1 2 1 からの状態情報を要求し、必要な場合これを更新するために、この独特な同定を用いる。S M B 応答のための必要な情報がプロト - S C S I モジュールから受信されたとき、S M B 送信制御エンジン 1 2 2 は、S M B フレーム生成エンジン 1 2 1 a に要求を送る。

【 0 0 7 6 】

S M B フレーム生成エンジン 1 2 1 a は、S M B 応答情報メモリ 1 0 3 内に含まれるデータ及び S M B 送信バッファメモリ 1 2 2 c に格納されたファイルデータから S M B 応答：フレームを構築する。それは、結果としてそれをネットワークサブシステムに転送する S M B 送信インターフェース 1 0 6 にそのフレームを送る。

【 0 0 7 7 】

図 1 3 は、図 1 1 のハードウェア加速サービスサブシステムの詳細なブロック図である。入力 1 0 5 を越えて供給される I P ブロックから入ってくる S M B フレームは、S M B 受信 F I F O 1 3 1 7 を介して、S M B 受信バッファメモリ 1 2 1 c 内のフリーバッファに書き込まれる。S M B 受信バッファメモリ 1 2 1 c は、一実施の形態では、2 K b の長さであり、一つの S M B フレームが多くの受信バッファにまたがり得る一連の受信バッファを含む。フレームが S M B 受信バッファメモリ 1 2 1 c に書き込まれるので、S M B 受信バッファ記述子は、S M B 受信制御メモリ 1 2 1 d で更新される。

【 0 0 7 8 】

3 2 ビット接続同定及び 3 2 ビットフレームバイトカウントは、フレームの始まりで I P ブロックから S M B ブロックに送られる。これら 2 つのフィールドは、受信バッファメモリ 1 2 1 c の受信バッファの最初の 2 つの位置に書き込まれる。

【 0 0 7 9 】

フレームが格納されるが、S M B ヘッダもまた、S M B 送信処理による後の使用のために S M B 応答情報メモリ 1 0 3 に書き込まれる。I P ブロックによって S M B ブロックに送られる独特な接続同定は、S M B 応答情報メモリ 1 0 3 内の適切な情報フィールドへのポインタとして用いられる。このメモリは、1 6 ワードのブロックで、各独特な接続同定のために一つのブロックとして配置される。1 2 8 M b の S D R A M で適合させて、これは 2 M の接続を可能にする。目下、S M B フレームの最初の 3 2 バイトは各情報フィールドに書き込まれる。

【 0 0 8 0 】

完全なフレームが受信バッファメモリ 1 2 1 c に書き込まれたとき、S M B バッファロケータは、S M B 受信イベント待ち行列 1 3 1 4 に書き込まれ、ホストプロセッサ 1 3 0 1 への割り込みが生成される。S M B バッファロケータは、バッファポインタと「最後の」ビットを含む S M B フレームに関する情報を含む。バッファポインタは、S M B フレームの始まりを含む受信バッファメモリ 1 2 1 c を指し示す。「最後の」ビットは、このバッファが S M B フレームの終わりも含むか否か（すなわち、S M B フレームが 2 K b の長さより小さいか否か、）を示す。

【 0 0 8 1 】

ホストプロセッサ 1 3 0 1 は、イベント待ち行列 1 3 1 4 に結び付けられた適切な S M B 受信イベントレジスタを読むことによって S M B 受信待ち行列 1 3 1 4 内の S M B バッファロケータを読むことができる。S M B バッファロケータから読み取られたバッファポインタから、ホストプロセッサ 1 3 0 1 は、受信バッファメモリ 1 2 1 c 内の S M B フレームの最初のバッファのアドレスを決定することができ、従って S M B ヘッダとフレームの最初の部分を読むことができる。

【 0 0 8 2 】

もし、S M B フレームが 2 K b よりも長く、S M B フレームの最初の 2 K b 以上を読む

10

20

30

40

50

必要があるならば、この受信バッファに結び付けられた受信バッファ記述子は、受信制御メモリ 121d から読まれるべきである。この受信バッファ記述子は、SMB フレームの次のバッファへのポインタを含む。受信バッファが SMB フレームの終わりを含むことを指摘することを示す「最後の」ビットを前のバッファの記述子が含まなければ、この次のバッファは、同様に、それに結び付けられた受信バッファ記述子を有する。

【0083】

受信された SMB フレームを呼んだ後、もし、フレーム内に含まれるデータがさらに使用されるべきではないならば、受信されたフレームのバッファは、受信バッファ制御メモリ 121d に含まれる受信フリーバッファ待ち行列にポインタを書き込むことによって、結び付けられた受信返却フリーバッファレジスタに書き込むことによって、再び使用する

10

【0084】

プロト - SCSI フレームを送信するために、ホストプロセッサ 1301 は、受信フェッチフリーバッファレジスタから読み込むことによってフリー SMB 受信バッファへのポインタを最初に得る。このアクションは、受信制御メモリ 121d に含まれるフリーバッファ待ち行列からフリーバッファへのポインタを持ってくる。このバッファでは、プロト - SCSI 要求フレームの始まりが構築され得る。

【0085】

プロト - SCSI フレームをプロト - SCSI エンティティに移すようにプロト - SCSI 送信エンティティに要求するために、ホストプロセッサ 1301 は、プロト - SCSI 送信イベント待ち行列 1315 に結び付けられた受信プロト - SCSI イベントレジスタにそれらを書き込むことによって、プロト - SCSI 送信イベント待ち行列 1315 にバッファロケータとバッファオフセット対を書き込む。

20

【0086】

バッファロケータは、プロト - SCSI フレームのためのデータを含むバッファへのポインタを含む。バッファオフセットは、バッファと長さフィールドの中のデータの始まりに対するオフセットを含む。バッファロケータは、また、さらなるバッファロケータ / バッファオフセット対が、このプロト - SCSI フレームのためのより多くのデータへのポインタを含むプロト - SCSI 送信イベント待ち行列 1315 に書き込まれるか否かを示すためのさごのビットを含む。

30

【0087】

もし、プロト - SCSI フレームがもう一つの SMB 受信バッファからのデータを含むべきならば、SMB WRITE コマンドのために典型的であるように、ホストプロセッサ 1301 は、この SMB 受信バッファを記述するもう一つのバッファロケータ / バッファオフセット対をプロト - SCSI 送信イベント待ち行列 1315 に書き込まなければならない。もし、プロト - SCSI フレームに含まれるべきデータが一つの SMB 受信バッファ以上にまたがるならば、プロト - SCSI 送信エンティティは、データとともにリンクするために、受信制御メモリ 121d に位置される結び付けられた SMB 受信バッファ記述子内のバッファポインタを用いることができる。もし、余分なデータが SMB 受信フレームからのものであるならば、これらの記述子は、SMB 受信エンティティによって予め満たされたであろう。

40

【0088】

SMB 受信バッファからのデータが一つ以上のプロト - SCSI フレームのために用いられ得るので、それらが用いられた後の SMB 受信バッファを自由化することは単純な処理ではない。プロト - SCSI 送信に含まれない受信 SMB フレームのセクションを含む SMB 受信バッファは、関連付けられた受信戻りフリーバッファレジスタを介して受信制御メモリに含まれるフリーバッファ待ち行列にそれらを書き戻すことによって自由にされ得る。プロト - SCSI フレームに含まれるべきデータを含む SMB 受信バッファは、それらの中のデータが送信されるまで自由にされないような同一の方法で自由にされ得ない。

50

【 0 0 8 9 】

それで、SMBデータを含む種々のプロト - SCSIフレームに対するバッファロケータ / バッファオフセット対がプロト - SCSI送信イベント待ち行列 1 3 1 5 に書き込まれた後、オリジナルのSMB受信バッファへのポインタも同じく、プロト - SCSIに書き込まれる。これらのポインタは、受信制御メモリに含まれるフリーバッファ待ち行列に戻されて自由にされるべきであることを示すために記録される。プロト - SCSI送信が次々と処理されるので、SMB受信バッファは、それらの中のあらゆるデータが送信された後にのみ自由にされる。

【 0 0 9 0 】

IPブロックから入ってくるプロト - SCSIフレームは、プロト - SCSI受信 FIFO 1 3 2 7 を介してSMB送信バッファメモリ 1 2 2 c 内のフリーバッファに書き込まれる。SMB送信バッファは2 Kb長であり、一つのプロト - SCSIフレームは、多くの送信バッファにまたがってもよい。フレームがSMB送信バッファメモリ 1 2 2 c に書き込まれるので、SMB送信バッファ記述子は、SMB送信制御メモリ 1 2 2 d 内で更新される。

10

【 0 0 9 1 】

完全なフレームがSMB送信バッファメモリ 1 2 2 c に書き込まれるとき、SMBバッファロケータは、プロト - SCSI受信イベント待ち行列 1 3 2 4 に書き込まれ、ホストプロセッサ 1 3 0 1 への割り込みが生成される。SMBバッファロケータは、バッファポインタと「最後の」ビットを含むプロト - SCSIフレームに関する情報を含む。バッファポインタは、プロト - SCSIフレームの始まりを含む送信バッファメモリ 1 2 1 c 内のバッファを指し示す。「最後の」ビットは、このバッファが同じくプロト - SCSIフレームの終わりを含むか否か（すなわち、フレームが長さで2 Kbより小さいか否か）を示す。

20

【 0 0 9 2 】

ホストプロセッサ 1 3 0 1 は、イベント待ち行列 1 3 2 4 に関連した適切なプロト - SCSI受信イベント待ち行列 1 3 2 4 を読むことによって、プロト - SCSI受信イベント待ち行列のバッファロケータを読むことができる。バッファロケータから読まれたバッファポインタから、ホストプロセッサ 1 3 0 1 は、送信バッファメモリ 1 2 2 c 内のプロト - SCSIフレームの最初のバッファのアドレスを決定することができ、従ってヘッダとそのフレームの最初の部分を読むことができる。

30

【 0 0 9 3 】

もし、プロト - SCSIフレームが2 Kbより長く、フレームの最初の2 Kbより多くを読む必要があるならば、この送信バッファに関連した送信記述子は、受信制御メモリ 1 2 1 d から読まれるべきである。バッファがプロト - SCSIフレームの終わりを含むことを指摘することを示す「最後の」ビットを前のバッファの記述子が含まないならば、記述子は、プロト - SCSIフレームの次のバッファへのポインタを含む。この次のバッファは、同様に、それに関連した送信記述子を有する。

【 0 0 9 4 】

受信されたプロト - SCSIフレームを読んだ後、もし、そのフレーム内に含まれるデータが更に用いられるべきでないならば、受信されたフレームのバッファは、それに関連した送信戻りフリーバッファレジスタに書き込むことによって、送信制御メモリ 1 2 2 d に含まれる送信フリーバッファ待ち行列に戻されるべきである。

40

【 0 0 9 5 】

SMBフレームを送信するために、ホストプロセッサは、関連したレジスタから読むことによって、送信制御メモリ 1 2 2 d に含まれる送信フリーバッファ待ち行列から送信バッファメモリ 1 2 2 c 内のフリーSMB送信バッファへのポインタを最初に得る。このバッファでは、SMB応答フレームの始まりが構築され得る。

【 0 0 9 6 】

3 2 ビット接続道程と3 2 ビットSMB送信制御フィールドは、バッファ内のSMBフ

50

フレーム前に置かれる。SMB送信制御フィールドは、24ビットフレームバイトカウントと前置ヘッダビットを含む。もし、前置ヘッダビットが設定されるならば、接続同定とSMB送信制御フィールドがIPブロックに送られた後に、応答情報メモリ103に格納されたSMBヘッダは、自動的に挿入される。

【0097】

SMB送信エンティティにSMBフレームをSMBエンティティに移すことを要求するために、ホストプロセッサ1301は、関連した送信SMB送信イベントレジスタにそれらを書き込むことによって、バッファローケータとバッファオフセット対をSMB送信イベント待ち行列1325に書き込む。

【0098】

バッファローケータは、SMBフレームのためのデータを含むバッファへのポインタを含む。バッファオフセットは、バッファと長さフィールド内のデータの始まりに対するオフセットを含む。バッファローケータは、また、バッファローケータ/バッファオフセット対がさらにこのSMBフレームのためのより多くのデータへのポインタを含んで書き込まれるか否かを示す最後のビットを含む。

【0099】

もし、SMBフレームがバッファメモリ122c内のもう一つのSMB送信バッファからのデータを含むべきであるならば、ホストプロセッサは、SMB送信イベント待ち行列1325にこのSMB送信バッファを記述するもう一つのバッファローケータ/バッファオフセット対を書き込まなければならない。もし、SMBフレームに含まれるべきデータが一つ以上のSMB送信バッファにまたがるならば、SMB送信エンティティは、データとともにリンクするために、関連した送信バッファ記述子内のバッファポインタを用い得る。もし、余分なデータがプロト-SCSI受信フレームからのものであるならば、これらの記述子は、プロト-SCSI受信エンティティによって以前に満たされたであろう。

【0100】

送信バッファメモリ122cのSMB送信バッファからのデータが一つ以上のSMBフレームのために用いられ得るので、それらが用いられた後にSMB送信バッファを自由化することは、単純な処理ではない。SMB送信に含まれない受信されたプロト-SCSIフレームのセクションを含むSMB送信バッファは、関連した送信戻りフリーバッファレジスタを介して、送信制御メモリに含まれる送信フリーバッファ待ち行列にそれらを書き戻すことによって自由にされ得る。SMBフレームに含まれるべきデータを含むSMB送信バッファは、その中のデータが送信されるまでそれらが自由にされ得ないので、同一の方法で自由にされ得ない。

【0101】

それで、プロト-SCSIデータを含む種々のSMBフレームへのバッファローケータ/バッファオフセット対がSMB送信イベント待ち行列1325に書き込まれた後、オリジナルのSMB送信バッファへのポインタは、また、SMB送信イベント待ち行列1325に書き込まれる。これらのポインタは、送信フリーバッファ待ち行列に戻って自由化されるべきことを示すことを記録する。SMB送信イベント待ち行列1325が次々に処理されるので、SMB送信バッファは、その中のあらゆるデータが送信された後にのみ自由にされる。

【0102】

図14は、多数のサービス要求を多数のスレッドとして扱うための、ソフトウェアで実行された典型的な先行技術のアプローチを表すフローチャートである。伝統的な多数スレッドの構成では、典型的に各クライアントにサービスを提供するための少なくとも一つのスレッドがある。スレッドは、クライアントがサーバーに接続詞、それから分離するように、始められ、終了させられる。各クライアントは、サービス要求を処理するためのサーバー上のスレッドとディスク要求を処理するためのスレッドとを有してもよい。サービス処理1400は、ボックス1401でクライアント接続要求の存在をテストする繰り返し

10

20

30

40

50

ループを含み、もし、テストが肯定的であるならば、処理は、ボックス 1 4 0 2、クライアント処理 1 4 3 0 を始める。クライアント処理 1 4 3 0 がボックス 1 4 3 5 でディスクアクセスを要求するとき、それは最初にディスクへアクセスするための適切なディスク処理を要求し、それから、ディスクアクセスが完了するまでボックス 1 4 3 6 で待機する。ディスク処理 1 4 0 2 は、それがサービス要求を発するクライアントにボックス 1 4 3 7 で応答を送信することを可能にするために、クライアント処理 1 4 3 0 を目覚めさせる。それで、ディスクアクセスを要求する各クライアント要求のための少なくとも二つの処理スイッチがある。これらの多数のスレッド処理をハードウェアで実行することは、問題である。なぜならば、通常、それらはマルチタスク処理オペレーティングシステムによって扱われるからである。

10

【 0 1 0 3 】

図 1 5 は、図 2 のサービスサブシステムと、例えば、図 1 2 及び 1 3 の実施の形態とを関連して使用するために、多数のサービス要求の処理を示すフローチャートである。単純なスレッドアーキテクチャでは、一つのサービス処理 1 5 0 0 が単一のスレッドで多数のクライアントからの要求を扱い、一つのディスク処理 1 5 0 2 がサービス処理 1 5 0 0 からのすべての要求を扱う。要求をなす各クライアントのために別々の処理を用いる先行技術のアプローチは無視され、その機能は、一つのサービス処理 1 5 0 0 によってここで扱われた。それに加えて、これら 2 つの処理、サービス処理とディスク処理は、示されるように、同一のスレッド内に含まれてもよく、あるいは、ロードバランシングを容易にするために 2 つの別々のスレッド間に分担されてもよい。

20

【 0 1 0 4 】

図 1 5 の単一スレッドサービス処理は、同時に傑出した多数のクライアントからのディスク要求を有することができる。単一スレッドは、2 つのテストを持つメインループを含む。ボックス 1 5 0 1 の第 1 のテストは、クライアントからの要求を受信したか否かである。ボックス 1 5 0 8 の第 2 のテストは、前に始められたディスクアクセス要求が完了したか否かである。その結果、ディスクアクセスが完了させられるべきボックス 1 5 0 8 で決定されたので、ボックス 1 5 0 7 でサービス処理は、クライアントに適切な応答を送り返す。サービス処理 1 5 0 0 がボックス 1 5 0 1 を介してディスクアクセス要求を扱い、要求をボックス 1 5 0 2 で処理させ、ボックス 1 5 0 4 でディスクアクセスの開始を起こすとすぐに、サービス処理は、前のディスクアクセスが完了するまで停止し及び待つことなく、ボックス 1 5 0 1 を介してもう一つのクライアントからのもう一つの要求を扱うために自由である。ディスクアクセスが完了されたボックス 1 5 0 8 での決定において、ボックス 1 5 0 7 のディスク処理は、結果のサービス処理を通知し、クライアントに応答を送信する。従って、サービス及びディスク処理は、クライアントから送られる要求がある限り、絶えず走っている。

30

【 0 1 0 5 】

図 1 6 は、ファイル記憶装置を有するコンピュータシステムに関連して、図 3 に示されるようなファイルシステムモジュールの使用を示すブロック図である。(図 1 6 のものに類似する実行は、ファイル記憶装置を有するコンピュータシステムに関連して、図 3 に示されるような記憶装置モジュールを提供するために用いられてもよい。) この実施の形態では、ファイルシステムモジュール 1 6 0 1 は、マイクロプロセッサ 1 6 0 5、メモリ 1 6 0 6、及び、ここでは従来のディスクドライブ制御部であるディスクサブシステム 1 6 0 2 を介してアクセスされるディスクドライブ 1 6 1 0 と同様に、ビデオ 1 6 0 9 のような周辺機器を含むコンピュータシステムに統合化される。ファイルシステムモジュール 1 6 0 1 は、また、P C I バス 1 6 0 7 上の P C I ブリッジ 1 6 0 4 を介してコンピュータマルチプロセッサ 1 6 0 5 とコンピュータメモリ 1 6 0 6 に接続される。P C I バス 1 6 0 7 は、また、マイクロプロセッサ 1 6 0 5 をコンピュータ周辺機器 1 6 0 9 に接続する。ファイルシステムモジュールの受信エンジン 1 6 1 0 は、図 1 0、1 1、1 2 B、及び 1 3 に関連して上述された方法と類似の方法で、マイクロプロセッサ 1 6 0 5 からのディスクアクセス要求を処理する。同じく、送信エンジン 1 6 1 1 は、図 1 0、1 1、1 2 B

40

50

、及び13に関連して上述された方法と類似の方法で、そのようなディスクアクセス要求に対する応答を供給する。

【0106】

図17Aは、図3の記憶装置モジュールにおけるデータフローのブロック図である。図17Aと17Bにおけるカリフォルニア州パロアルトのヒューレットパッカード社（Hewlett Packard Co., Palo Alto, California）から利用可能なタキオンXL光ファイバチャネル制御部がI/O装置として用いられ得、本発明の実施の形態が他のI/O装置を等しく用い得ることに気付かれない。プロト-SCSI要求は、プロト-SCSI要求プロセッサ1702によってプロト-SCSI入力を越えて受信される。この要求に関する情報はSEST情報テーブルに格納され、もし、これがWRITE要求ならば、同じくプロト-SCSI入力1700を介して供給されるWRITEデータは、WRITEバッファメモリ1736に格納される。

10

【0107】

交換要求（ERQ）発生器1716は、WRITEバッファメモリ1736からの情報を取得する。もし、書き込まれるべきすべてのバッファが現在蓄えられ、あるいは書き込まれるべきデータが書き込まれるべきバッファを完全に満たすならば、WRITEは、すぐに実行され得る。書き込まれるべきデータは、キャッシュメモリ1740の適切な領域にWRITEバッファメモリ1736からコピーされる。ファイバチャネルI/O制御部1720は、その制御部1720と通信するディスク記憶装置の適切な領域にそのデータを書き込むように構成される。さもなければ、ディスクからのREADは、適切なディスクから要求されるデータを得るために、WRITEの前になされなくてはならない。

20

【0108】

プロト-SCSI承認発生器1730は、プロト-SCSI要求を生成する責任を負う。プロト-SCSI応答を生成することができる3つの情報源、プロセッサ1738、ファイバチャネルI/O制御部1720、及びキャッシュメモリ1740があり、そのそれぞれは、SESTインデックスを供給する。すべての移動のために、ステータス情報とともに、プロト-SCSI要求が承認に関連付けられることを可能にする同定は、プロト-SCSI承認インターフェース1734に戻される。

【0109】

図17Bは、図3の記憶装置モジュールの制御フローを示す詳細なブロック図である。プロト-SCSI要求がプロト-SCSI要求プロセッサ1702によってプロト-SCSI入力1700を越えて受信されるとき、それは、独特な識別子（SESTインデックスと呼ばれる）を割り当てられる。この要求に関する情報は、SEST情報テーブルに格納され、もし、これがWRITE要求ならば、プロト-SCSI入力1700に同じく供給されるWRITEデータは、WRITEバッファメモリ1736に格納される。SESTインデックスは、それからプロト-SCSI要求待ち行列1704に書き込まれる。

30

【0110】

キャッシュ制御部1706は、プロト-SCSI要求待ち行列1704と使用されたバッファ待ち行列1708からのエントリを取得する。エントリがプロト-SCSI要求待ち行列1704から取得されるとき、このSESTインデックスに関する情報は、SEST情報テーブルから読み込まれる。キャッシュ制御部1706は、どのディスクブロックがこの転送のために要求されるかを計算し、ディスクブロック番号とアクセスされるべきディスク装置のハッシュルックアップを用いるキャッシュバッファ位置にこれを翻訳する。もし、この転送のために要求される書き込みバッファメモリ1736のいずれのバッファも現在他の転送によって使用されているならば、SESTインデックスは、他の転送の完了を持つ顕著な要求待ち行列1710に入れられる。さもなければ、もし、これがREAD転送であり、要求されたバッファのすべてがキャッシュにあるならば、SESTインデックスは、蓄えられたREAD待ち行列1712に入れられる。さもなければ、SESTインデックスは、記憶装置要求待ち行列1714に書き込まれる。このアルゴリズムに対する可能な高度化は、バッファが現在蓄えられているとすれば、同一のバッファの多数のREADが進行中であることを

40

50

可能にするべきである。

【0111】

エントリが使用されたバッファ待ち行列1708から取得されるとき、チェックは、いずれの要求が利用可能になるためにこのバッファを待っているか否かについてなされる。これは、最も古い要求で始まる顕著な要求待ち行列を探すことによってなされる。もし、利用可能になるためにこのバッファを待っていた要求が見出されるならば、バッファは、その要求に割り当てられる。もし、その要求がこの転送のために要求されたすべてのバッファを有するならば、SESTインデックスは、記憶装置要求待ち行列1714に書き込まれ、この要求は、顕著な要求待ち行列1710から取り除かれる。さもなければ、その要求は、顕著な要求待ち行列1710に残されたままにされる。

10

【0112】

交換要求(ERQ)発生器1716は、記憶装置要求待ち行列1714と部分的なWRITE待ち行列1718からのエントリを取得する。SESTインデックスがいずれかの待ち行列から読み出されるとき、このSESTインデックスに関する情報は、SEST情報テーブルから読み出される。もし、それがREAD転送ならば、ファイバチャネル1/0制御部1720は、適切なディスクからのデータを読み込むように構成される。もし、それがWRITE転送であり、書き込まれるべきすべてのバッファが現在蓄えられているか、あるいは、書き込まれるべきデータが書き込まれるべきバッファを完全に満たしているならば、WRITEは、すぐに実行され得る。書き込まれるべきデータは、WRITEバッファメモリ1736からキャッシュバッファの適切な領域にコピーされる。ファイバチャネルI/O制御部1720は、適切なディスクにデータを書き込むように構成される。さもなければ、図17Aに関連して上述したように、我々は、適切なディスクからの要求されたデータのREADを始めるためにWRITEをする前に、ディスクからのREADをする必要がある。

20

【0113】

IMQプロセッサ1722は、入ってくるメッセージ待ち行列1724からのメッセージを取得する。これは、ファイバチャネルI/O制御部1720が完了した変換、または問題に遭遇した変換の待ち行列である。もし、ファイバチャネル変換に問題があったならば、IMQプロセッサ1722は、適切なエラー回復を可能にするために、プロセッサメッセージ待ち行列1726を介してプロセッサにメッセージを送る。もし、変換が許容できたならば、SEST情報は、このSESTインデックスのために読み出される。もし、この変換がWRITE変換の始まりにおけるREAD変換であったならば、SESTインデックスは、部分的なWRITE待ち行列1718に書き込まれる。さもなければ、それは、記憶装置承認待ち行列1728に書き込まれる。

30

【0114】

図17Aに関連して上述されるように、プロト-SCSI承認発生器1730は、プロト-SCSI要求を生成する責任を負う。再び、プロト-SCSI応答を生成することができ、それぞれがSESTインデックスを供給する3つの可能な情報源がある。

【0115】

プロセッサ承認待ち行列1732は、エラーを発生し、エラーが解決されるとすぐにハードウェアに戻ってプロセッサ1738によって解決されなければならない要求を送るために、プロセッサ1738によって用いられる。記憶装置承認待ち行列1728は、通常完了したファイバチャネル要求に送り戻すために用いられる。

40

【0116】

これらの待ち行列のいずれかにエントリがあるとき、SESTインデックスは、読み出される。このインデックスのためのSEST情報はそれから読み込まれる。転送のために、ステータス情報に沿って、プロト-SCSI要求が承認を関連付けるのを可能にする同定は、プロト-SCSI承認インターフェース1734の向こう側に戻される。READのために、読まれたデータは、また、プロト-SCSI承認インターフェース1734の向こう側へ戻される。

【0117】

50

プロト - S C S I 転送が完了されるとすぐに、この転送に関連したすべてのバッファのアドレスは、使用されたバッファ待ち行列 1708 に書き込まれる。この転送に用いられるあらゆる WRITE バッファメモリは、また、フリー WRITE バッファメモリのプールに戻される。

【0118】

図 18 は、ファイル記憶装置を有するコンピュータシステムに関連して、図 3 に示されるような記憶装置モジュールの使用を示すブロック図である。ここで、記憶装置モジュール 1801 は、マイクロプロセッサ 1802、メモリ 1803、ビデオシステム 1805 のような周辺機器、並びに、記憶装置 1809、1810、及び 1811 を含むコンピュータシステムのためのファイバチャネルホストバスアダプタ及びドライバとして行動する。記憶装置モジュール 1801 は、P C I バス 1807 上で P C I ブリッジ 1804 を介してマイクロプロセッサ 1802 とコンピュータメモリ 1803 に接続される。記憶装置モジュール 1801 は、P C I バスから要求を受信し、図 17A 及び 17B に関連して上述された方法でその要求を処理する。記憶装置モジュール 1801 は、記憶装置アクセスインターフェース 1808 を介して記憶装置 1809、1810、及び 1811 にアクセスする。

【0119】

図 19 は、本発明の実施の形態、特に、複数のネットワークサブシステムとサービスサブシステムが連続するサブシステム及び / 又はモジュールのポート間で通信を確立するために拡張スイッチを利用して使用される一実施の形態の拡大縮小可能性を示すブロック図である。余分なネットワーク接続がユニットの帯域幅能力を増加し、より多くの記憶要素をサポートすることを可能にするために、この実施の形態では、拡張スイッチ 1901、1902、1903 は、多くのモジュールとともに相互接続するために用いられる。拡張スイッチは、拡張スイッチの一面におけるモジュールから他面におけるあらゆるモジュールまでのあらゆる接続の経路を定める。拡張スイッチは、非ブロック化であり、多数の入力を取得し、特定の接続のために最良なルートを決める情報処理能力を持つ拡張スイッチ制御モジュールによって制御されてもよい。

【0120】

図 19 の実施の形態では、示された全システムは、ネットワークサブシステム 1904 と類似のサブシステム 1908 及び 1912 を含むコラム 1921 に示される複数のネットワークサブシステムを利用する。それは、ここでは (コラム 1922 の) ファイルアクセスモジュール、(コラム 1923 の) ファイルシステムモジュール、及び (コラム 1924 の) 記憶装置モジュールの組み合わせとして実行される複数のサービスサブシステムである。モジュールの各コラム間 (及びネットワークサブシステムコラムとファイルアクセスモジュールコラムの間) にスイッチ配置があり、ファイルアクセスプロトコル拡張スイッチ 1901、記憶装置アクセス拡張スイッチ 1902、及びプロト - S C S I プロトコル拡張スイッチ 1903 として実行される。ファイルアクセスプロトコルレベルでは、拡張スイッチ 1901 は、各ファイルアクセスモジュール 1905 の存在する家業負荷を含む基準に依存して、ネットワークサブシステム 1904 から特定のファイルアクセスモジュール 1905 に入ってくるネットワーク接続を動的に割り当てる。

【0121】

記憶装置アクセスプロトコルレベルでは、拡張スイッチ 1902 は、ファイルシステムモジュール 1906 の存在する作業負荷を含む基準に依存して、ファイルアクセスモジュール 1905 から特定のファイルシステムモジュール 1906 に入ってくるファイルアクセス接続を動的に割り当てる。

【0122】

プロト - S C S I プロトコルレベルでは、拡張スイッチ 1903 は、記憶要素の物理的位置を含む基準に依存して、特定の記憶装置モジュール 1907 へ入ってくるファイルシステム接続を動的に割り当てる。

【0123】

10

20

30

40

50

その代わりに、アイテム 1901、1902、及び 1903 は、バスとして実装されてもよい。この場合、入力信号を受けるコラム内の各モジュールは、信号の重複処理を避けるためにコラム内の他のモジュールと通信し、それによって、他の信号を扱うために他のモジュールを自由にする。アイテム 1901、1902、及び 1903 がバスまたはスイッチのいずれとして実現されるときも、ファイル要求への応答が含まれるとき、対応する要求からの適切なヘッダ情報が応答ヘッダの便利なフォーマット化を許すために利用可能であるように、システムを通して信号処理バスを追跡することは本発明の範囲内である。

【図面の簡単な説明】

【図 1】 図 1 は、ファイルサーバーやウェブサーバーのようなネットワークサービスを提供するために配置された本発明の一実施の形態の概略表示である。

10

【図 2】 図 2 は、図 1 に示された実施の形態のブロック図である。

【図 3】 図 3 は、ファイルサーバーとして配置される実施の形態のブロック図である。

【図 4】 図 4 は、ウェブサーバーとして配置される実施の形態のブロック図である。

【図 5】 図 5 は、図 2 ~ 4 の実施の形態のネットワークサブシステムである。

【図 6】 図 6 は、図 5 のネットワークサブシステムのブロック図である。

【図 7】 図 7 は、図 6 のネットワークサブシステムの受信モジュールのブロック図である。

【図 8】 図 8 は、図 6 のネットワークサブシステムの送信モジュールのブロック図である。

20

【図 9】 図 9 は、ワークステーション又はサーバーのようなネットワークノードで使用するネットワークインターフェースアダプタとして、図 5 のネットワークサブシステムの使用を例証するブロック図である。

【図 10】 図 10 は、図 3 に例証されるような一実施の形態で使用する図 3 の S M B サービスモジュール 33 とファイルシステムモジュール 34 のハードウェアで実行される組み合わせのブロック図である。

【図 11】 図 11 は、図 3 に例証されるような一実施の形態で使用する図 3 の S M B サービスモジュール 33 とファイルシステムモジュール 34 のハードウェア加速の組み合わせのブロック図である。

【図 12】 図 12 A は、それぞれ図 3 又は図 4 のアイテム 33 又は 34 めようなハードウェアで実行されるサービルモジュールのブロック図である。図 12 B は、それぞれ図 3 又は図 4 のアイテム 34 又は 44 のようなハードウェアで実行されるファイルモジュールのブロック図である。図 12 C は、組み合わせられたサービスモジュールとファイルモジュールを提供する図 10 のハードウェアで実行されるサービスサブシステムの詳細なブロック図である。

30

【図 13】 図 13 は、図 11 のハードウェア加速サービスサブシステムの詳細なブロック図である。

【図 14】 図 14 は、多数のサービス要求を多数のスレッドとして扱うために、ソフトウェアで実行される典型的な先行技術アプローチを表すフローチャートである。

【図 15】 図 15 は、図 2 のサービスサブシステムと、例えば、図 12 及び 13 の実施の形態とを関連して使用するために、多数のサービス要求の処理を示すフローチャートである。

40

【図 16】 図 16 は、ファイル記憶装置を有するコンピュータシステムに関連して、図 3 に示されるようなファイルシステムモジュールの使用を示すブロック図である。

【図 17】 図 17 A は、図 3 の記憶装置モジュールにおけるデータフローのブロック図である。図 17 B は、図 3 の記憶装置モジュールにおける制御フローのブロック図である。

【図 18】 図 18 は、ファイル記憶装置を有するコンピュータシステムに関連して、図 3 に示されるような記憶装置モジュールの使用を示すブロック図である。

【図 19】 図 19 は、本発明の実施の形態、特に、複数のネットワークサブシステムと

50

サービスサブシステムが連続するサブシステム及び／又はモジュールのポート間で通信するための拡張スイッチを利用して使用される一実地の形態の拡大縮小可能性を示すブロック図である。

【図 1】

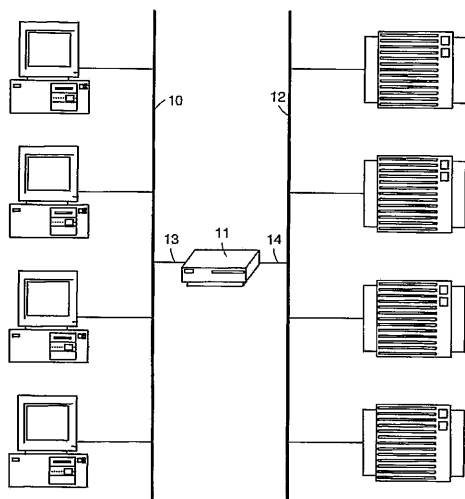


FIG. 1

【図 2】

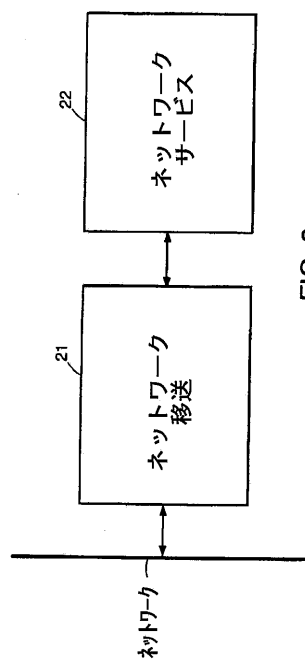


FIG. 2

【図 3】

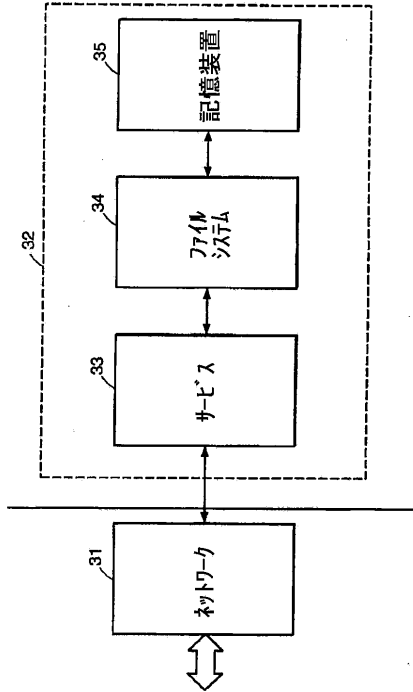


FIG. 3

【図 4】

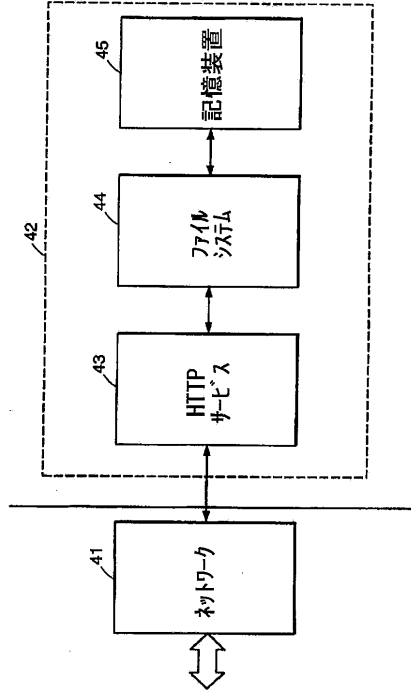


FIG. 4

【図 5】

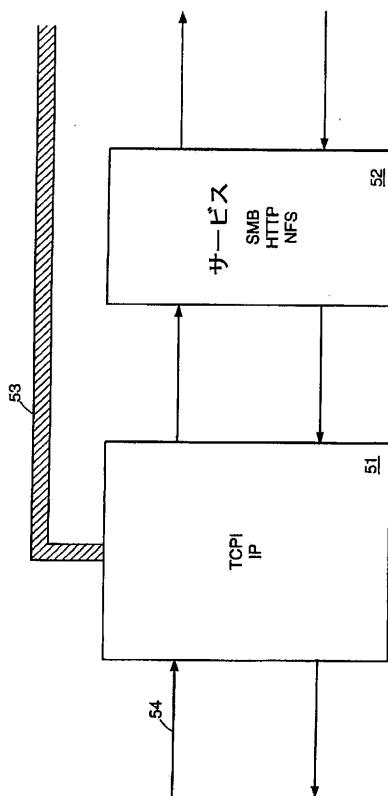


FIG. 5

【図 6】

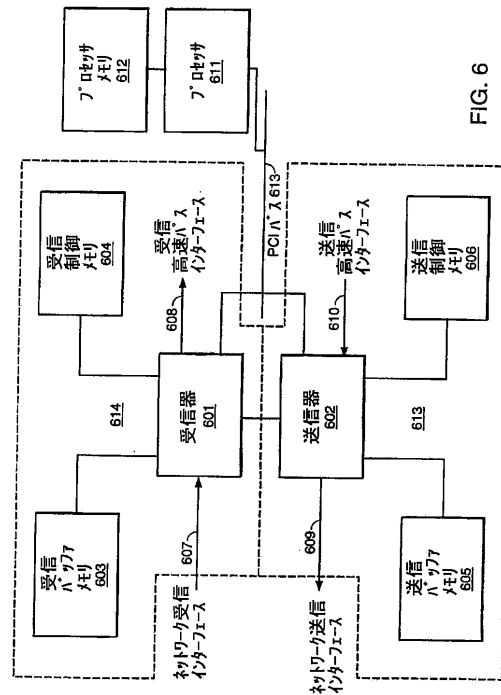


FIG. 6

【図 7】

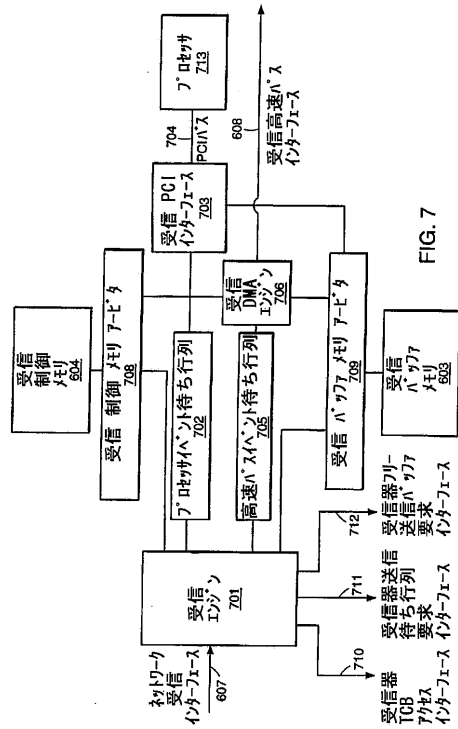


FIG. 7

【図 8】

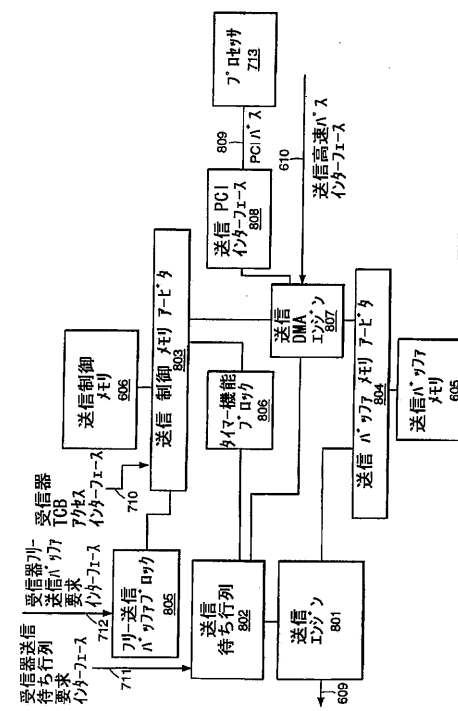


FIG. 8

【図 9】

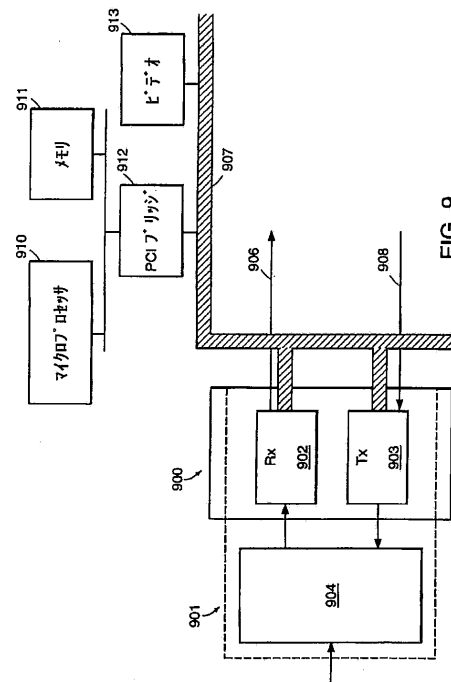


FIG. 9

【図 10】

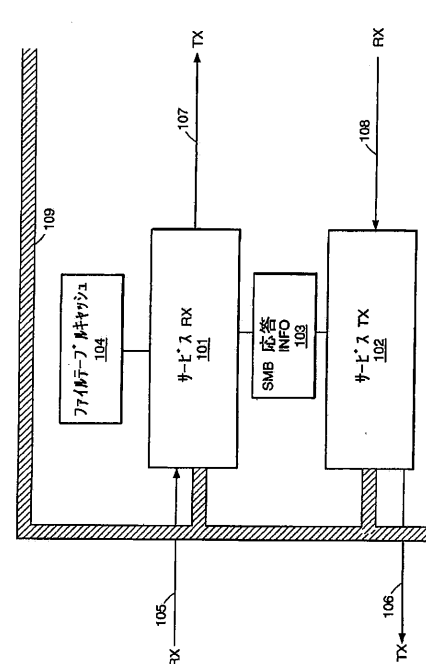
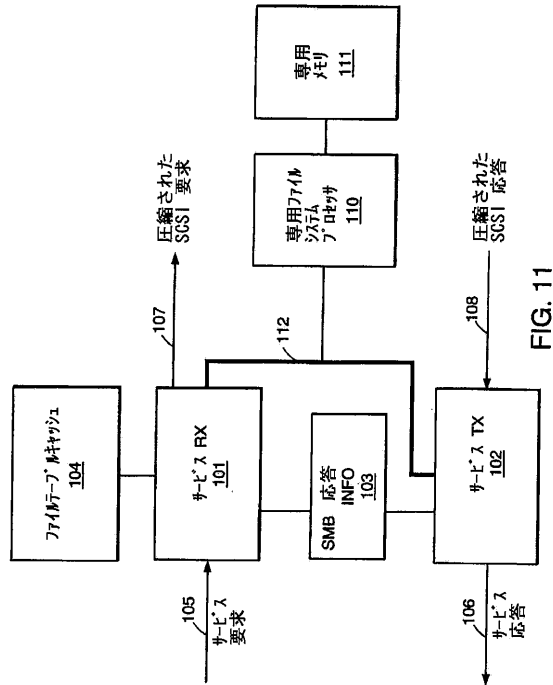
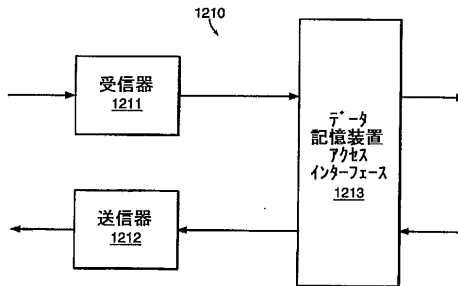
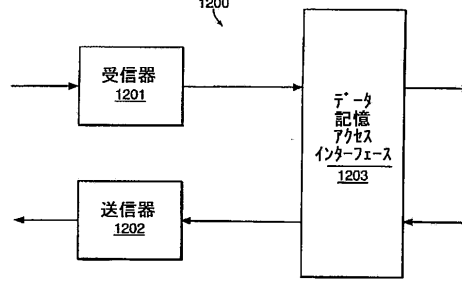


FIG. 10

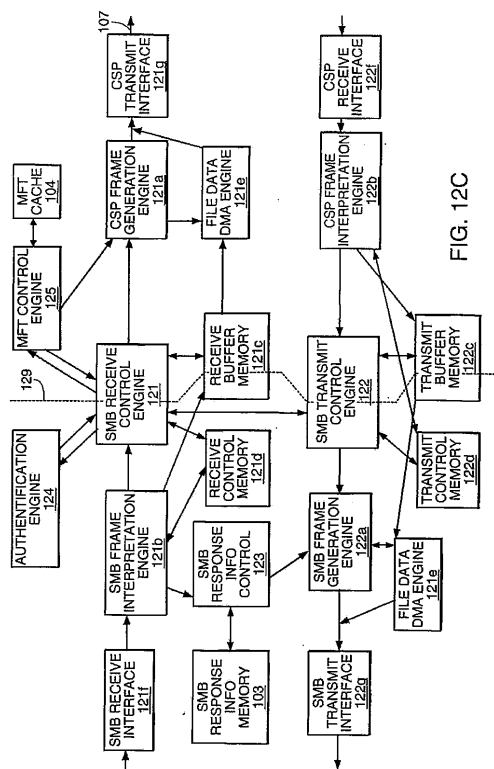
【図 11】



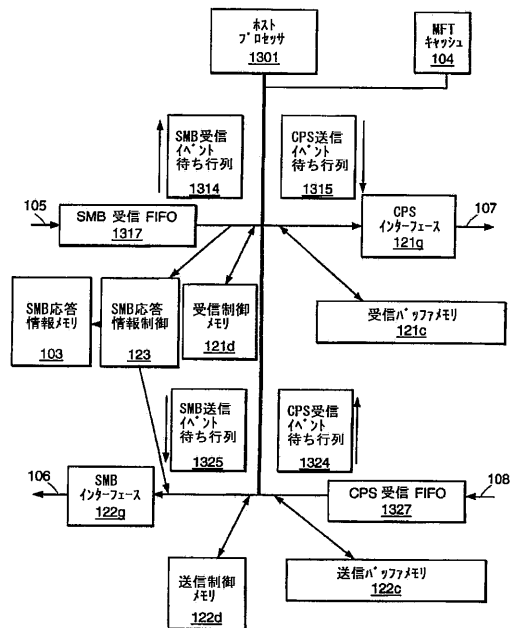
【図 12 A - B】



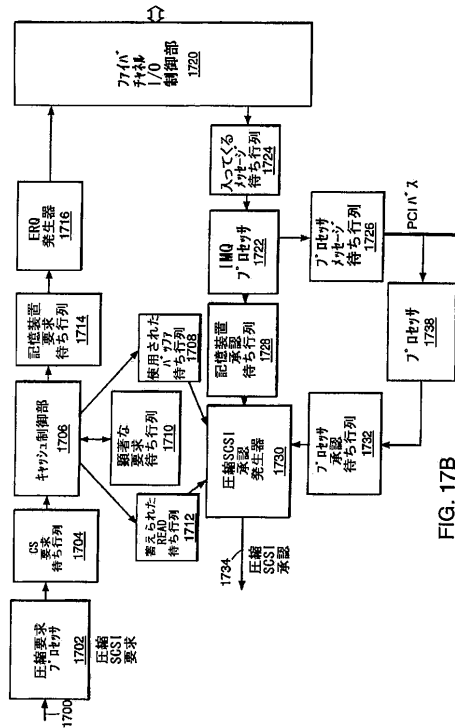
【図 12 C】



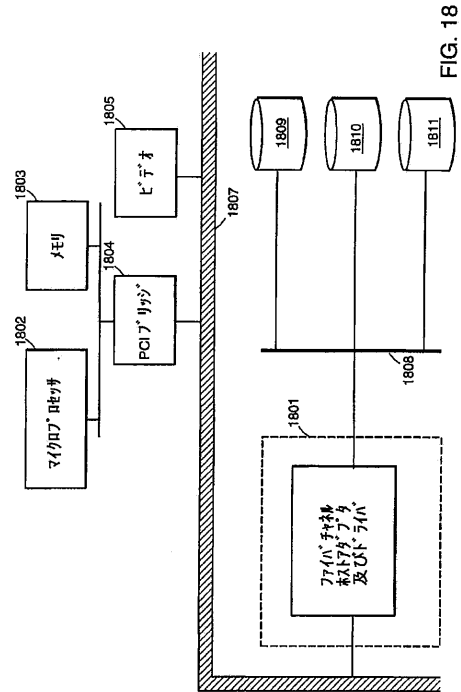
【図 13】



【 図 1 7 B 】



【 図 1 8 】



【 圖 1 9 】

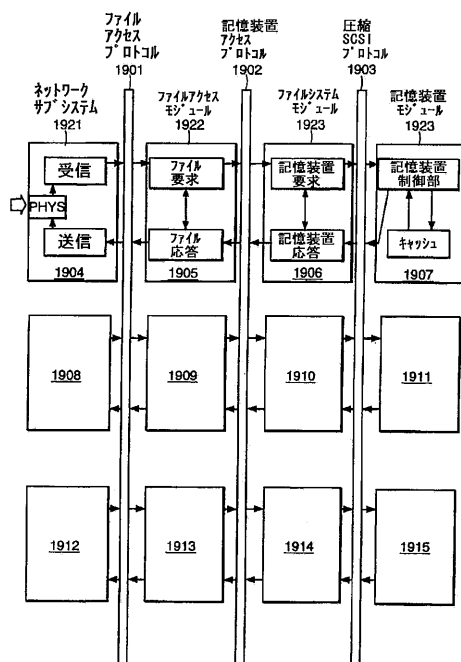


FIG. 19

フロントページの続き

- (72)発明者 ウィリス、トレバー
英国、バッキンガムシャー・エイチピー 22・6 エイチジー、エイリズベリー、イクフィールド・クロース 9
- (72)発明者 ベンハム、サイモン
英国、バークシャー・アールジー 12・2 エヌエイ、ブラックネル、ロックバイ・クロース 4
- (72)発明者 クーパー、マイケル
英国、バークシャー・アールジー 1・3 ピーキュー、リーディング、リバプール・ロード 102
- (72)発明者 マイヤー、ジョナサン
英国、サリー・ジーユー 10・3 エヌジー、ファーマム、ホワイト・ローズ・レーン 2エイ
- (72)発明者 アストン、クリストファー・ジェイ
英国、バッキンガムシャー・エイチピー 13・6 ジェイアール、ウィカム、ザ・クレセント 40
- (72)発明者 ウィンフィールド、ジョン
英国、バークシャー・アールジー 2・0 ビーエヌ、リーディング、エルガー・ロード 126

合議体

審判長 水野 恵雄

審判官 山田 正文

審判官 稲葉 和生

- (56)参考文献 特表平5 - 502525 (JP, A)
特開平3 - 273350 (JP, A)
特開平9 - 26970 (JP, A)
特開平11 - 203308 (JP, A)
特開平10 - 232788 (JP, A)

- (58)調査した分野(Int.Cl., DB名)

G06F13/00