

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4262888号  
(P4262888)

(45) 発行日 平成21年5月13日(2009.5.13)

(24) 登録日 平成21年2月20日(2009.2.20)

(51) Int.Cl.

F I

G 0 6 F 9/50 (2006.01)

G 0 6 F 9/46 4 6 5 A

請求項の数 28 (全 23 頁)

(21) 出願番号	特願2000-553887 (P2000-553887)	(73) 特許権者	500046438
(86) (22) 出願日	平成11年5月11日 (1999.5.11)		マイクロソフト コーポレーション
(65) 公表番号	特表2002-517855 (P2002-517855A)		アメリカ合衆国 ワシントン州 9805
(43) 公表日	平成14年6月18日 (2002.6.18)		2-6399 レッドモンド ワン マイ
(86) 国際出願番号	PCT/US1999/010273		クロソフト ウェイ
(87) 国際公開番号	W01999/064952	(74) 代理人	100089705
(87) 国際公開日	平成11年12月16日 (1999.12.16)		弁理士 社本 一夫
審査請求日	平成18年4月27日 (2006.4.27)	(74) 代理人	100071124
(31) 優先権主張番号	09/097, 169		弁理士 今井 庄亮
(32) 優先日	平成10年6月12日 (1998.6.12)	(74) 代理人	100076691
(33) 優先権主張国	米国 (US)		弁理士 増井 忠武
		(74) 代理人	100075270
			弁理士 小林 泰
		(74) 代理人	100096013
			弁理士 富田 博行

最終頁に続く

(54) 【発明の名称】 処理タスクをソフトウェアからハードウェアにオフロードする方法およびコンピュータ・プログラム製品

(57) 【特許請求の範囲】

【請求項 1】

コンピュータと、少なくとも1つのソフトウェア・コンポーネントと、少なくとも1つの周辺ハードウェア・デバイスとを有するコンピュータ・システム環境において、パケット毎に、前記コンピュータ・システムのその時々必要性に応じて、前記コンピュータから前記周辺ハードウェア・デバイスに動作タスクを動的にオフロードすることによって、前記コンピュータのプロセッサ資源を解放し、前記コンピュータ・システム全体の効率を向上させる方法であって、

前記コンピュータが、前記周辺ハードウェア・デバイスに問い合わせ、前記周辺ハードウェア・デバイスのタスク・オフロード能力を判定するステップと、

前記コンピュータが、前記周辺ハードウェア・デバイスの選択したタスク・オフロード能力をイネーブルして、1つ以上のパケットが必要とするイネーブルされたタスク・オフロード能力を提供するステップと、

前記コンピュータによってあるパケットに対して実行する動作タスクが前記周辺ハードウェア・デバイス上のイネーブルされたタスク・オフロード能力に対応する場合、前記コンピュータ・システムのその時々必要性に応じて、前記コンピュータが、データ・パケットを前記周辺ハードウェア・デバイスに送り、前記周辺ハードウェア・デバイスが前記動作タスクを実行することを示すことにより、前記コンピュータから前記周辺ハードウェア・デバイスに前記動作タスクを選択的にオフロードする作用を行なうステップと、

前記周辺ハードウェア・デバイスが、前記オフロードされた動作タスクを実行するステ

10

20

ップと、  
を含む方法。

【請求項 2】

請求項 1 に記載した前記ステップおよび作用を実行するためのコンピュータ実行可能命令を有するコンピュータ読取可能媒体。

【請求項 3】

請求項 1 に記載の方法において、前記周辺ハードウェア・デバイスが、前記コンピュータ・システムに動作的に接続してあるネットワーク・インターフェース・カード (NIC) であること、を特徴とする方法。

【請求項 4】

請求項 1 に記載の方法において、前記ソフトウェア・コンポーネントが、レイヤ状ネットワーク・モデルにおいて実行するネットワーク・ソフトウェア・アプリケーションであること、を特徴とする方法。

【請求項 5】

請求項 1 に記載の方法において、前記周辺ハードウェア・デバイスのタスク・オフロード能力を問い合わせることは、前記周辺ハードウェア・デバイスに関連する少なくとも 1 つのタスク・オフロード・バッファ内に格納してあるタスク・データを読み取ることにより行い、該タスク・データが前記周辺ハードウェア・デバイスの特定のタスク・オフロード能力を示すこと、を特徴とする方法。

【請求項 6】

請求項 5 に記載の方法において、前記周辺ハードウェア・デバイスには、複数のタスク・オフロード・バッファを関連付けてあり、これによって、当該周辺ハードウェア・デバイスの複数のタスク・オフロード能力を定めること、を特徴とする方法。

【請求項 7】

請求項 1 に記載の方法において、前記周辺ハードウェア・デバイスの選択したタスク・オフロード能力をイネーブルすることは、当該周辺ハードウェア・デバイスに関連するタスク・オフロード・バッファ内に、少なくとも 1 つのフラグ・インディケータをセットすることにより行う、を特徴とする方法。

【請求項 8】

請求項 1 に記載の方法において、前記パケットが、ネットワーク・データおよびパケット拡張データを含むネットワーク・データ・パケットであり、前記パケット拡張データが、前記周辺ハードウェア・デバイスが実行する少なくとも 1 つの動作タスクを示す少なくとも 1 つのデータ・フィールドを含むこと、を特徴とする方法。

【請求項 9】

請求項 8 に記載の方法において、前記周辺ハードウェア・デバイスが、ネットワーク・インターフェース・カード (NIC) であること、を特徴とする方法。

【請求項 10】

請求項 1 に記載の方法において、チェックサム動作、暗号化動作、メッセージ・ダイジェスト計算動作、TCP 区分動作、および解読動作である動作タスクのうちの 1 つ以上から前記動作タスクを選択すること、を特徴とする方法。

【請求項 11】

請求項 1 に記載の方法において、前記ソフトウェア・コンポーネントがトランスポート・プロトコル・ドライバ・ソフトウェア・コンポーネントを備え、前記コンピュータ・システムがネットワーク・インターフェース・カード (NIC) デバイス・ドライバを含み、前記周辺ハードウェア・デバイスが、前記コンピュータ・システムに接続した対応のネットワーク・インターフェース・カード (NIC) を備え、

前記周辺ハードウェア・デバイスに問い合わせる当該周辺ハードウェア・デバイスのタスク・オフロード能力を判定する前記ステップが、

前記 NIC のタスク・オフロード能力を識別するタスク・データを含むデータ構造を作成する作用と、

10

20

30

40

50

前記データ構造内に収容したタスク・データを問い合わせ、前記NICのタスク・オフロード能力を識別する作用と、  
を含み、

前記動作タスクを選択的にオフロードする前記作用が、タスク・オフロード・データを前記パケットに添付する作用を含み、

前記オフロードされた動作タスクを実行する前記ステップが、前記添付したタスク・オフロード・データにしたがって、前記NICにおいて前記少なくとも1つのオフロードされたタスクを実行する作用を含むこと、を特徴とする方法。

【請求項12】

請求項11に記載のステップおよび作用を実行するためのコンピュータ実行可能命令を有するコンピュータ読取可能媒体。

10

【請求項13】

請求項11記載の方法において、前記データ構造が、  
タスク・オフロード能力のタイプを表すデータを収容する第1データ・フィールドと、  
前記第1データ・フィールドにおいて定めたタスク・オフロードを実行するために必要なデータを収容する第2データ・フィールドと、  
を備えること、を特徴とする方法。

【請求項14】

請求項11記載の方法において、前記データ構造内に少なくとも1つのデータ・フィールドを書き込むことによって、前記NICを選択的にイネーブルし、前記識別したタスク・オフロード能力の少なくとも1つを実行すること、を特徴とする方法。

20

【請求項15】

請求項11記載の方法において、前記パケットに添付した前記タスク・オフロード・データが、前記NICが前記パケットを受信したときに、当該NICが実行するオフロードされたタスクのタイプを表す第1データ・フィールドを備えること、を特徴とする方法。

【請求項16】

請求項15記載の方法において、前記タスク・オフロード・データが、更に、メモリ位置への少なくとも1つのポインタを収容する第2データ・フィールドを備え、前記メモリ位置が、前記タスク・オフロード・データの第1データ・フィールドにおいて指定した前記オフロードされたタスクを実行するために必要なデータを収容すること、を特徴とする方法。

30

【請求項17】

コンピュータを備えたコンピュータ・システムのその時々必要性に応じて、パケット毎に、計算タスクを動的にオフロードして、前記コンピュータに接続してあるネットワーク・インターフェース・カード(NIC)において実行されるようにすることにより、前記コンピュータのプロセッサ資源を解放し、前記コンピュータ・システム全体の効率を向上させるコンピュータ実行可能命令を有するコンピュータ読取可能媒体であって、

前記コンピュータ上で実行するアプリケーションに、前記NICのあらゆるタスク・オフロード処理能力について通知するステップと、

前記コンピュータ上で実行するアプリケーションから前記NICに向けたネットワーク・データ・パケット毎に、該データ・パケットに添付するパケット拡張データ構造を検査するステップであって、前記パケット拡張データ構造が、前記NICにオフロードする1つ以上の特定の動作タスクの識別に対する情報を収容する、ステップと、

40

前記ネットワーク・データ・パケットに添付する前記パケット拡張データ構造内に収容してあるデータにしたがって、前記NICに動作タスクを実行させるステップと、  
を含むステップを実行させるコンピュータ実行可能命令を更に含むコンピュータ読取可能媒体。

【請求項18】

請求項17記載のコンピュータ読取可能媒体において、前記通知するステップを実行する前記コンピュータ実行可能命令が、前記NICのタスク・オフロード能力のタイプを表

50

すデータを収容する第1データ・フィールドを備えるタスク・オフロード・バッファ・データ構造を作成するステップを含むこと、を特徴とするコンピュータ読取可能媒体。

【請求項19】

請求項18記載のコンピュータ読取可能媒体において、前記タスク・オフロード・バッファ・データ構造が、更に、前記第1データ・フィールドにおいて定めた前記タスク・オフロードを実行するために必要なデータを収容する少なくとも1つの追加データ・フィールドを備えること、を特徴とするコンピュータ読取可能媒体。

【請求項20】

請求項17記載のコンピュータ読取可能媒体において、前記ネットワーク・データ・パケットに添付する前記パケット拡張データ構造が、前記NICが実行する動作タスクのタイプを表す第1データ・フィールドを備えること、を特徴とするコンピュータ読取可能媒体。

10

【請求項21】

請求項20記載のコンピュータ読取可能媒体において、前記パケット拡張データ構造が、更に、メモリ位置への少なくとも1つのポインタを収容する少なくとも1つの追加データ・フィールドを備え、前記メモリ位置が、前記パケット拡張データ構造の第1データ・フィールドにおいて指定した動作タスクを実行するために必要なデータを収容すること、を特徴とするコンピュータ読取可能媒体。

【請求項22】

コンピュータと、少なくとも1つのソフトウェア・コンポーネントと、少なくとも1つの周辺ハードウェア・デバイスとを有するコンピュータ・システム環境において、パケット毎に、前記コンピュータ・システムのその時々の必要性に応じて、前記コンピュータから前記周辺ハードウェア・デバイスに動作タスクを動的にオフロードすることによって、前記コンピュータのプロセッサ資源を解放し、前記コンピュータ・システム全体の効率を向上させる方法であって、

20

前記周辺ハードウェア・デバイスに関連した少なくとも1つのタスク・オフロード・バッファに格納してあるタスク・データを読み取る作用であって、前記タスク・データが、前記周辺ハードウェア・デバイスの特定のタスク・オフロード能力を示す、作用と、

前記少なくとも1つのタスク・オフロード・バッファ内に収容してあるタスク・データを問い合わせ、前記周辺ハードウェア・デバイスのタスク・オフロード能力を識別する作用と、

30

前記周辺ハードウェア・デバイスに関連したタスク・オフロード・バッファ内に、少なくとも1つのフラグ・インディケータをセットする作用であって、前記フラグ・インディケータが1つ以上のパケットに必要な関連するタスク・オフロード能力のイネーブルを表す、作用と、

前記コンピュータによってパケットに対して実行する動作タスクが前記周辺ハードウェア・デバイス上のイネーブルされたタスク・オフロード能力に対応する場合、前記コンピュータ・システムのその時々の必要性に応じて、前記コンピュータが、タスク・オフロード・データをデータ・パケットに添付し、該データ・パケットを前記周辺ハードウェア・デバイスに送ることによって、前記コンピュータから前記周辺ハードウェア・デバイスに

40

前記動作タスクを選択的にオフロードする作用を実行し、  
前記添付したタスク・オフロード・データにしたがって、前記周辺ハードウェア・デバイスにおいて前記少なくとも1つのオフロードされたタスクを実行する作用と、  
を含む方法。

【請求項23】

請求項22に記載の作用を実行するためのコンピュータ実行可能命令を有するコンピュータ読取可能媒体。

【請求項24】

請求項22記載の方法において、前記周辺ハードウェア・デバイスが、前記コンピュータ・システムに動作的に接続してあるネットワーク・インターフェース・カード(NIC

50

）であること、を特徴とする方法。

【請求項 2 5】

請求項 2 2 記載の方法において、前記ソフトウェア・コンポーネントが、レイヤ状ネットワーク・モデルにおいて実行するネットワーク・ソフトウェア・アプリケーションであること、を特徴とする方法。

【請求項 2 6】

請求項 2 2 記載の方法において、前記周辺ハードウェア・デバイスには、複数のタスク・オフロード・バッファを関連付けてあり、これによって前記周辺ハードウェア・デバイスの複数のタスク・オフロード能力を定めること、を特徴とする方法。

【請求項 2 7】

請求項 2 2 記載の方法において、前記周辺ハードウェア・デバイスがネットワーク・インターフェース・カード（NIC）であること、を特徴とする方法。

【請求項 2 8】

請求項 2 2 記載の方法において、チェックサム動作、暗号化動作、メッセージ・ダイジェスト計算動作、TCP 区分動作、および解読動作である動作タスクのうちの 1 つ以上から前記動作タスクを選択すること、を特徴とする方法。

【発明の詳細な説明】

【0001】

（発明の背景）

（1．発明の分野）

本発明は、一般的に、コンピュータ・システムの効率、速度および／またはスループットを向上させる方法に関する。更に特定すれば、本発明は、典型的にホスト・プロセッサがソフトウェアで実行する計算タスクを特定のハードウェア・コンポーネントにオフロード（offload）することにより、ホスト・コンピュータ資源を開放し、コンピュータ・システムの全体的効率を高める方法に関する。

【0002】

（2．従来技術の状態）

機能的コンピュータ・システムは、概略的に、3つの基本コンポーネントで構成されている。第1のコンポーネントは、ホスト・コンピュータおよびそれに付随する周辺ハードウェア・コンポーネントである。ホスト・コンピュータは、典型的に、中央演算装置（CPU）を含み、バスを介して、例えば、RAMまたはROMのようなシステム・メモリと相互接続してある。また、システムは、必要な機能性に応じて、磁気または光のディスク記憶装置、キーボードまたはその他の入力デバイス、ディスプレイまたはその他の出力デバイス、ならびにモデムおよび／またはネットワーク・インターフェース・カード（NIC）のような通信機器といった、多数の周辺ハードウェア・デバイスも含む。別の機能的コンピュータ・コンポーネントに、アプリケーション・ソフトウェアがある。このようなソフトウェアは、周知のワード・プロセッサ・アプリケーション、スプレッド・シート・アプリケーション、データベース・アプリケーション、通信およびネットワーク・アプリケーション等を含む。

【0003】

最近の機能的コンピュータ・システムの最後のコンポーネントは、オペレーティング・システムである。コンピュータのオペレーティング・システムは、ユーザにアプリケーション・プログラムの実行を開始させるといような、多くの機能を実行する。加えて、最新のオペレーティング・システムは、アプリケーション・ソフトウェアとホスト・コンピュータおよびその周辺ハードウェアとの間のインターフェースも備えている。したがって、以前にはアプリケーション・プログラムが直接コンピュータ・システムのハードウェアにアクセスすることが当たり前であったが、最新のオペレーティング・システムは標準化された一貫性のあるインターフェースを備え、ユーザ・アプリケーションに、標準化された方法でコンピュータ・ハードウェア周辺機器とインターフェースし、これにアクセスさせるようにしている。一貫性のあるインターフェースを備えるために、オペレーティング・

10

20

30

40

50

システムのアーキテクチャの設計において、実際のハードウェア周辺機器およびアプリケーション・プログラム間にいくつものソフトウェア・レイヤが存在し得るようにする場合が増えている。例えば、アプリケーションは、オペレーティング・システムにコールすることができる。一方、オペレーティング・システムは、ハードウェア・デバイス・ドライバ・レイヤが与えるサービスを利用することができる。そのときには、デバイス・ドライバ・レイヤは、特定のハードウェア周辺機器と直接インターフェースする。このようなレイヤ型手法の主な利点は、他のレイヤに影響を及ぼすことなく、レイヤの追加または交換が可能なことである。

#### 【 0 0 0 4 】

認識されるように、このようなオペレーティング・システム、アプリケーション・ソフトウェア、ならびにネットワークおよび通信は、増々複雑化および精巧化し続けている。勿論、その結果、高機能化し有用性が高いコンピュータ・システムが得られる。しかしながら、この機能性の増大は、コストを伴わない訳ではない。オペレーティング・システムおよびソフトウェア・アプリケーションでは、機能が更に増加するに連れて、このようなシステム機能および／またはアプリケーションを実行するときのプロセッサ／CPUが実行しなければならない動作が増加し、その結果プロセッサのオーバーヘッド増大を招く場合が多い。この現象は、特に、ネットワーク通信型ソフトウェア・アプリケーションのような、特定の種類のアプリケーションと比較すると、特に明白となる。広帯域幅媒体が増々普及するに連れて、ネットワーク速度は、ホスト・コンピュータのCPUのプロセッサ速度やメモリ帯域幅に匹敵するかあるいはこれを凌駕することも多い。したがって、このようなネットワークを通じて効率的に通信するためには、ネットワーク接続するホスト・コンピュータのCPU利用およびメモリ帯域幅利用を極力抑えなければならない。

#### 【 0 0 0 5 】

加えて、ネットワーク・アプリケーションは、7層ISOモデルのような殆どが用いているレイヤ状アーキテクチャ、またはWindows NTオペレーティング・システムが用いているレイヤ・モデルにより、ホスト・プロセッサに更に負担をかける。公知のように、このようなモデルは、ネットワークへの物理的接続とエンド・ユーザ・アプリケーションとの間のデータ・フローを記述するために用いられる。ネットワーク・ケーブル上にデータ・ビットを置くというような、最も基本的な機能はボトム・レイヤにおいて行われ、一方アプリケーションの詳細に関与する機能はトップ・レイヤにおいて行われる。本質的に、各レイヤの目的は、次に高いレイヤにサービスを提供し、サービスが実際にどのように実施されるかに関する詳細から、高いレイヤを遮蔽することである。レイヤは、各レイヤが、ネットワークを介して通信している別のコンピュータ上にある同じレイヤと通信していると考えるように、抽象化されている。

#### 【 0 0 0 6 】

データ・パケットがレイヤ間を移動する際にこれに対して実行される種々の機能はソフトウェア集中的となり得るものであり、したがってかなりの量のCPUプロセッサおよびメモリ資源を要求する可能性があることは認められよう。例えば、Windows NTネットワーク・モデルでは、種々のレイヤにおいてパケットに対して行われるある機能は、パケット・チェックサムの算出および検証、データの暗号化および解読、メッセージ・ダイジェスト計算およびTCPセグメント化のように、極度にCPU集中的である。これらの機能の各々を実行する際、その結果としてのCPU／メモリに対する要求は、コンピュータ・システム全体のスループットおよび性能に大きく影響する可能性がある。

#### 【 0 0 0 7 】

ソフトウェア・アプリケーションおよびオペレーティング・システムの機能がコンピュータ・システム資源に対して行なう要求は増大するが、同時に、ネットワーク・インターフェース・カード(NIC)のような、多くのコンピュータ・ハードウェア周辺機器の能力、効率およびスループットも向上する。これらのコンピュータ・システム周辺機器には、多くの場合専用プロセッサおよびメモリが搭載されており、典型的に非常に精巧化され複雑なタスク、即ち、コンピュータ・システムのプロセッサがソフトウェアで実行するよう

なタスクを実行することができる。例えば、多くのNICは、チェックサム計算／検証、データ暗号化／解読、メッセージ・ダイジェスト計算、TCPセグメント化およびその他というような、適切なネットワーク・レイヤにおいてCPUがソフトウェアで実行するタスクを、独立して実行することができる。したがって、このようなCPU集約タスクを周辺のハードウェア・デバイスにオフロードすることは効果的である。これが可能であれば、ホスト・コンピュータにおけるプロセッサの利用度およびメモリ帯域幅の使用度が低下することにより、システム全体の効率、速度およびスループットが向上するであろう。

#### 【0008】

しかしながら、異なる周辺デバイスの処理能力は大きく異なる。したがって、コンピュータ・システム／オペレーティング・システムがこのような周辺デバイスの処理能力を識別し、次いで必要に応じて特定の処理タスクをデバイスに割り当てオフロードすることを可能にする効率的な方法が必要となる。また、その時々プロセッサの必要性に応じて動的にタスクを識別し割り当てることができれば望ましいであろう。これが可能であれば、コンピュータ・システムのプロセッサは、必要に応じてハードウェア周辺機器の能力を利用できるようになるであろう。

#### (発明の概要)

従来技術の現状における前述の問題は、本発明が見事に解決した。本発明は、以前はプロセッサ・ソフトウェア・レベルで行われていた機能およびタスクを、コンピュータ・システムに接続してある適切なハードウェア周辺機器にオフロードするシステムおよび方法に関する。本発明は、特に、コンピュータのCPUがソフトウェアで実行するタスクの多くを実行できる場合が多い、ネットワーク・インターフェース・カード(NIC)周辺デバイスに対する、タスクのオフロードにおいて特に有用である。

#### 【0009】

本発明の好適な実施形態の1つでは、例えば、オペレーティング・システム(OS)が、コンピュータ・システムに接続してあるあらゆるハードウェア周辺機器(NIC等)のデバイス・ドライバ(「MAC」ドライバと呼ばれる場合が多い)に「問い合わせ」することを可能にする、ソフトウェア実装方法およびプロトコルを提供する。種々のデバイス・ドライバは、各々、それぞれのハードウェア周辺機器処理能力を識別することによって応答する。この能力のことをここでは「タスク・オフロード能力」と呼ぶ。好適な実施形態では、一旦個々の周辺機器のタスク・オフロード能力を識別したなら、OSは選択した周辺機器に、OSによって潜在的に使用可能なあるタスクを実行させることができる。その後、OSは、その時々コンピュータ・システムの処理の必要性にしたがって、動的に、要求に応じて、先にイネーブルしたタスクまたは複数のタスクを周辺機器が実行するように要求する。

#### 【0010】

この概略的な発明概念は、他のアプリケーションまたはオペレーティング・システム環境にも適用可能であるが、本発明の実施形態は、ここでは、Windows NTのレイヤ状ネットワークング・モデルと共に実施し利用するものとして説明する。勿論、本発明は、ネットワーク通信を管理し制御するための同様の種類のアーキテクチャであれば、本質的にいずれとでも実施可能である。即ち、本発明は、典型的にネットワーク・パケットに対して例えば種々のネットワーク・レイヤにおいて行われ、典型的に専用CPUおよびメモリ資源を必要とするタスクまたは機能をオフロードする機能を提供する。これらのオフロードしたタスクは、任意に、ネットワークに実際の物理的通信チャネルを提供するハードウェア周辺機器、即ち、NICによって代わりに実行することができる。例えば、データ・パケットがそれぞれのネットワーク・レイヤを通過する際に、例えば、チェックサム計算／検証、暗号化／解読、メッセージ・ダイジェスト計算およびTCPセグメント化というようなデータ・パケットに対するCPU集約動作の一部を実行せずに、これらのタスクをオフロードし、NICハードウェアにおいて代わりに実行することができる。

#### 【0011】

本発明の好適な実施形態では、Windows NTレイヤ状ネットワークング・アーキ

10

20

30

40

50

テクチャにおいて、トランスポート・プロトコル・ドライバ、即ち、トランスポートを適切なプログラム方法によって実現し、コンピュータに接続してある対応のNIC（複数のNIC）と関連するデバイス・ドライバ（複数のデバイス・ドライバ）の各々に問い合わせができるようにする。問い合わせされた各デバイス・ドライバも、同様に、その具体的な処理能力、即ち、「タスク・オフロード」能力を識別することによって、応答可能なように実現する。好適な実施形態の1つでは、一旦個々の周辺デバイスのタスク・オフロード能力を識別したなら、トランスポートが、これら具体的な能力のどれをイネーブルするかについて設定を行なう。これは、本質的に、周辺デバイスに、続くデータ・パケットの送信および/または受信の間にどのような種類のタスクを実行することになっているのかについて知らせることになる。その後、トランスポートは、要求に応じて、周辺デバイスの機能をイネーブルし、これを利用することができる。好ましくは、イネーブルした機能は、ネットワーク・チャンネルに宛てられた実際のデータ・パケットに添付されている適切なデータによって呼び出す。このようにして、タスクを動的にオフロードすることができ、一度に1つより多いタスクをオフロードすることができる。

#### 【0012】

したがって、ネットワーク・パケットを特定の下位デバイス・ドライバ（例えば、Windows NT環境におけるMACサブレイヤに常駐するデバイス・ドライバ）に送る前に、トランスポートは、最初に、対応するNICの能力がどれくらいなのかについて判定を行なう。特定の機能または複数の機能が可能な場合、トランスポートは所望の機能をイネーブルする。続くパケットの送信中、トランスポートが特定のタスクをハードウェアにオフロードすることを望む場合、当該所望の機能（複数の機能）をNICハードウェアにおいてそのパケットに対して実行することを意味する情報を、パケットに動的に添付することができる。例えば、トランスポートは、データ・パケット内にデータ・フラグをセットすることによって、対応するデバイス・ドライバに、NICはチェックサムを計算し該当の発信パケットに添付することを通知する。このときには、対応するNIC上のハードウェア/ソフトウェアは、それ自体でこの特定のパケット処理に対処し、システムCPUからの介入や補助は全くない。したがって、システム・プロセッサは解放され、別の処理タスクを実行するので、システム全体の効率やスループットが向上する。

#### 【0013】

先に注記したように、好適な実施形態の1つでは、タスクを動的にダウンロードする。即ち、その時々コンピュータ・システムに対する必要性に応じて、NICの能力を選択的にパケット毎に用いることができる。更に、以前にはネットワーク・スタックの種々のレベルで実行していたタスクが、今や単一点、即ち、NIC自体で実行されるので、この手法は統合化および効率化を一層推し進め、システム全体のスループットを更に高めることになる。好ましくは、本発明の実施形態は、動作を「バッチ」する。即ち、多数のタスクを単一のNICにオフロードする機能を、トランスポートに設ける。例えば、単一のNICがチェックサム計算および暗号化双方を1つのパケットに対して実行することにより、同じ機能をそれぞれのソフトウェア・レイヤにおいてソフトウェアで実施する場合に必要な多数のCPUサイクルを不要とすることができる。

#### 【0014】

したがって、本発明は、コンピュータ・システム・プロセッサから、当該コンピュータ・システムに接続してあるハードウェア周辺機器に計算タスクをオフロードするシステムおよび方法を提供する。また、本発明は、個々の周辺機器の処理能力を識別するシステムおよび方法を提供する。本発明は、レイヤ状ネットワーク・アーキテクチャと共に効率的に使用することにより、ネットワーク内の種々のレイヤにおいて典型的に実行するタスクを、代わりに、適切なネットワーク・インターフェース・カード（NIC）にオフロードすることを可能にする、タスク・オフロード・システムおよび方法を提供する。本発明は、コンピュータ・プロセッサの現処理状態にしたがって、動的に、要求に基づいて計算タスクをオフロードすることができるシステムおよび方法を提供する。更に、本発明は、多数のタスクを共にバッチし、次いでNICのような単一の周辺デバイスにオフロードするこ

10

20

30

40

50



とができるシステムおよび方法を提供する。

【 0 0 1 5 】

本発明のその他の利点は、以下の説明に明記されており、その記載から部分的に明らかとなり、あるいは本発明の実施によって習得することができよう。本発明の利点は、添付した特許請求の範囲に特定の指摘した手段および組み合わせによって実現し、得ることができる。本発明のこれらおよびその他の特徴は、以下の説明および添付した特許請求の範囲から一層明らかとなり、以下に明記する本発明の実施によって習得することができよう。

【 0 0 1 6 】

( 図面の簡単な説明 )

先に引用した本発明の利点およびその他の利点が得られるようにするために、先に端的に説明した本発明の更に特定の説明を、添付図面に示す、その具体的な実施形態を参照しながら行なう。これらの図面は本発明の典型的な実施形態のみを図示するのであり、したがってその範囲を限定するものとはみなされないという理解の下で、添付図面を用いて、本発明を更に具体的かつ詳細に、記載し説明する。

【 0 0 1 7 】

( 好適な実施形態の詳細な説明 )

以下、本発明を説明するにあたって、本発明のシステムおよび方法を実現するために用いる実施形態の構造または処理を示す図を用いることにする。このように図面を用いて本発明を提示することは、その範囲を限定するものとして解釈されるべきでない。本発明は、パーソナル・コンピュータのようなホスト・コンピュータから、ネットワーク・インターフェース・カード ( N I C ) のような、コンピュータに接続してあるハードウェア周辺機器に処理タスクをオフロードする方法およびシステム双方を想定している。図 1 および以下の論述は、本発明を実現可能な適当な計算機環境の端的な概説を与えることを意図したものである。必須ではないが、本発明の実施形態の説明は、総じて、パーソナル・コンピュータが実行するプログラム・モジュールのような、コンピュータ実行可能命令に関して行なう。一般に、プログラム・モジュールは、ルーティン、プログラム、オブジェクト、コンポーネント、データ構造等、特定のタスクを実行するか、あるいは特定の抽象データ・タイプを実施するものを含む。更に、本発明は、ハンドヘルド・デバイス、マルチプロセッサ・システム、マイクロプロセッサを用いたまたはプログラム可能な消費者用電子機器、ネットワーク P C、ミニコンピュータ、メインフレーム・コンピュータ等を含む、その他のコンピュータ・システム・コンフィギュレーションとも実施可能である。また、本発明の実施形態は、通信ネットワークを通じてリンクしたりリモート処理デバイスによってタスクを実行する分散計算機環境においても実施可能である。分散計算機環境では、プログラム・モジュールは、ローカルおよびリモートのメモリ記憶装置双方に位置することができる。

【 0 0 1 8 】

図 1 を参照すると、本発明を実施するシステム例は、従来のパーソナル・コンピュータ 20 の形態の汎用計算機を含み、演算装置 21 ( C P U と呼ぶ )、システム・メモリ 22、およびシステム・メモリを含む種々のシステム・コンポーネントを演算装置 21 に結合するシステム・バス 23 を含む。システム・バス 23 は、メモリ・バスまたはメモリ・コントローラ、周辺バス、および種々のバス・アーキテクチャのいずれかを用いたローカル・バスのいずれでもよい。システム・メモリは、リード・オンリ・メモリ ( R O M ) 24 およびランダム・アクセス・メモリ ( R A M ) 25 を含む。基本入出力システム 26 ( B I O S ) は、起動中のように、パーソナル・コンピュータ 20 内のエレメント間におけるデータ転送を補助する基本的なルーティンを含み、R A M 24 内に格納されている。更に、パーソナル・コンピュータ 20 は、図示しないハード・ディスクの読み書きを行なうハード・ディスク・ドライブ 27、リムーバブル磁気ディスク 29 の読み書きを行なう磁気ディスク・ドライブ 28、C D R O M またはその他の光媒体のようなリムーバブル光ディスク 29 の読み書きを行なう光ディスク・ドライブ 30 のような種々の周辺ハードウエ

10

20

30

40

50

ア・デバイスも含む。ハード・ディスク・ドライブ 27、磁気ディスク・ドライブ 28、および光ディスク・ドライブ 30 は、それぞれ、ハード・ディスク・ドライブ・インターフェース 32、磁気ディスク・ドライブ・インターフェース 33、および光ドライブ・インターフェース 34 を介して、システム・バス 23 に接続されている。ドライブおよびそれに関連するコンピュータ読取可能媒体は、コンピュータ読取可能命令、データ構造、プログラム・モジュールおよびパーソナル・コンピュータ 20 のその他のデータの揮発性格納を行なう。ここに記載する環境の一例は、ハード・ディスク、リムーバブル磁気ディスク 29 およびリムーバブル光ディスク 31 を採用するが、磁気カセット、フラッシュ・メモリ・カード、デジタル・ビデオ・ディスク、ベルヌーイ・カートリッジ、ランダム・アクセス・メモリ (RAM)、リード・オンリ・メモリ (ROM) 等のように、コンピュータによるアクセスが可能なデータを格納することができる、別の形式のコンピュータ読取可能媒体も、動作環境例では使用可能であることは、当業者には認められよう。

#### 【0019】

ハード・ディスク、磁気ディスク 29、光ディスク 31、ROM 24 または RAM 25 上には、多数のプログラム・モジュールを格納可能であり、オペレーティング・システム 35、1 つ以上のアプリケーション・プログラム 36、その他のプログラム・モジュール 37、およびプログラム・データ 38 を含む。ユーザは、キーボード 40 およびポインティング・デバイス 42 のような入力デバイスによって、コマンドおよび情報をパーソナル・コンピュータ 20 に入力することができる。他の入力デバイス (図示せず) は、マイクロフォン、ジョイスティック、ゲーム・パッド、衛星パラボラアンテナ、スキャナ等を含むことができる。これらおよびその他の入力デバイスは、多くの場合、システム・バスに結合するシリアル・ポート・インターフェース 46 のような周辺ハードウェア・デバイスを介して、演算装置 21 に接続されるが、パラレル・ポート、ゲーム・ポートまたはユニバーサル・シリアル・バス (USB) のようなその他のインターフェースによって接続することも可能である。また、ビデオ・アダプタ 48 のような周辺ハードウェア・インターフェース・デバイスを介して、モニタ 47 またはその他の種類のディスプレイ装置もシステム・バス 23 に接続してある。モニタに加えて、パーソナル・コンピュータは、典型的に、スピーカおよびプリンタのような、その他の周辺出力デバイス (図示せず) を含む。

#### 【0020】

パーソナル・コンピュータ 20 は、リモート・コンピュータ 49 のような 1 つ以上のリモート・コンピュータへの論理接続を用いれば、ネットワーク環境においても動作可能である。リモート・コンピュータ 49 は、別のパーソナル・コンピュータ、サーバ、ルータ、ネットワーク PC、ピア・デバイス、またはその他の共通ネットワーク・ノードとすることができ、典型的に、パーソナル・コンピュータ 20 に関して先に述べた要素の多くまたは全てを含むが、図 1 にはメモリ記憶装置 50 のみを図示している。図 1 に示す論理接続は、ローカル・エリア・ネットワーク (LAN) 51 およびワイド・エリア・ネットワーク (WAN) 52 を含む。このようなネットワーク環境は、会社全域に及ぶコンピュータ・ネットワーク、イントラネットおよびインターネットでは一般的である。

#### 【0021】

LAN ネットワーク環境で用いる場合、パーソナル・コンピュータ 20 は、ネットワーク・インターフェース・カード (NIC) またはアダプタ 53 と多くの場合呼ばれている、周辺ハードウェア・デバイスを介してローカル・ネットワーク 51 に接続する。WAN ネットワーク環境で用いる場合、パーソナル・コンピュータ 20 は、典型的に、モデム 54、またはインターネットのようなワイド・エリア・ネットワーク 52 を通じて通信を確立するその他の手段を含む。モデムは、内蔵でも外付けでもよく、典型的に、シリアル・ポート・インターフェース 46 を介してシステム・バス 23 に接続する。ネットワーク環境では、パーソナル・コンピュータ 20 に関して図示したプログラム・モジュールまたはその一部は、リモート・メモリ記憶装置 50 に格納することも可能である。尚、図示のネットワーク接続は一例であり、コンピュータ間に通信リンクを確立する別の手段も使用可能であることは認められよう。

## 【 0 0 2 2 】

また、本発明の範囲内の実施形態は、実行可能な命令を有するコンピュータ読取可能媒体も含む。このようなコンピュータ読取可能媒体は、汎用コンピュータまたは特殊目的コンピュータがアクセスすることができる、入手可能なあらゆる媒体とすることができる。限定ではなく、一例として、このようなコンピュータ読取可能媒体は、RAM、ROM、EEPROM、CD-ROMまたはその他の光ディスク・ストレージ、磁気ディスク・ストレージまたはその他の磁気記憶ストレージ、あるいは所望の実行可能命令を格納するために使用でき、汎用コンピュータまたは特殊目的コンピュータがアクセスすることができるその他のあらゆる媒体とすることができる。前述の組み合わせも、コンピュータ読取可能媒体の範囲に含まれて当然である。実行可能命令は、例えば、汎用コンピュータ、特殊目的コンピュータ、または特殊目的処理デバイスに、ある機能または機能群を実行させる命令およびデータから成る。最後に、本発明の範囲内の実施形態は、データ構造を表す複数のデータ・フィールドが格納されているコンピュータ読取可能媒体を備えている。

10

## 【 0 0 2 3 】

本発明の実施形態は、演算装置21の処理オーバーヘッドおよびメモリ使用量を削減する機能を備えることを目的とする。これを達成するには、例えば、オペレーティング・システム、アプリケーション・プログラムおよび/または演算装置/CPU21上で実行するその他のプログラム・モジュールによって実行する特定の計算タスクを、コンピュータ・システム20に接続してある適切な周辺ハードウェア・デバイスにオフロードする。このような周辺デバイスの多くは、増々専用のプロセッサやメモリが搭載されるようになっており、典型的にCPU21のみが実行する同じタスクの多くを完全に実行することができる。このようなデバイスの例は、例えば、ネットワーク・インターフェース・カード(図1における53)、ディスク・ドライブ・インターフェース・カード(例えば、図1における32、33、34)、小型コンピュータ・システム・インターフェース(SCSI)デバイス、インテリジェント・シリアル・インターフェース・カード、またはデータの暗号化/解読のためのデバイスというような特定用途周辺機器を含む。

20

## 【 0 0 2 4 】

ここで説明する概略的な発明の概念は、前述の周辺ハードウェア・デバイスのいずれとでも、計算タスクをオフロードするために使用可能であるが、本発明の説明は、現時点において好適な一実施形態の例に関して行なう。ここでは、計算タスクを、図1に示すNIC53のようなネットワーク通信デバイスにオフロードする。即ち、例示の実施形態は、Microsoft Corporation(マイクロソフト社)から入手可能なWindows NTオペレーティング・システムのネットワーキング環境およびアーキテクチャにおいて実施したものとして説明する。とは言え、Windows NTの概念および用語を具体的に参照するものの、殆どではなくとも多くのオペレーティング・システムおよびネットワーキング・アーキテクチャが、本発明の環境に関連する類似性を共有することを当業者は認めよう。

30

## 【 0 0 2 5 】

本発明の内容をより良く理解するために、次に図2を参照する。図2は、Windows NTネットワーキング・モデルを構成するコンポーネントの一部の簡略図を示す。例示の目的のために、種々のWindows NTコンポーネントに対応するOSレイヤも示す。OSモデルにおける物理レイヤに対応するボトム・レイヤに、実際のNIC(ネットワーク・カード、またはネットワーク・アダプタと呼ぶこともある)100~104が常駐している。NICは、物理媒体(ネットワーク・ケーブル)との物理的相互接続、および特定のネットワーク・トポロジにしたがって、上位レイヤ全てが発生するデータを搬送する信号の伝送を行なうハードウェア・デバイスである。先に注記したように、多くのNICは専用プロセッサおよびメモリを搭載しており、精巧な計算タスクも実行可能である。そのタスクには、ホスト・プロセッサCPUが通常であれば処理するようなタスクが含まれる。NICは、物理的に、コンピュータ内のスロットに配置するプリント回路ボード・カードとして、マザー・ボード上のコンピュータ・シャーシ内に配置した専用チップ

40

50

プとして、またはその他のいずれかの適当な方法で実施することができる。

【0026】

各NICは、双方向ライン108～112によって概略的に表すように、対応するネットワーク・ドライバ116～120を介して論理的にWindows NT ネットワーキング・モデルと相互接続している。ネットワーク・ドライバは、ネットワーク・モデルのMACサブレイヤ内に常駐し、対応するNICを介してWindows NTを物理ネットワーク・チャンネルにリンクする。各ドライバは、典型的に対応するNICのベンダが供給するソフトウェア・コンポーネントとして実施され、その対応するネットワーク接続を通じてパケットを送受信する役割、およびオペレーティング・システムの代わりにNICを管理する役割を果たす。また、各ドライバは、対応するNIC上でI/Oを開始し、これらから割込を受け取り、上流側のプロトコル・ドライバにコールし、アウトバウンド・データ転送の完了をこれらに通知する。また、デバイス・ドライバは、対応するNICの追加処理能力全てに対する呼び出し、制御および/または監視を行なう役割も担う。

10

【0027】

環境によっては、TCP/IPまたはXNSのような単一のネットワーク特定プロトコルを実施するように、ドライバ・コンポーネントに書き込む。ここに記述し特許を請求する本発明の基本的発明は、このような環境に適用可能である。しかしながら、例示の目的のために、本発明は、Windows NTネットワーク・アーキテクチャに関連付けて説明することにする。この場合、ネットワーク・ドライバ・インターフェース仕様(NDIS)と呼ばれるインターフェースおよび環境が与えられる。NDISインターフェースは、機能的には、図2の126に示す通りである。NDISは、ネットワーク・ドライバ116～120の各々を、種々のトランスポート・プロトコル(その例を128～134で示す)の詳細から、およびその逆に、遮蔽する。即ち、NDISは、1つ以上のNICドライバ(116～120)が1つまたは多数の下位NIC(100～104)、1つまたは多数の上位トランスポート・プロトコル・ドライバ、あるいはトランスポート(図2の128～134で表す)、およびオペレーティング・システムと通信するためのインターフェースを記述する。

20

【0028】

本質的に、NDISはNICドライバの開発のために、完全に抽象化した環境を定義する。したがって、NICドライバが実行すべき各外部機能毎に、NICハードウェア割込の登録および代行受信(intercept)からトランスポート・プロトコル・ドライバとの通信、登録操作およびポートI/Oを介した下位のNICとの通信まで、機能を実行するためにNDIS APIに頼ることができる。このレベルの抽象化およびその結果としての移植性を得るために、NDISは、NDISインターフェース・ライブラリ・ラップ(NDIS interface Library Wrapper)(図示せず)と呼ぶ、エクスポート・ライブラリを用いる。NICドライバおよびプロトコル・ドライバ間、NICドライバおよびオペレーティング・システム間、ならびにNICドライバおよびNIC間の相互作用は全て、ラップ関数へのコールによって実行する。したがって、Windows NTのトランスポート特定ドライバに書き込む代わりに、ネットワーク・ベンダは、NDISインターフェースを、単一ネットワーク・ドライバの最上位レイヤとして与える。こうすることにより、いずれのプロトコル・ドライバも、このインターフェースをコールすることにより、そのネットワーク要求をネットワーク・カードに宛てて送ることが可能となる。したがって、ユーザは、TCP/IPネットワークおよびDL(または、NWLINK、またはDECnet、VINES、NetBEUI等)ネットワークを通じて、1つのネットワーク・カードおよび単一のネットワーク・ドライバを用いて通信することができる。

30

40

【0029】

ネットワークおよびデータ・リンク・レイヤは、一例として図2の128～134で示す、トランスポート、プロトコルおよび関連するドライバである。Windows NTでは、トランスポート・プロトコル・ドライバは、トランスポート・ドライバ・インターフェース(TDI)、または恐らくその上縁における別の特定用途インターフェースを実施

50

し、ネットワークのユーザにサービスを提供するソフトウェア・コンポーネントである。Windows NTでは、TDIは、機能ブロック138に示すリディレクトおよびサーバ機能のような、セッション・レイヤにおいて通信するネットワーク・コンポーネントのために共通インターフェースを備える。公知のように、トランスポート・プロトコルは、ネットワークのためのデータ・オーガナイザとして作用し、本質的にどのようにデータを次の受信レイヤに提示し、それに応じてデータをパッケージ化するかを定義する。これらは、パケット（Windows NT関連ではNDISパケットとも呼ぶ）を割り当て、データを送出側アプリケーションからパケットにコピーし、NDISをコールすることによってパケットを下位デバイス・ドライバに送り、対応するNICを介してデータをネットワーク上に送出できるようにする。

10

#### 【0030】

尚、データ・パケットが種々のネットワーク・レイヤ、典型的に、ネットワーク・モデルのレイヤ3および4を通過する際に、追加の機能またはタスクもデータ・パケットに対して実行可能であることは認められよう。例えば、トランスポート・プロトコル・ドライバは、チェックサム値を計算し、これをパケットに添付する。一般に、この動作には、ネットワーク・パケットの送付側に対応するトランスポート・プロトコルが、パケットを構成するデータ・エレメントを全て加算することによって計算した数値を、それに添付する必要がある。次いで、パケットの受信側が、添付されているチェックサム数値をデータと比較することによって、データが移動中に変化しなかったことを確認する。

20

#### 【0031】

NICハードウェアにおいて最適に実行可能な別の関連タスクに、データ・パケットに対するメッセージ・ダイジェストの計算がある。チェックサムと同様、メッセージ・ダイジェストは、パケット内のデータの完全性を保証するために用いられる。加えて、メッセージ・ダイジェストは、メッセージを送った側が意図した者であることを確かめることによって、データの認証を保証するために用いることもできる。メッセージ・ダイジェストの計算は非常にCPU集中的であり、ソフトウェアで実施するには不経済な機能である。

#### 【0032】

別の望ましい機能に、パケット内のデータの暗号化がある。暗号化とは、パッケージ内のメッセージを変換することにより、パケットの無許可の読み手が、暗号化キーを予め知らなければ、メッセージの内容を実際に見ることができないようにする暗号化プロセスのことを言う。勿論、暗号アルゴリズムも非常にCPUおよびメモリ集中的となりがちであり、ソフトウェアで実行すると、法外な程不経済となる可能性がある。

30

#### 【0033】

データ・パケット上で実行可能な別のタスクに、TCPセグメント化がある。公知のように、TCPは大きなデータ・パケットを、下位ネットワークが許可する最大データ・サイズに合ったセグメントに区分する。例えば、イーサネットはネットワーク上で最大1514バイトのパケットを許可する。したがって、TCPが例えば64Kバイトを送らなければならない場合、データを1514バイトのセグメントに分解しなければならない。

#### 【0034】

これらおよびその他の機能は、典型的に、種々のネットワーク・レイヤに常駐するソフトウェア・コンポーネントにおいて、コンピュータのCPU20が実行するので、かなりのコンピュータ資源を利用し、コンピュータ・システムの性能の全体的な低下を招く可能性がある。したがって、これらまたはその他の同様のタスクをオフロードし、対応するNICにおいてこれらを代わりに実行するようにすれば、コンピュータ・システムの全体的な速度および効率を大幅に向上させることが可能となる。

40

#### 【0035】

先に注記したように、Windows NTまたは同様のレイヤ状ネットワーク・モデルにおけるデータ伝送の基本単位はデータ・パケットである。Windows NT環境では、データ・パケットのことをNDISパケットと呼ぶ。各パケットは、スタックの最上位（即ち、ISOスタックにおけるレイヤ5）から最下位のソフトウェア・レイヤ（

50

即ち、ISOスタックにおけるレイヤ2)まで移動する。したがって、パケットは、データの送信および受信の間にレイヤを通過する際、各レベルを通じて共通のデータ構造を定義する。一例として、図3は、パケットが各レイヤを通過して、イーサネットNICとして100で示すNICまで達する際に、パケットが辿る経路を示す。先に注記したように、トランスポート・ドライバ128は、送出側アプリケーションからデータを受け取り、下位プロトコルと一貫性のあるパケットにこれをパッケージ化し、次いでNDISインターフェース126を介して下位のデバイス・ドライバ116にパケットを転送する。加えて、トランスポート・プロトコルは、パケットに対して他の機能(例えば、チェックサムの計算等)も行なうことができる。あるいは、他の機能コンポーネントが、図3に示すIPセキュリティ機能144(例えば、暗号化および/またはパッケージ・ダイジェストの計算)のように、パケットに追加の機能を実行するネットワーク・レイヤまたはデータ・リンク・レイヤに常駐することも可能である。

10

#### 【0036】

本発明の好適な実施形態の1つでは、データ・パケット142は、計算タスクを、NICハードウェア100のような周辺デバイスにオフロードするための手段である。例えば、図3において、アプリケーション・データ140は、ネットワーク・モデルの上位レイヤから、TCP/IP128のような適切なトランスポート・プロトコル・ドライバに渡される。ドライバはデータを適切なデータ・パケット142にパッケージ化し直す。次いで、この特定のデータ・パケット142にどの追加機能を実行するかに応じて、パケット拡張部と呼ぶ、既定のデータ構造をデータ・パケットに添付する機能コンポーネントを含ませる。以下で更に詳細に説明するが、パケット拡張部の内容は、データ・パケットがNIC100に到達したときに、どのタスクまたは複数のタスクをこのデータ・パケットに対して実行するかを示す。データ・パケット142がネットワーク・ドライバ116に到達すると、ネットワーク・ドライバ116がこのパケット拡張部の内容を問い合わせ、どのタスク(複数のタスク)をNIC100が実行するのかを確認する。次いで、ドライバ116は、NIC上のハードウェアを制御/操作し、特許拡張部の内容を通じて要求されたいかなる機能タスクをも実行する。

20

#### 【0037】

例えば、図3において、データ・パケット142をソフトウェア・コンポーネント144に渡す。ソフトウェア・コンポーネント144は、別個に実施することも、あるいはトランスポート・プロトコル・ドライバ自体の一部として実施することも可能であり、パケット拡張部をパケット142に添付する。オフロードする特定のタスクに応じて、パケット拡張部にデータを含ませる。例えば、IPセキュリティ機能を実施する場合、NICが指定された暗号化キーにしたがってデータ・パケットを暗号化することを示すデータを含ませる。勿論、ソフトウェア・コンポーネント144は、既定のデータを添付し、先に説明したような多数の機能のいずれか1つを、ネットワーク・レイヤ内に常駐するソフトウェア・コンポーネントによって実行する代わりに、ハードウェア・レベルで実行することができる。デバイス・ドライバ116は、パケット拡張部から情報を抽出し、次いでNIC100において指定されたタスク(複数のタスク)を呼び出す。

30

#### 【0038】

図4は、パケット142の概略構造の現時点における好適な実施形態の1つを示す。パケット142は、Windows NT環境において用いられる正確なネットワーク環境に応じて、いかなるフォーマットでも可能であるが、パケットはNDISにしたがってフォーマット化してあり、パケット記述子、協働するデバイス・ドライバ(ドライバ群)およびプロトコル・ドライバ(ドライバ群)が意味を規定するフラグ、当該パケットに関連する帯域外データ(OOB)の記憶領域、パケット長に関する情報、ならびにパケットのデータ内容に係るメモリ位置に対するポインタというような情報を含む。

40

#### 【0039】

図4は、更に、タスクのオフロードを識別するためにNDISデータ・パケットに添付する、追加のデータ構造フィールド、即ち、パケット拡張部150も示す。先に説明したよ

50

うに、宛先のNICにオフロードしている特定のタスクまたは複数のタスクの識別に必要な情報を収容するデータ構造を定義するのは、このパケット拡張部150である。好適な実施形態では、各タスク・オフロード・タイプ（例えば、チェックサム、暗号化／解読、等）毎に、既定のデータ・フィールドをパケット拡張部150内に含ませる。このデータ・フィールドは、単に、制御フラグまたは複数のフラグの形式とすることができる。制御フラグは、単に、特定の機能を実行する（チェックサム等）ことを示すに過ぎない。または、情報は、タスクをどのように遂行すべきかについて更に定義するデータ構造へのポインタの形式とすることも可能である。例えば、図4に示す例では、パケット拡張部150は、NICは152で示すチェックサム動作を実行することを示すフラグを収容する。このタイプのタスクでは、パケット拡張部は、受信局におけるNICが、送出側のNICが計算したチェックサムの有効性をチェックするようにも設定する。

10

【0040】

限定ではなく一例として、NICがチェックサム動作を実行することを意味するパケット拡張部データ構造の好適な実施形態は、以下の構造を有する。

【0041】

【表1】

```

typedef struct _NDIS_TCP_IP_CHECKSUM_PACKET_INFO
{
    union
    {
        struct
        {
            ULONG NdisPacketChecksumV4:1;
            ULONG NdisPacketChecksumV6:1;
        }
        Transmit;
        struct
        {
            ULONG NdisPacketTcpChecksumFailed:1;
            ULONG NdisPacketUdpChecksumFailed:1;
            ULONG NdisPacketIpChecksumFailed:1;
            ULONG NdisPacketTcpChecksumSucceeded:1;
            ULONG NdisPacketUdpChecksumSucceeded:1;
            ULONG NdisPacketIpChecksumSucceeded:1;
            ULONG NdisPacketLoopback:1;
        }
        Receive;
        ULONG Value;
    };
}

NDIS_TCP_IP_CHECKSUM_PACKET_INFO,
*PNDIS_TCP_IP_CHECKSUM_PACKET_INFO;

```

この特定の packets 拡張部データ構造では、変数NdisPacketChecksumV4およびNdisPacketChecksumV6双方をセットしない場合、デバイス・ドライバは、データ・パケットに対してチェックサムを全く行なうことなく、これを送る。

【 0 0 4 2 】

また、図 4 は、パケット・データの暗号化および / またはメッセージ・ダイジェストの計算に関連して実行する、セキュリティ機能 1 5 4 も、送出側 N I C が実行すべきことを指定する。このタイプのタスクでは、フィールド 1 5 4 は、データ構造を収容するメモリ位置に対するポインタを収容することが好ましい。一方、このデータ構造は、暗号化および / またはメッセージ・ダイジェスト機能の実行に関連する情報を収容する。状況によって



は、関連データを有するメモリ位置に対するポインタを含ませることは、パケット拡張部自体の中に実際のデータを格納するよりも有効な場合がある。

【 0 0 4 3 】

このような利点の1つを図5に示す。これは、好適なプログラム・シーケンスの一例を示し、多数の連続するデータ・パケットに同じタイプの暗号化またはダイジェスト計算動作を実行しなければならない状況に対応する。プログラム・ステップ302において開始した後、プログラム・ステップ304は、現パケットがパケット・シーケンス内の最初のパケットであるか否かについて判定を行なう。そうである場合、この最初のパケットに、続くパケットに対する動作にも同様に用いる情報即ちコンテキストを与える。例えば、最初のパケットは、用いる特定の暗号化キーを明記するパケット拡張部を有する。この値即ちコンテキストは、別個のメモリ位置即ちハンドルに格納する。連続するパケットは、この情報を含む必要がなく、例えば、プログラム・ステップ308に示すように、暗号化キー（またはその他のコンテキスト情報）が格納されているメモリ位置に対するポインタのみを有する。この手法は、パケットのシーケンスにおける後続のデータ・パケット全体のサイズを縮小し、更にタスク・オフロード方法の効率および移植性を高める。

10

【 0 0 4 4 】

好適な実施形態では、パケット拡張部150内に収容する情報は、パケット142を送る特定のデバイス・ドライバによる問い合わせを受ける。例示の実施形態において記載しているWindows NT環境では、このタイプの機能は、適切なNDIS関数をコールすることによって実行することが好ましい。例えば、パケットに対するパケット拡張部150のメモリ位置へのポインタを戻す既定のNDIS関数を実行することができる。このときには、デバイス・ドライバ・ソフトウェアは、どのタスクを実行するのか識別し、オフロードするタスク（複数のタスク）に応じて、適切にドライバの対応するNICハードウェアを動作させ/操作する。

20

【 0 0 4 5 】

実際のデータ・パケットを利用して計算タスクをコンピュータ・プロセッサからハードウェア周辺機器にオフロードすることは、多くの理由から有利である。例えば、ポート・ドライバは、パケット毎に、周辺機器の能力を利用することができる。これによって、タスクを動的にダウンロードすることができ、周辺機器の能力を必要に応じて用いることができる。したがって、特定の時点において、コンピュータ・システムの処理オーバーヘッドが低い場合、あるタスクを従来通りコンピュータ・プロセッサ上で実行することが望ましいこともある。あるいは、CPUが他の計算タスクを多く負担する場合、単にデータ・パケットに必要なパケット拡張部を添付することによって、タスクを周辺デバイスにオフロードすることができる。

30

【 0 0 4 6 】

他の利点は、単一のパケットによって多数のタスクをオフロードできることであり、本質的に、多数の動作を一度に「バッチ処理」する。例えば、コンピュータ・プロセッサがチェックサム動作または暗号化動作を実行する場合、その動作、即ち、チェックサムの計算またはパケット・データの暗号化を完了することができるようにするには、その前にデータ・フィールド全体をメモリ位置にロードしなければならない。更に、レイヤ状ネットワーク・モデルのため、一度に実行できる動作は1つのみであり、データをメモリに多数回コピーする必要がある。しかしながら、パケット単位の手法では、多数のタスクを1つのパケットにオフロードすることができる。したがって、ハードウェア周辺機器は、ハードウェアの能力に応じて、データ上での単一のパスにおいて2つ以上の動作を実行することができ、コンピュータ・システムのスループットおよび効率を大幅に向上させることができる。

40

【 0 0 4 7 】

前述の方式は、特定のNICにオフロードするタスクを指定する機能を用いる場合特に有用であるが、パケット単位の情報転送は他の方法でも同様に使用可能であることは認められよう。例えば、特定のNICが所定の時点にパケットの配信をスケジュールすることが

50

できる場合、パケット拡張データ構造は、N I Cハードウェアがパケットをどのようにおよび/またはいつ送るかを識別する情報を渡すためにも使用することができる。

【 0 0 4 8 】

本発明の好適な実施形態では、トランスポート・プロトコル・ドライバがパケット拡張部をデータ・パケットに添付し、特定のタスクをN I Cにオフロードする前に、2つの追加の機能を最初に実行する。異なるタイプのハードウェア周辺機器が数多くあり、各々が異なる処理能力を有する限りにおいて、本発明の実施形態は、トランスポート・ドライバが最初にコンピュータ・システムに接続してある周辺機器のタスク・オフロード能力を問い合わせることができるようにする手段を備えることが好ましい。一旦これらの機能を確認し終われば、トランスポート・プロトコル・ドライバは、対象のタスクを設定またはイネーブルすることができる。一旦イネーブルすれば、前述のようにパケット毎に指定のタスクを続いて利用することができる。

10

【 0 0 4 9 】

図6は、N I Cのような周辺機器のタスク・オフロード機能を問い合わせるこの機能を実施し、次いで必要と思われるタスクを設定/イネーブルする現時点において好適な1つのプログラム・ステップ・セットを示す。好ましくは、図示のプログラム・ステップは、トランスポート・プロトコル・ドライバ内に統合化したソフトウェア・コンポーネントまたはモジュールとして実施する。更に、Windows NT環境において実施する場合、機能的動作およびドライブ間通信の多くは、N D I Sインターフェースを介して実行することが好ましい。勿論、異なるオペレーティング・システム環境において実施する場合、または非ネットワーク・タイプの周辺機器と共に実施する場合、図示のプログラム・ステップはそれに応じて変更しなければならない。

20

【 0 0 5 0 】

実行可能なプログラム命令は、コンピュータ・システムのC P U (例えば、図1の演算装置21)が実行する。プログラム・ステップ202において開始し、当該コンピュータ・システムに接続してある周辺機器(複数の周辺機器)のタスク・オフロード能力について問い合わせを行なう。例示の環境では、周辺機器(複数の周辺機器)は、システムに付属するN I C (複数のN I C)である。好ましくは、各N I Cデバイス・ドライバ(図2の116~120)には、既定のタスク・オフロード・バッファ位置(複数の位置)が関連付けられており、その各々が当該デバイス・ドライバおよびその対応するN I Cのタスク・オフロード能力を収容する。少なくとも、タスク・オフロード・バッファは、N I Cおよびそのデバイス・ドライバが対応する個々のタスク(複数のタスク)を識別し、更に、対応する個々のタスク毎に特定のあらゆる情報を含む。好適な実施形態では、デバイス・ドライバのタスク・オフロード・バッファの内容は、N D I S関数コールによって検索する。

30

【 0 0 5 1 】

限定ではなく一例として、ドライバのタスク・オフロード・バッファの好適な実施形態は以下の構造を有する。

【 0 0 5 2 】

【表2】

40

```

typedef
_NDIS_TASK_OFFLOAD
{
    ULONG                Size;
    NDIS_TASK            Task;
    ULONG                OffsetNextTask;
    ULONG                TaskBufferLength;
    UCHAR                TaskBuffer[1];
}

NDIS_TASK_OFFLOAD,
*PNDIS_TASK_OFFLOAD;

```

10

ここで、  
【 0 0 5 3 】  
【 表 3 】

20

Size	この変数は、この構造のサイズにセットし、使用する構造のバージョンを判定するために用いる。
Task	これは、構造において定義したタスク・オフロードのタイプを示す。
OffsetNextTask	次のタスク・オフロード・バッファに対するオフセット。これが、このドライバ／NICが対応する最後のタスクである場合、値は0である。
TaskBufferLength	この構造に従うタスク・オフロード・バッファの長さ。
TaskBuffer	これは、この構造で定義した特定のタスク・オフロードを記述するために必要な情報／データを収容するタスク特定バッファである。

30

40

一旦各周辺機器の個々のタスク・オフロード能力を確認したなら、プロセッサ 2 1 は、図 6 のプログラム・ステップ 2 0 4 に示す機能に対応するコンピュータ実行可能命令に進む。ここで、各周辺NICが対応する個々のタスク・オフロード能力のいずれかが、トランスポートに該当するものであるか否かについて判定を行なう。当該周辺機器に対応するタスクがない場合、または対応するタスクが当該トランスポートには有用ではない場合、プロセッサは直接プログラム・ステップ 2 0 8 に進み、この処理の特定のラインを停止し、トランスポートはその通常動作に進む。あるいは、識別したタスクの 1 つ以上が該当する場合、プログラム・ステップ 2 0 6 を実行する。このステップにおいて、実行可能命令を実行し、プロトコル・ドライバが望む特定のタスク・オフロード能力をイネーブルにセッ

50

トする。これは、デバイス・ドライバおよびその対応するNICに、どのタイプのタスク・オフロードをパケット毎に予期する可能性があるのかについて知らせる。好適な実施形態では、タスク・オフロード能力をセットするには、データ構造内に適切なデータ値をセットし、次いで、NDISインターフェース関数コールによって、対応するデバイス・ドライバに受け渡す。例えば、個々のタスクについて、当該タスクのタスク・オフロード・バッファ構造内の適切なビットをトランスポートによってセットすることにより、当該ドライバ/NICのタスクをイネーブルすることができる。こうして、トランスポートは、後に使用したくなる可能性がある各NIC/NICドライバのタスク・オフロード能力をいくつでもイネーブルすることができる。

#### 【0054】

一旦所望のオフロード能力を各デバイス・ドライバおよびその対応するNIC毎にイネーブルしたなら、コンピュータ・システム・プロセッサ21はプログラム・ステップ208に進み、この特定の機能に対する処理が終了する。この時点において、トランスポート・ドライバは、前述のように各パケット毎に、イネーブルしたタスク・オフロード能力の各々を利用することができる。

#### 【0055】

要約すれば、本発明の実施形態は、現在従来技術において利用可能なものと比較して、格別な利点を与える。コンピュータ・システムのプロセッサおよびメモリにおいて本来実行すべき具体的な処理タスクを、代わりに、当該コンピュータに接続してある特定の周辺デバイスまたは複数の周辺デバイスにダウンロードする。次いで、周辺機器によって計算タスクを実行することにより、コンピュータ・システムの資源を他の計算タスクのために保存する。特に、オフロードした処理タスクがCPUおよび/またはメモリ集中的な場合、このタスク・オフロード方式は、コンピュータ・システム全体の計算効率を劇的に向上させる。開示したコンピュータ・タスク・オフロード方法は、動的に要求に応じてタスクをオフロードすることを可能にする手段を備えるという利点がある。したがって、プロセッサは、他の計算タスクの処理に忙しく、プロセッサのオーバーヘッドが高いときに、タスクをオフロードすることができる。逆に、コンピュータ・システムの処理資源に対する要求が低くなった場合には、代わりにプロセッサはタスクをそれ自体で実行すればよい。加えて、特定の周辺機器に多数のタスクをバッチでオフロードすることができる。多くの場合、周辺機器は、このような多数のタスクを、コンピュータ・プロセッサよりはるかに効率的に実行するように最適化されているので、システム全体の効率が更に改善する結果となる。最後に、本発明の方法は、周辺機器の特定のタスク・オフロード能力について問い合わせ、その後を選択的にイネーブルすることを可能にする処理方式を提供する。このように、コンピュータ・システムは、その種々の周辺デバイスの能力を容易に識別し、後に必要となる処理能力またはその可能性がある処理能力のみを利用することができる。

#### 【0056】

本発明は、その精神または本質的な特徴から逸脱することなく、他の特定形態でも具体化が可能である。前述の実施形態は、あらゆる観点においても、例示であり限定として解釈すべきではない。したがって、本発明の範囲は、前述の説明ではなく、添付した特許請求の範囲によって示すこととする。特許請求の範囲の均等の意味および範囲に該当するあらゆる変更は、その範囲に含まれるものとする。

#### 【図面の簡単な説明】

【図1】 本発明と共に使用可能な典型的なコンピュータ・システムおよび付属周辺デバイスの一例を示す図である。

【図2】 レイヤ状ネットワーク・アーキテクチャ内にある機能コンポーネントの一部を示す図である。

【図3】 本発明の現時点における好適な実施形態の1つに応じたプログラム・コンポーネントによるデータ・パケットのフローを示す機能ブロック図である。

【図4】 データ・パケットおよびパケット拡張の現時点における好適な実施形態の1つを示す図である。

10

20

30

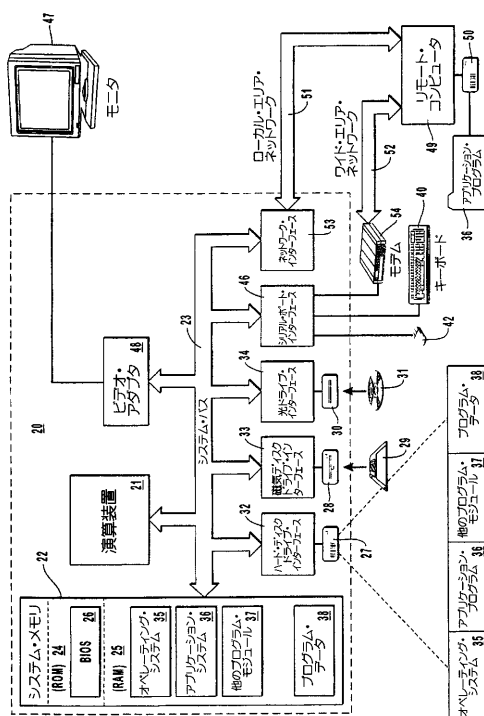
40

50

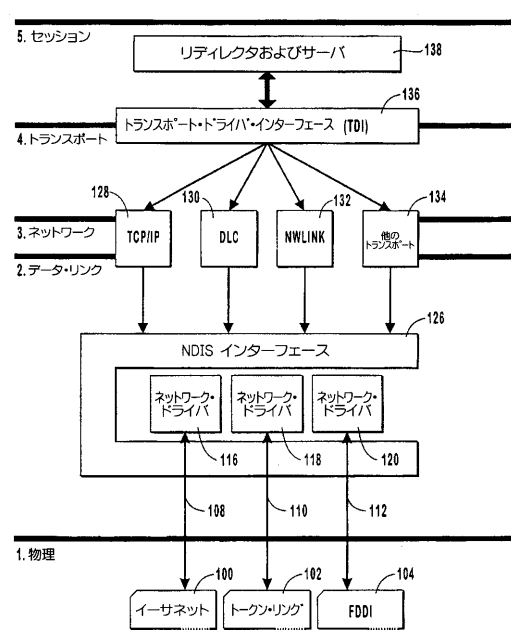
【図 5】 パケット毎にタスクをオフロードする際に用いるプログラム・ステップの現時点における好適な実施形態の 1 つを示すフロー・チャートである。

【図 6】 周辺デバイスのタスク・オフロード能力を問い合わせそして設定するために用いるプログラム・ステップの現時点における好適な実施形態の 1 つを示すフロー・チャートである。

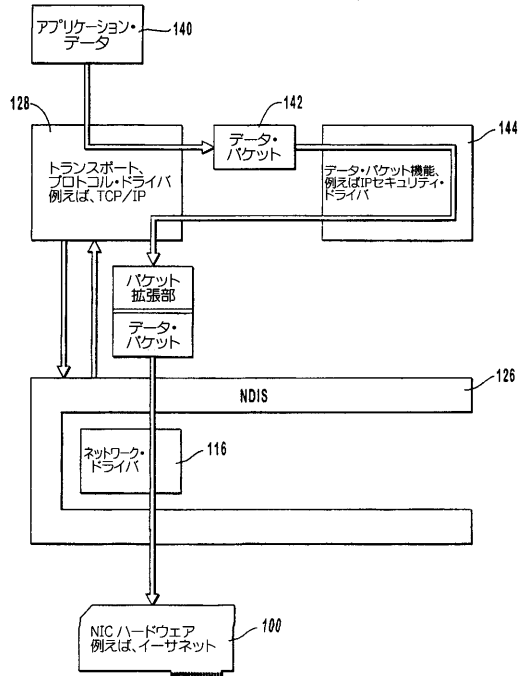
【 図 1 】



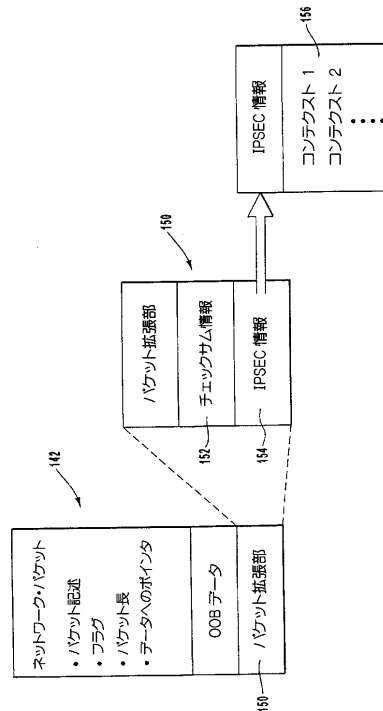
【圖 2】



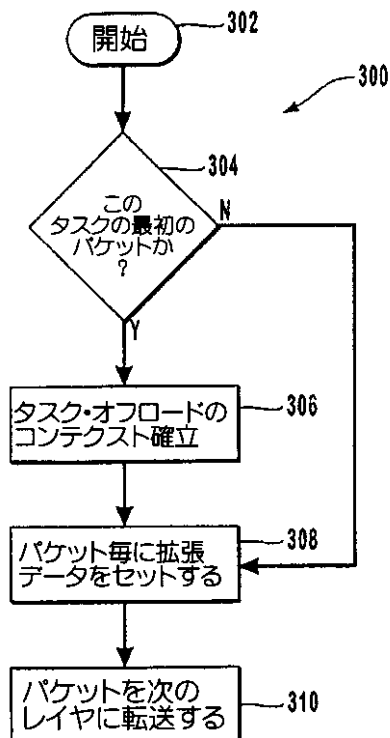
【図 3】



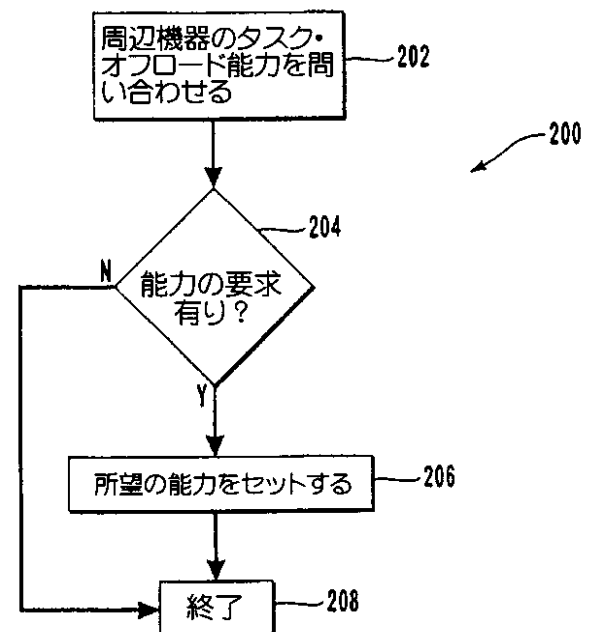
【図 4】



【図 5】



【図 6】



---

フロントページの続き

- (72)発明者 アナンド, サンジェイ  
アメリカ合衆国ワシントン州 9 8 0 5 2, レッドモンド, オールド・レッドモンド・ロード 7 0  
0 1, アpartment エヌ 1 5 5
- (72)発明者 ブランドン, カイル  
アメリカ合衆国ウィスコンシン州 5 3 7 0 3, マディソン, ウィリアムソン・ストリート 7 5 4
- (72)発明者 スリニヴァス, シク  
アメリカ合衆国ワシントン州 9 8 0 2 9, イサクア, サウス・イースト・フォーティセカンド・ス  
トリート 2 5 0 4 1
- (72)発明者 ハイダー, ジャミール  
アメリカ合衆国ワシントン州 9 8 0 5 3, レッドモンド, ノース・イースト・シックスティーンス  
・プレイス 2 3 2 9 2

審査官 鈴木 修治

- (56)参考文献 米国特許第 0 5 6 3 4 0 7 0 ( U S , A )  
特開平 0 9 - 2 3 4 8 5 3 ( J P , A )  
特許第 2 9 5 0 3 6 6 ( J P , B 2 )  
特開平 1 0 - 1 2 4 4 6 8 ( J P , A )  
特開平 0 6 - 0 1 9 8 5 5 ( J P , A )

(58)調査した分野(Int.Cl. , D B 名)

G06F 9/46-9/54  
G06F 15/16-15/177  
G06F 13/10-13/14