

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5689526号
(P5689526)

(45) 発行日 平成27年3月25日 (2015. 3. 25)

(24) 登録日 平成27年2月6日 (2015. 2. 6)

(51) Int. Cl.	F I
HO 4 L 29/10 (2006. 01)	HO 4 L 13/00 3 O 9 Z
GO 6 F 13/10 (2006. 01)	GO 6 F 13/10 3 3 O C
GO 6 F 9/50 (2006. 01)	GO 6 F 9/46 4 6 2 Z
HO 4 L 13/08 (2006. 01)	HO 4 L 13/08

請求項の数 15 (全 18 頁)

(21) 出願番号	特願2013-505372 (P2013-505372)	(73) 特許権者	390009531
(86) (22) 出願日	平成23年3月1日 (2011. 3. 1)		インターナショナル・ビジネス・マシーンズ・コーポレーション
(65) 公表番号	特表2013-530573 (P2013-530573A)		INTERNATIONAL BUSIN
(43) 公表日	平成25年7月25日 (2013. 7. 25)		ESS MACHINES CORPOR
(86) 国際出願番号	PCT/EP2011/052992		ATION
(87) 国際公開番号	W02011/131400		アメリカ合衆国10504 ニューヨーク
(87) 国際公開日	平成23年10月27日 (2011. 10. 27)		州 アーモンク ニュー オーチャード
審査請求日	平成25年11月11日 (2013. 11. 11)		ロード
(31) 優先権主張番号	12/766, 282	(74) 代理人	100108501
(32) 優先日	平成22年4月23日 (2010. 4. 23)		弁理士 上野 剛史
(33) 優先権主張国	米国 (US)	(74) 代理人	100112690
			弁理士 太佐 種一
		(74) 代理人	100091568
			弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 マルチキュー・ネットワーク・アダプタの動的再構成によるリソース・アフィニティ

(57) 【特許請求の範囲】

【請求項 1】

データ処理システムにおいて、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供する方法であって、

前記データ処理システム内のデバイス・ドライバによって、メモリ内に初期キュー・ペアを割り当てるステップと、

前記デバイス・ドライバによって、前記データ処理システムのワークロードが所定の高閾値より高くなったかどうかを判断するステップと、

前記ワークロードが前記所定の高閾値より高くなることに応じて、前記デバイス・ドライバによって、前記メモリ内に追加のキュー・ペアを割り当て、前記追加のキュー・ペアを初期化するステップと、

前記追加のキュー・ペアに関連する追加の処理エンジンの動的挿入を可能にするために、前記デバイス・ドライバによって、ネットワーク・アダプタ内の受信側スケーリング (RSS) メカニズムをプログラムするステップと、

前記デバイス・ドライバによって、前記追加のキュー・ペアへの送信タプル・ハッシングを有効化するステップと、

を含む方法。

【請求項 2】

前記追加のキュー・ペアを割り当て、前記追加のキュー・ペアを初期化するステップ、前記プログラムするステップ、および前記有効化するステップが、前記デバイス・ドライ

10

20

バによって、前記データ処理システムの前記ワークロードが前記所定の高閾値より高くなるたびに実行される、請求項 1 に記載の方法。

【請求項 3】

前記追加のキュー・ペアに関連する前記追加の処理エンジンの動的挿入を可能にするために、前記ネットワーク・アダプタ内の前記受信側スケーリング (RSS) メカニズムをプログラムするステップは、

前記デバイス・ドライバによって、前記初期キュー・ペアへの送信タブル・ハッシングを有効化するステップ

を含む、請求項 1 に記載の方法。

【請求項 4】

前記データ処理システムの前記ワークロードが前記所定の高閾値より高くなったかどうかを判断するステップは、前記デバイス・ドライバによって、データ・フローおよびリソース利用可能性を介して前記データ処理システムの前記ワークロードを監視することによって実行される、請求項 1 に記載の方法。

【請求項 5】

前記デバイス・ドライバは、前記ワークロードに関連する少なくとも 1 つのパラメータを監視し、前記少なくとも 1 つのパラメータは、送信 / 受信バイト・パー・セカンド、前記ネットワーク・アダプタによって送信および受信されているフロー制御フレームの数、前記ネットワーク・アダプタによって検出される DMA オーバーランの数、前記デバイス・ドライバによって検出される送信タイムアウト・イベントの数、割り込み毎に前記デバイス・ドライバによって処理される受信パケットの数、またはソフトウェア・キュー上の送信パケットの数のうちの少なくとも 1 つである、請求項 4 に記載の方法。

【請求項 6】

前記デバイス・ドライバによって、前記ワークロードが所定の低閾値を下回ったかどうかを判断するステップと、

前記ワークロードが前記所定の低閾値を下回るのに応じて、前記メモリ内に割り当てられているキュー・ペアが残り 1 つのみあるかどうかを、前記デバイス・ドライバによって判断するステップと、

前記メモリ内に割り当てられているキュー・ペアが 2 つ以上残っていることに応じて、割り当てられているキュー・ペアの削除を可能にするために、前記デバイス・ドライバによって前記ネットワーク・アダプタ内の前記 RSS メカニズムを再プログラムするステップと、

特定されたキュー・ペアへの送信タブル・ハッシングを、前記デバイス・ドライバによって無効化するステップと、

前記特定されたキュー・ペアへの前記ワークロードが休止したかどうかを、前記デバイス・ドライバによって判断するステップと、

前記特定されたキュー・ペアへの前記ワークロードが休止することに応じて、前記特定されたキュー・ペアを前記デバイス・ドライバによってメモリから除去し、その結果、前記特定されたキュー・ペアによって使用されたメモリを解放するステップと、

をさらに含む、請求項 1 に記載の方法。

【請求項 7】

前記特定されたキュー・ペアをメモリから除去し、その結果、前記特定されたキュー・ペアによって使用されたメモリを解放するステップより前に、前記特定されたキュー・ペアへの前記ワークロードが休止するのを、前記特定されたキュー・ペアへの前記ワークロードが休止しないことに応じて、前記デバイス・ドライバによって待つステップ

をさらに含む、請求項 6 に記載の方法。

【請求項 8】

請求項 1 ~ 7 の何れか 1 項に記載の方法の各ステップをコンピュータに実行させる、コンピュータ・プログラム。

【請求項 9】

プロセッサと、
前記プロセッサに結合されたメモリと、
を含む装置であって、前記メモリは、命令を含み、前記命令は、前記プロセッサによって実行されると、前記プロセッサに、
メモリ内に初期キュー・ペアを割り当てること、
前記データ処理システムのワークロードが所定の高閾値より高くなったかどうかを判断すること、
前記ワークロードが前記所定の高閾値より高くなることに応じて、前記メモリ内に追加のキュー・ペアを割り当て、前記追加のキュー・ペアを初期化すること、
前記追加のキュー・ペアに関連する追加の処理エンジンの動的挿入を可能にするために、ネットワーク・アダプタ内の受信側スケーリング（RSS）メカニズムをプログラムすること、および
前記追加のキュー・ペアへの送信タプル・ハッシングを有効化すること
を行わせる、装置。

【請求項 10】

前記命令は、前記プロセッサに、
前記追加のキュー・ペアを割り当て、前記追加のキュー・ペアを初期化すること、前記プログラムすること、および前記有効化することを行わせる命令を、前記データ処理システムの前記ワークロードが前記所定の高閾値より高くなるたびに実行させる、請求項 9 に記載の装置。

【請求項 11】

前記追加のキュー・ペアに関連する前記追加の処理エンジンの動的挿入を可能にするために、前記ネットワーク・アダプタ内の前記受信側スケーリング（RSS）メカニズムをプログラムするための前記命令は、前記プロセッサに、
前記初期キュー・ペアへの送信タプル・ハッシングを有効化すること
をさらに行わせる、請求項 9 に記載の装置。

【請求項 12】

前記データ処理システムの前記ワークロードが前記所定の高閾値より高くなったかどうかを判断するための前記命令は、データ・フローおよびリソース利用可能性を介して前記データ処理システムの前記ワークロードを監視することを前記プロセッサによって実行される、請求項 9 に記載の装置。

【請求項 13】

前記命令は、前記ワークロードに関連する少なくとも 1 つのパラメータを監視することを前記プロセッサにさらに行わせ、前記少なくとも 1 つのパラメータは、送信 / 受信バイト・パー・セカンド、前記ネットワーク・アダプタによって送信および受信されているフロー制御フレームの数、前記ネットワーク・アダプタによって検出される DMA オーバーランの数、前記デバイス・ドライバによって検出される送信タイムアウト・イベントの数、割り込み毎に前記デバイス・ドライバによって処理される受信パケットの数、またはソフトウェア・キュー上の送信パケットの数のうちの少なくとも 1 つである、請求項 12 に記載の装置。

【請求項 14】

前記命令は、前記プロセッサに、
前記ワークロードが所定の低閾値を下回ったかどうかを判断すること、
前記ワークロードが前記所定の低閾値を下回るのに応じて、前記メモリ内に割り当てられているキュー・ペアが残り 1 つのみあるかどうかを判断すること、
前記メモリ内に割り当てられているキュー・ペアが 2 つ以上残っていることに応じて、割り当てられているキュー・ペアの削除を可能にするために前記ネットワーク・アダプタ内の前記 RSS メカニズムを再プログラムすること、
特定されたキュー・ペアへの送信タプル・ハッシングを無効化すること、
前記特定されたキュー・ペアへの前記ワークロードが休止したかどうかを判断すること

10

20

30

40

50

、および、

前記特定されたキュー・ペアへの前記ワークロードが休止することに応じて、前記特定されたキュー・ペアをメモリから除去し、その結果、前記特定されたキュー・ペアによって使用されたメモリを解放すること

をさらに行わせる、請求項 9 に記載の装置。

【請求項 15】

前記命令は、前記プロセッサに、

前記特定されたキュー・ペアをメモリから除去し、その結果、前記特定されたキュー・ペアによって使用されたメモリを解放することより前に、前記特定されたキュー・ペアへの前記ワークロードが休止するのを、前記特定されたキュー・ペアへの前記ワークロードが休止しないことに応じて待つこと

をさらに行わせる、請求項 14 に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本願は、全般的に、改善されたデータ処理装置および方法に関し、特に、動的再構成 (dynamic reconfiguration) によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供するメカニズムに関する。

【背景技術】

【0002】

イーサネット (R) 媒体の速度が上がり続けるにつれて、可能な限り低いレイテンシで理論上最大限のパフォーマンスを実現するために、より多くのシステム・リソースを使用する必要性が高まっている。イーサネット (R) ・アダプタ要件に関するシステム・リソースには、多数の送信 / 受信ディスクリプタおよびバッファの必要性から、大きな物理メモリ・フットプリントおよび対応する直接メモリ・アクセス入出力メモリ・マッピング・リソースが含まれる。10 Gbps イーサネット (R) ・ドライバは、典型的には、アダプタ毎に約 150 ~ 300 MB の物理システム・メモリおよび直接メモリ・アクセス入出力メモリを消費する。

【発明の概要】

【発明が解決しようとする課題】

【0003】

従来のドライバ・モデルのもとでは、デバイス・ドライバは、アダプタがその理論上最大限のパフォーマンスを実現できるであろう量に送信 / 受信リソースの割り当てを行う。しかし、最大パフォーマンス限界が必要とされないようなワークロード (workload) またはネットワーク・トラフィックであれば、ドライバは必要以上に多くのリソースを消費していることになり、これはシステム・リソースの無駄である。さらに、このモデルには、変化するワークロードに効率的に対処する能力がない。

【課題を解決するための手段】

【0004】

一例示実施形態では、データ処理システムにおいて、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供する方法が提供される。例示実施形態は、メモリ内に初期キュー・ペア (initial queue pair) を割り当てる。例示実施形態は、データ処理システムのワークロードが所定の高閾値より高くなったかどうかを判断する。例示実施形態は、ワークロードが所定の高閾値より高くなることに応じて、メモリ内に追加のキュー・ペアを割り当て、初期化する。例示実施形態は、追加のキュー・ペアに関連する追加の処理エンジン (processing engine) の動的挿入 (dynamic insertion) を可能にするために、ネットワーク・アダプタ内の受信側スケーリング (RSS: receive side scaling) メカニズムをプログラムする。例示実施形態は、追加のキュー・ペアへの送信タプル・ハッシング (transmit tuple hashing) を有効化する

10

20

30

40

50

。

【0005】

別の例示実施形態では、コンピュータ可読プログラムを有するコンピュータ使用可能または可読媒体を含むコンピュータ・プログラム製品が提供される。コンピュータ可読プログラムは、コンピューティング・デバイス上で実行されると、方法の例示実施形態に関して上記に要点を記載した動作のうちの様々な動作、およびその組み合わせを、コンピューティング・デバイスに実行させる。

【0006】

さらに別の例示実施形態では、システム／装置が提供される。システム／装置は、1つ以上のプロセッサと、該1つ以上のプロセッサに結合されたメモリとを含むとよい。メモリは、該1つ以上のプロセッサにより実行されると、方法の例示実施形態に関して上記に要点を記載した動作のうちの様々な動作、およびその組み合わせを該1つ以上のプロセッサに実行させる、命令を含む。

10

【0007】

本発明のこれらの特徴および利点、ならびにその他の特徴および利点が、以下の本発明の例示の実施形態の詳細な説明に記載されるか、またはそれを考慮することで当業者には明らかとなる。

【0008】

例示実施形態の以下の詳細な説明を添付の図面と併せて読み、参照することで、本発明、ならびにその好適な使用方法、ならびにさらなる目的および利点が、最も深く理解される。

20

【図面の簡単な説明】

【0009】

【図1】例示実施形態の各側面が実装され得る例示の分散型データ処理システムの図的表現を示す。

【図2】例示実施形態の側面がともに有利に利用され得るデータ処理システムのブロック図を示す。

【図3】例示実施形態が実装され得る、例示的な論理分割されたプラットフォームのブロック図を示す。

【図4】例示実施形態による、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供するメカニズムの例示的な実装を示す。

30

【図5】例示実施形態による、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供する例示の動作の概要を示すフローチャートを提供する。

【発明を実施するための形態】

【0010】

例示実施形態は、基礎をなすハードウェアの動的再構成によって、アクティブ・メモリ・シェアリング (AMS: active memory sharing) および中央処理ユニット (CPU: central processing unit) 利用のためのリソース・アフィニティを提供するメカニズムを提供し、パフォーマンスまたはサービス
40
の中断なしに、変化するワークロードの要求に応じる。現代のアダプタは、最大限のパフォーマンスのために複数のパケット・キュー・ペア (QP: queue pair) を提供することもある。こうしたアダプタは、インターフェイス毎に複数の送信／受信キュー (QP) を使用することによって、並列ネットワーク・データ処理を実行できることもあり、これは、高トランザクション・ワークロード、および小さなパケット・サイズでより高速な回線速度を実現するためには必須の機能である。イングレス、すなわち受信トラフィックは、アダプタにより、オペレーティング・システム処理のために、適切なQPおよび関連する割り込みにハッシュされるタプルとされ得る。イーグレス、すなわち送信トラフィックは、アダプタに配信するためにオペレーティング・システム (OS: operating system) ドライバによってハッシュされるタプルとされ得る。アダプタ
50

およびOSドライバは、典型的には、最大限のパフォーマンスを実現するために、十分なディスクリプタおよびバッファとともに複数のQP、通常はアダプタ毎に約250MBの平均メモリ・フットプリントに対し2~4つのQPを割り当てる。各QPが、関連する受信割り込みを有することもあるため、複数のQPが使用されていて、トラフィックが少なく、単一のQPで容易に対処可能であるのに、増加した割り込み送付が原因で追加のCPU利用オーバーヘッドが存在することになる。複数のQPは、多くの通常の使用事例に関するパフォーマンスには負の影響を有するということが、しかし特定の高ストレス・高トラザクシオン・ワークロードに関しては、理論上最大限のパフォーマンスを実現するために必須であるということが、既知のアダプタの分析により指摘されているため、この問題に対するパフォーマンス・チームの関心は次第に高まっている。

10

【0011】

したがって、例示実施形態は、分散型データ処理環境、単一データ処理デバイス、または同様のものを含む、多数の異なるタイプのデータ処理環境において利用され得る。例示実施形態の具体的な構成要素および機能性を説明するための背景を提供するために、以下、図1~3が、例示実施形態の各側面が実装されるとよい例示の環境として提供される。図1~3に続く説明は、主として、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供するメカニズムの単一データ処理デバイス実装に焦点を当てるが、これは例でしかなく、本発明の特徴に関していかなる制限を記載することも、示唆することも目的としていない。反対に、例示実施形態は、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティが提供されるとよい分散型データ処理の環境および実施形態を含むものとする。

20

【0012】

以下、図面、特に図1~3を参照する。本発明の例示実施形態が実装され得るデータ処理環境の例示の図が提供されている。当然のことながら、図1~3は例でしかなく、本発明の側面または実施形態が実装され得る環境に関していかなる制限を主張することも、示唆することも目的としていない。示される環境に対して多数の変更が、本発明の意図および範囲から逸脱することなく加えられ得る。

【0013】

以下、図面を参照する。図1は、例示実施形態の各側面が実装され得る例示の分散型データ処理システムの図的表現を示す。分散型データ処理システム100は、例示実施形態の各側面が実装され得るコンピュータのネットワークを含むとよい。分散型データ処理システム100は、少なくとも1つのネットワーク102を含み、ネットワーク102は、分散型データ処理システム100内で互いに接続された様々なデバイスとコンピュータとの通信リンクを提供するために使用される媒体である。ネットワーク102は、有線、無線通信リンク、または光ファイバ・ケーブルなどの接続を含み得る。

30

【0014】

示されている例では、サーバ104およびサーバ106が、ストレージ・ユニット108とともに、ネットワーク102に接続されている。さらに、クライアント110、112、および114もネットワーク102に接続されている。これらのクライアント110、112、および114は、例えばパーソナル・コンピュータ、ネットワーク・コンピュータ、または同様のものであるとよい。示されている例では、サーバ104が、ブート・ファイル、オペレーティング・システム・イメージ、およびアプリケーションなどのデータを、クライアント110、112、および114に提供する。クライアント110、112、および114は、示されている例では、サーバ104のクライアントである。分散型データ処理システム100は、示されていない追加のサーバ、クライアント、およびその他のデバイスを含んでもよい。

40

【0015】

示されている例では、分散型データ処理システム100は、インターネットであり、ネットワーク102は、互いに通信するために伝送制御プロトコル/インターネット・プロトコル(TCP/IP: Transmission Control Protocol

50

/ Internet Protocol) プロトコル・スイートを使用する、ネットワークおよびゲートウェイの世界的規模の集合を表現している。インターネットの中心には、データおよびメッセージをルーティングする数千の商用、政府機関用、教育用、およびその他のコンピュータ・システムから成る主要ノードまたはホスト・コンピュータ間の高速データ通信回線のバックボーンがある。当然、分散型データ処理システム100は、例えばイントラネット、ローカル・エリア・ネットワーク(LAN: local area network)、広域ネットワーク(WAN: wide area network)、または同様のものなど、いくつかの異なるタイプのネットワークを含むようにも実装され得る。上記のように、図1は、本発明の種々の実施形態のアーキテクチャ上の制限としてではなく、例として意図されており、したがって、図1に示されている特定の構成要素は、本発明の例示実施形態が実装され得る環境に関する制限的なものと見なされてはならない。

10

【0016】

例示実施形態では、コンピュータ・アーキテクチャは、ハードウェアおよびソフトウェアの組み合わせとして実装される。コンピュータ・アーキテクチャのソフトウェア部分は、マイクロコードまたはミリコードと呼ばれることもある。ハードウェアおよびソフトウェアの組み合わせは、命令セットおよびシステム・アーキテクチャを作成し、これに基づき、基本入出力システム(BIOS: Basic Input/Output System)、仮想マシン・モニタ(VMM: Virtual Machine Monitors)、ハイパーバイザ、アプリケーションなどのコンピュータのソフトウェアの残りが動作する。初期の組み合わせによって作成されるコンピュータ・アーキテクチャは、ごく少数しかないと考えられる定義済みインターフェイスを介する場合を除き、コンピュータ・ソフトウェア(BIOSなど)に対し不変である。

20

【0017】

以下、図2を参照する。例示実施形態の側面がともに有利に利用され得るデータ処理システムのブロック図が示されている。図のように、データ処理システム200は、プロセッサ・ユニット211a~211nを含む。プロセッサ・ユニット211a~211nはそれぞれ、プロセッサおよびキャッシュ・メモリを含む。例えば、プロセッサ・ユニット211aは、プロセッサ212aおよびキャッシュ・メモリ213aを含み、プロセッサ・ユニット211nは、プロセッサ212nおよびキャッシュ・メモリ213nを含む。

30

【0018】

プロセッサ・ユニット211a~211nは、メイン・バス215に接続されている。メイン・バス215は、プロセッサ・ユニット211a~211nおよびメモリ・カード223を含むシステム・プレーナ220をサポートしている。システム・プレーナ220はさらに、データ・スイッチ221およびメモリ・コントローラ/キャッシュ222を含む。メモリ・コントローラ/キャッシュ222は、複数のデュアル・インライン・メモリ・モジュール(DIMM: dual in-line memory module)を有するローカル・メモリ216を含むメモリ・カード223をサポートしている。

【0019】

データ・スイッチ221は、ネイティブI/O(NIO: native I/O)プレーナ224内に位置するバス・ブリッジ217およびバス・ブリッジ218と接続している。図のように、バス・ブリッジ218は、システム・バス219を経由して、ペリフェラル・コンポーネント・インターコネクト(PCI: peripheral components interconnect)ブリッジ225および226と接続している。PCIブリッジ225は、PCIバス228を経由して様々なI/Oデバイスと接続している。図のように、ハード・ディスク236が、小型コンピュータ・システム・インターフェイス(SCSI: small computer system interface)ホスト・アダプタ230を経由してPCIバス228に接続されているとよい。グラフィックス・アダプタ231が、直接的または間接的にPCIバス228に接続されているとよい。PCIブリッジ226は、PCIバス227を経由して、ネットワーク・アダ

40

50

プタ234およびアダプタ・カード・スロット235a~235nを介して外部データ・ストリームに対する接続を提供する。

【0020】

業界標準アーキテクチャ (ISA: Industry standard architecture) バス229が、ISAブリッジ232を経由してPCIバス228と接続している。ISAブリッジ232は、シリアル接続シリアル1およびシリアル2を有するNIOコントローラ233を介して相互接続機能を提供する。データ処理システム200が対応する入力デバイスを経由してユーザからのデータ入力を受け取ることができるために、フレキシブル・ディスク・ドライブ接続、キーボード接続、およびマウス接続が、NIOコントローラ233によって提供されている。さらに、ISAバス229に接続されている不揮発性RAM (NVRAM: non-volatile random-access memory) 240は、電力供給問題などのシステム障害またはシステム故障から特定のタイプのデータを保護する不揮発性メモリを提供する。システム・ファームウェア241もISAバス229に接続されており、初期基本入出力システム (BIOS) 機能を実装する。サービス・プロセッサ244が、ISAバス229に接続しており、システム診断またはシステム・サービスの機能性を提供する。

10

【0021】

オペレーティング・システム (OS) は、ハード・ディスク236に記憶され、ハード・ディスク236はさらに、データ処理システムによって実行されるさらなるアプリケーション・ソフトウェアのストレージを提供するとよい。NVRAM240は、フィールド交換可能ユニット (FRU: field replaceable unit) の分離のために、システム変数およびエラー情報を記憶するように使用される。システム起動中、ブートストラップ・プログラムが、オペレーティング・システムをロードし、オペレーティング・システムの実行を開始する。オペレーティング・システムをロードするために、ブートストラップ・プログラムはまず、オペレーティング・システム・カーネル・イメージをハード・ディスク236上に置き、OSカーネル・イメージをメモリ内にロードし、オペレーティング・システム・カーネルによって提供される先頭アドレスにジャンプする。典型的には、オペレーティング・システムは、データ処理システム内のランダム・アクセス・メモリ (RAM: random-access memory) 内にロードされる。ロードおよび初期化が行われると、オペレーティング・システムは、プログラムの実行を制御し、リソース割り当て、スケジューリング、入出力制御、およびデータ管理などのサービスを提供するとよい。

20

30

【0022】

例示実施形態は、いくつかの異なるハードウェア構成、ならびにブートストラップ・プログラムおよびオペレーティング・システムなどのソフトウェアを利用する、様々なデータ処理システムにおいて具現化され得る。データ処理システム200は、例えば、スタンダードアロン・システム、またはローカル・エリア・ネットワーク (LAN) もしくは広域ネットワーク (WAN) などのネットワークの一部としてもよい。上記のように、図2は、本発明の種々の実施形態のアーキテクチャ上の制限としてではなく、例として意図されており、したがって、図2に示されている特定の構成要素は、本発明の例示実施形態が実装され得る環境に関する制限的なものと見なされてはならない。

40

【0023】

以下、図3を参照する。例示実施形態が実装され得る、例示的な論理分割されたプラットフォームのブロック図が示されている。論理分割されたプラットフォーム300内のハードウェアは、例えば、図2のデータ処理システム200のハードウェアを使用して実装されてもよい。

【0024】

論理分割されたプラットフォーム300は、分割されたハードウェア330、オペレーティング・システム302、304、306、308、および仮想マシン・モニタ310を含む。オペレーティング・システム302、304、306、および308は、単一の

50

オペレーティング・システムの複数のコピーであってもよく、または論理分割されたプラットフォーム 300 上で同時に実行する複数の異種オペレーティング・システムであってもよい。これらのオペレーティング・システムは、例えば、ハイパーバイザなどのパーティション管理ファームウェアなど、仮想化メカニズムとインターフェイスをとるよう設計されている、z / OS (IBM 社の登録商標) を使用して実装されてもよい。z / OS は、これらの例示実施形態において単なる例として使用されている。当然、特定の実装に応じて、OS / 400 (IBM 社の登録商標)、AIX (R)、および Linux (R) など、他のタイプのオペレーティング・システムが使用されてもよい。オペレーティング・システム 302、304、306、および 308 は、それぞれ論理パーティション 303、305、307、および 309 に位置する。

10

【0025】

ハイパーバイザ・ソフトウェアは、プラットフォームの実装に使用され得るソフトウェアの例であり(この例では、仮想マシン・モニタ 310)、インターナショナル・ビジネス・マシーンズ・コーポレーションから入手できる。ファームウェアは、例えば読み取り専用メモリ (ROM: read-only memory)、プログラム可能 ROM (PROM: programmable ROM)、消去可能プログラム可能 ROM (EPROM: erasable programmable ROM)、および電氣的消去可能プログラム可能 ROM (EEPROM: electrically erasable programmable ROM) など、電力がなくてもそのコンテンツを保持するメモリ・チップに記憶された、「ソフトウェア」である。

20

【0026】

論理分割されたプラットフォーム 300 はさらに、システム・メモリ仮想化機能を IBM Power Systems に提供する IBM (R) Power VM (IBM 社の登録商標) の高度なメモリ仮想化技術である、IBM (R) の Power VM (IBM 社の登録商標) Active Memory (TM) Sharing を使用して、複数の論理パーティションが物理メモリの共通のプールを共有できるようにすることもできる。IBM (R) Power Systems (TM) の物理メモリは、専用または共有モードのいずれかで、複数の論理パーティションに割り振られることが可能である。システム管理者は、一部の物理メモリをある論理パーティションに割り振り、一部の物理メモリを他の複数の論理パーティションにより共有されるプールに割り振る能力を有する。単一のパーティションは、専用または共有メモリのいずれかを有し得る。Active Memory (TM) Sharing は、必要システム・メモリを減らすこと、または既存システム上での追加の論理パーティションの作成を可能にすることによって、システム上でのメモリ利用度を上げるよう利用され得る。

30

【0027】

論理パーティション 303、305、307、および 309 はさらに、パーティション・ファームウェア・ローダ 311、313、315、および 317 を含む。パーティション・ファームウェア・ローダ 311、313、315、および 317 は、IPL または初期ブートストラップ・コード、IEEE-1275 標準 Open Firmware、およびインターナショナル・ビジネス・マシーンズ・コーポレーションから入手できるランタイム抽象化ソフトウェア (RTAS: runtime abstraction software) を使用して実装され得る。

40

【0028】

論理パーティション 303、305、307、および 309 がインスタンス化されると、ブートストラップ・コードのコピーが、仮想マシン・モニタ 310 によって論理パーティション 303、305、307、および 309 内へロードされる。その後、ブートストラップ・コードに制御が移され、次にブートストラップ・コードは、open firmware および RTAS をロードする。続いて、論理パーティション 303、305、307、および 309 に関連する、または割り振られているプロセッサが、論理パーティション・ファームウェアを実行するよう論理パーティションのメモリにディスパッチされる

50

。

【 0 0 2 9 】

分割されたハードウェア 3 3 0 は、複数のプロセッサ 3 3 2 ~ 3 3 8、複数のシステム・メモリ・ユニット 3 4 0 ~ 3 4 6、複数の入出力 (I / O : i n p u t / o u t p u t) アダプタ 3 4 8 ~ 3 6 2、およびストレージ・ユニット 3 7 0 を含む。プロセッサ 3 3 2 ~ 3 3 8、メモリ・ユニット 3 4 0 ~ 3 4 6、N V R A M ストレージ 3 9 8、および I / O アダプタ 3 4 8 ~ 3 6 2 はそれぞれ、論理分割されたプラットフォーム 3 0 0 内の複数の論理パーティション 3 0 3、3 0 5、3 0 7、および 3 0 9 のうちの 1 つに割り振られるとよく、それぞれがオペレーティング・システム 3 0 2、3 0 4、3 0 6、および 3 0 8 のうちの 1 つに対応する。

10

【 0 0 3 0 】

仮想マシン・モニタ 3 1 0 は、論理パーティション 3 0 3、3 0 5、3 0 7、および 3 0 9 に関するいくつかの機能およびサービスを実行して、論理分割されたプラットフォーム 3 0 0 の分割をもたらし、実行する。仮想マシン・モニタ 3 1 0 は、基礎をなすハードウェアと等しい、ファームウェアで実装される仮想マシンである。したがって、仮想マシン・モニタ 3 1 0 は、独立した O S イメージ 3 0 2、3 0 4、3 0 6、および 3 0 8 の同時実行を、論理分割されたプラットフォーム 3 0 0 のすべてのハードウェア・リソースを仮想化することによって可能にする。

【 0 0 3 1 】

サービス・プロセッサ 3 9 0 は、論理パーティション 3 0 3、3 0 5、3 0 7、および 3 0 9 におけるプラットフォーム・エラーの処理などの様々なサービスを提供するために使用され得る。サービス・プロセッサ 3 9 0 はさらに、インターナショナル・ビジネス・マシーンズ・コーポレーションなどのベンダにエラーのレポートを返すサービス・エージェントとしての機能を果たしてもよい。種々の論理パーティションの動作は、ハードウェア・システム・コンソール 3 8 0 を介して制御され得る。ハードウェア・システム・コンソール 3 8 0 は、別個のデータ処理システムであり、このシステムからシステム管理者は、リソースを異なる論理パーティションに再割り当てすることを含む、様々な機能を実行することもできる。

20

【 0 0 3 2 】

例示実施形態は、通常動作を実現するために最低限必要なリソースを少し上回るリソースを用いて、単一のキュー・ペア (Q P)、すなわち送信 / 受信ペアのみを最初に割り当てるよう、オペレーティング・システム (O S) ドライバを規定する。トラフィック・フローまたはワークロードが所定の閾値超に増えるのに従って、O S ドライバは、必要に応じて追加の Q P を動的に割り当てる。Q P が実行中の O S ドライバに追加され、アダプタに対して利用可能になるため、トラフィックは、イングレス、すなわち受信パス、およびイーグレス、すなわち送信パスの両方において、より多くの中央処理ユニット (C P U) にハッシュされ、パフォーマンスおよび C P U / メモリ・リソース使用量が効果的に拡大される。パフォーマンス拡大は、静的なディスクリプタ・カウントではなく追加の Q P によって実現されるため、システム・リソースは、必要なくなれば削減されることもできる。トラフィック・フローおよびワークロードが、規定の最小閾値未満に下がるのに従って、O S ドライバは、Q P を除去し、通常動作の最小限のリソース利用まで再び低下するとよい。ワークロードまたはトラフィックが増え、このサイクルが繰り返される。したがって、例示実施形態は、C P U 利用およびアクティブ・メモリ・シェアリングのアフィニティを提供しながら、パフォーマンスのために調整する動的再構成を効果的に実現する。

30

40

【 0 0 3 3 】

図 4 は、例示実施形態による、主要な動作コンポーネントおよびその相互作用を示す例示のブロック図である。図 4 に示されている構成要素は、ハードウェア、ソフトウェア、またはハードウェアおよびソフトウェアの任意の組み合わせにおいて実装され得る。一例示実施形態では、図 4 の構成要素は、1 つ以上のデータ処理デバイスまたはシステムの 1

50

つ以上のプロセッサ上で実行されているソフトウェアとして実装される。

【 0 0 3 4 】

図 4 は、例示実施形態による、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供するメカニズムの例示的な実装を示す。データ処理システム 4 0 0 は、オペレーティング・システム 4 0 4 内にデバイス・ドライバ 4 0 2 を含む。デバイス・ドライバ 4 0 2 は、アプリケーション 4 1 0 およびネットワーク・アダプタ 4 1 2 が使用するメモリ 4 0 8 内に 1 つ以上のキュー・ペア 4 0 6 a ~ 4 0 6 n を提供する。例示実施形態は、1 つのデバイス・ドライバ 4 0 2、1 つのアプリケーション 4 1 0、および 1 つのネットワーク・アダプタ 4 1 2 のみを示すが、当業者には当然のことながら、データ処理システム 4 0 0 は、複数のデバイス・ドライバ、複数のアプリケーション、および複数のネットワーク・アダプタをデータ処理システム 4 0 0 内に含むこともできる。

10

【 0 0 3 5 】

オペレーティング・システム 4 0 4 の初期化時、オペレーティング・システム 4 0 4 は、デバイス・ドライバ 4 0 2 を構成および初期化する。次に、デバイス・ドライバ 4 0 2 は、メモリ 4 0 8 内に、いくつかの受信ディスクリプタ/バッファおよびいくつかの送信ディスクリプタ/バッファを含むとよい初期キュー・ペア 4 0 6 a を割り当てる。続いてデバイス・ドライバ 4 0 2 は、トラフィックがネットワーク・アダプタ 4 1 2 に送られることができるように、ネットワーク・アダプタ 4 1 2 を起動する。動作中、キュー・ペア 4 0 6 a のみが初期化されているとき、デバイス・ドライバ 4 0 2 は、ネットワーク・アダプタ 4 1 2 へ送信されるトラフィックを受信し、トラフィックを、配信のためにネットワーク・アダプタ 4 1 2 上へ送る。続いてネットワーク・アダプタ 4 1 2 は、送信トラフィックを、サーバ 4 1 6、サーバ 4 1 8、クライアント 4 2 0、クライアント 4 2 2、または同様のものなどのデバイス上へ、ネットワーク 4 2 4 を経由して送る。反対に、ネットワーク・アダプタ 4 1 2 が、サーバ 4 1 6、サーバ 4 1 8、クライアント 4 2 0、クライアント 4 2 2、または同様のものなどのデバイスからデバイス・ドライバ 4 0 2 にネットワーク 4 2 4 を経由して送信されるトラフィックを受信すると、ネットワーク・アダプタ 4 1 2 は、オペレーティング・システム 4 0 4 による処理の割り込みを起こし、トラフィックは、デバイス・ドライバ 4 0 2 に送られる。キュー・ペア 4 0 6 a、および下記のようにさらに割り当てられるとよい 4 0 6 b ~ 4 0 6 n はそれぞれ、それ自体の、関連する受信割り込みを有する。

20

30

【 0 0 3 6 】

デバイス・ドライバ 4 0 2 がキュー・ペア 4 0 6 a を割り当ててインスタンス化した後、デバイス・ドライバ 4 0 2 は、データ・フローおよびリソース利用可能性を介して、オペレーティング・システム 4 0 4 のワークロードを継続的に監視する。デバイス・ドライバ 4 0 2 は、送信/受信バイト・パー・セカンド、ネットワーク・アダプタ 4 1 2 によって送信および受信されているフロー制御フレームの数、ネットワーク・アダプタ 4 1 2 によって検出される DMA オーバーラン (DMA over run) の数、デバイス・ドライバ 4 0 2 によって検出される送信タイムアウト・イベントの数、割り込み毎にデバイス・ドライバ 4 0 2 によって処理される受信パケットの数、ソフトウェア・キュー上の送信パケットの数、または同様のものなどのパラメータを監視するとよい。デバイス・ドライバ 4 0 2 が、所定の高閾値を上回ることなどによって、ワークロード監視によりリソース不足の状況を検出すると、デバイス・ドライバ 4 0 2 は、キュー・ペア 4 0 6 b ~ 4 0 6 n のうちの追加のキュー・ペアを動的に割り当てて初期化する。続いてデバイス・ドライバ 4 0 2 は、ネットワーク・アダプタ 4 1 2 内の受信側スケールリング (RSS) メカニズム 4 1 4 をプログラムし、キュー・ペア 4 0 6 a およびキュー・ペア 4 0 6 b ~ 4 0 6 n のうちの追加のキュー・ペアに関連する追加の処理エンジンの動的挿入を可能にする。続いてデバイス・ドライバ 4 0 2 は、キュー・ペア 4 0 6 a およびキュー・ペア 4 0 6 b ~ 4 0 6 n のうちの追加のキュー・ペアへの、送信タプル・ハッシングを有効化する。ネットワーク・アダプタ 4 1 2 による受信タプル・ハッシングおよび処理は、RSS メカニズ

40

50

ム 4 1 4 をプログラムすることにより自動的に有効化される。デバイス・ドライバ 4 0 2 は、キュー・ペア 4 0 6 b ~ 4 0 6 n のうちの利用可能なキュー・ペアがすべて消費されるまで、またはキュー・ペア 4 0 6 a ~ 4 0 6 n がデータ処理システム 4 0 0 内の中央処理ユニットの数を上回るまで、キュー・ペア 4 0 6 b ~ 4 0 6 n のうちのキュー・ペアを、まだ割り当ておよび初期化されていなければ、ワークロードによる必要に応じて追加し続ける。デバイス・ドライバ 4 0 2 はさらに、ネットワーク・アダプタ 4 1 2 内の R S S メカニズム 4 1 4 を再プログラムし、キュー・ペア 4 0 6 b ~ 4 0 6 n のうちの新たなキュー・ペアが追加されるたびに、追加の処理エンジンの動的挿入を可能にし、さらにキュー・ペア 4 0 6 b ~ 4 0 6 n のうちの新たなキュー・ペアへの送信タプル・ハッシングを有効化する。

10

【 0 0 3 7 】

デバイス・ドライバ 4 0 2 が、ワークロードが所定の低閾値を下回ることなどによってワークロードの減少を認識すると、デバイス・ドライバ 4 0 2 は、ネットワーク・アダプタ 4 1 2 内の R S S メカニズム 4 1 4 を動的に再プログラムし、キュー・ペア 4 0 6 a ~ 4 0 6 n のうちの割り当てられているキュー・ペアの削除を可能にする。続いてデバイス・ドライバ 4 0 2 は、キュー・ペア 4 0 6 a ~ 4 0 6 n のうちの削除されたキュー・ペアへの送信タプル・ハッシングを無効化する。キュー・ペア 4 0 6 a ~ 4 0 6 n のうちの削除されたキュー・ペアが休止する (q u i e s c e) と、デバイス・ドライバ 4 0 2 は、キュー・ペア 4 0 6 a ~ 4 0 6 n のうちの削除されたキュー・ペアを除去し、その結果、キュー・ペア 4 0 6 a ~ 4 0 6 n のうちの削除されたキュー・ペアによって使用されたメモリが解放される。ネットワーク・アダプタ 4 1 2 における受信タプル・ハッシングの有効化のように、ネットワーク・アダプタ 4 1 2 による受信タプル・ハッシングおよび処理は、R S S メカニズム 4 1 4 を再プログラムすることにより自動的に無効化される。

20

【 0 0 3 8 】

したがって、例示実施形態は、基礎をなすハードウェアの動的再構成によって、アクティブ・メモリ・シェアリング (A M S) および中央処理ユニット (C P U) 利用のためのリソース・アフィニティを提供するメカニズムを提供し、パフォーマンスまたはサービスの中断なしに、変化するワークロードの要求に応じる。トラフィック・フローまたはワークロードが所定の閾値超に増えるのに従って、デバイス・ドライバは、必要に応じて追加のキュー・ペアを動的に割り当てる。トラフィック・フローおよびワークロードが所定の最小閾値未満に下がるのに従って、デバイス・ドライバは、キュー・ペアを除去し、通常動作の最小限のリソース利用まで再び低下するとよい。ワークロードまたはトラフィックが増加し、続いて減少すると、このサイクルが繰り返される。

30

【 0 0 3 9 】

当業者であれば当然のことであるが、本発明は、システム、方法、またはコンピュータ・プログラム製品として具現化され得る。したがって、本発明の側面は、完全にハードウェアの実施形態、完全にソフトウェアの実施形態 (ファームウェア、常駐ソフトウェア、マイクロコードなどを含む)、または本願明細書においてすべて概して「回路」、「モジュール」もしくは「システム」と呼ばれ得る、ソフトウェアおよびハードウェアの側面を兼ね備えた実施形態の形態をとり得る。さらに、本発明の側面は、コンピュータ使用可能プログラム・コードが具現化された任意の 1 つ以上のコンピュータ可読媒体 (単数または複数) において具現化されたコンピュータ・プログラム製品の形態をとることもできる。

40

【 0 0 4 0 】

1 つ以上のコンピュータ可読媒体 (単数または複数) の任意の組み合わせが利用され得る。コンピュータ可読媒体は、コンピュータ可読信号媒体またはコンピュータ可読ストレージ媒体とされ得る。コンピュータ可読ストレージ媒体は、例えば、限定はされないが、電子、磁気、光学、電磁気、赤外線、もしくは半導体のシステム、装置、デバイス、または前述のものの任意の適切な組み合わせとされ得る。コンピュータ可読媒体のより具体的な例 (包括的でないリスト) には、1 つ以上のワイヤを有する電氣的接続、ポータブル・コンピュータ・ディスク、ハード・ディスク、ランダム・アクセス・メモリ (R A M

50

）、読み取り専用メモリ（ROM）、消去可能プログラム可能読み取り専用メモリ（EPROMまたはフラッシュ・メモリ）、光ファイバ、ポータブル・コンパクト・ディスク読み取り専用メモリ（CDROM: compact disc read-only memory）、光学式ストレージ・デバイス、磁気ストレージ・デバイス、または前述のものの任意の適切な組み合わせが含まれるであろう。この文書の文脈では、コンピュータ可読ストレージ媒体は、命令実行システム、装置、もしくはデバイスによって、またはそれに関連して使用されるプログラムを含むこと、または記憶することができる任意の有形の媒体とされ得る。

【0041】

コンピュータ可読信号媒体は、例えば、ベースバンドに、または搬送波の一部として、コンピュータ可読プログラム・コードが具現化された伝播データ信号を含み得る。そのような伝播信号は、電磁気、光学、またはその任意の適切な組み合わせを含むがこれらに限定はされない、様々な形態のいずれかをとりよい。コンピュータ可読信号媒体は、コンピュータ可読ストレージ媒体でなく、命令実行システム、装置もしくはデバイスによって、またはそれに関連して使用されるプログラムの伝達、伝播、または搬送をすることができる、任意のコンピュータ可読媒体としてよい。

【0042】

コンピュータ可読媒体上に具現化されたコンピュータ・コードは、無線、有線、光ファイバ・ケーブル、無線周波数（RF: radio frequency）など、またはその任意の適切な組み合わせを含むがこれらに限定はされない、任意の適切な媒体を使用して送られてもよい。

【0043】

本発明の側面の動作を実行するコンピュータ・プログラム・コードは、Java（R）、Smalltalk（R）、C++または同様のものなどのオブジェクト指向プログラミング言語、および「C」プログラミング言語もしくは同様のプログラミング言語などの従来の手続きプログラミング言語を含む、1つ以上のプログラミング言語の任意の組み合わせで書かれていてよい。プログラム・コードは、スタンド・アロン・ソフトウェア・パッケージとして、完全にユーザのコンピュータ上で実行されることも、部分的にユーザのコンピュータ上で実行されることも、または部分的にユーザのコンピュータ上で、かつ部分的にリモート・コンピュータ上で実行されることも、または完全にリモート・コンピュータもしくはサーバ上で実行されることもできる。後者のシナリオでは、ローカル・エリア・ネットワーク（LAN）もしくは広域ネットワーク（WAN）を含む任意の種類のネットワークを介してリモート・コンピュータがユーザのコンピュータに接続されてもよく、または、（例えば、インターネット・サービス・プロバイダを使用しインターネットを介して）外部コンピュータに接続されてもよい。

【0044】

本発明の側面について、本発明の例示実施形態による方法、装置（システム）およびコンピュータ・プログラム製品のフローチャート図もしくはブロック図またはその両方を参照して以下に記載する。当然のことながら、フローチャート図もしくはブロック図またはその両方の各ブロック、およびフローチャート図もしくはブロック図またはその両方の複数ブロックの組み合わせは、コンピュータ・プログラム命令により実装可能である。マシンを生じるよう、こうしたコンピュータ・プログラム命令が、汎用コンピュータ、専用コンピュータ、またはその他のプログラム可能データ処理装置のプロセッサに提供されて、この命令が、コンピュータまたはその他のプログラム可能データ処理装置のプロセッサにより実行されて、フローチャートもしくはブロック図またはその両方のブロックもしくは複数ブロックにおいて指定された機能／動作を実装する手段を作り出すようにすることもできる。

【0045】

さらに、特定の形で機能するようコンピュータ、その他のプログラム可能データ処理装置、またはその他のデバイスに指示することができるこうしたコンピュータ・プログラム

10

20

30

40

50

命令は、コンピュータ可読媒体に記憶されて、コンピュータ可読媒体に記憶されたこの命令が、フローチャートもしくはブロック図またはその両方のブロックもしくは複数ブロックにおいて指定された機能／動作を実装する命令を含む製品を生じるようにすることもできる。

【 0 0 4 6 】

さらに、コンピュータ・プログラム命令は、コンピュータ、その他のプログラム可能データ処理装置、またはその他のデバイスにロードされて、コンピュータ、その他のプログラム可能装置、またはその他のデバイス上で一連の動作ステップが実行されるようにしてコンピュータで実装されるプロセスを生じさせ、コンピュータまたはその他のプログラム可能装置上で実行される命令が、フローチャートもしくはブロック図またはその両方のブロックもしくは複数ブロックにおいて指定された機能／動作を実装するためのプロセスを提供するようにすることもできる。

10

【 0 0 4 7 】

以下、図5を参照する。この図面は、例示実施形態による、動的再構成によってマルチキュー・ネットワーク・アダプタのリソース・アフィニティを提供する例示の動作の概要を示すフローチャートを提供する。動作が開始するとき、構成および初期化されたデバイス・ドライバが、メモリ内に初期キュー・ペアを割り当てる（ステップ502）。続いてデバイス・ドライバは、トラフィックがネットワーク・アダプタに送られることができるように、ネットワーク・アダプタを起動する（ステップ504）。

【 0 0 4 8 】

20

デバイス・ドライバが、キュー・ペアを割り当ててインスタンス化した後、デバイス・ドライバは、データ・フローおよびリソース利用可能性を介して、オペレーティング・システムのワークロードを継続的に監視する（ステップ506）。デバイス・ドライバは、送信／受信バイト・パー・セカンド、ネットワーク・アダプタによって送信および受信されているフロー制御フレームの数、ネットワーク・アダプタによって検出されるDMAオーバーランの数、デバイス・ドライバによって検出される送信タイムアウト・イベントの数、割り込み毎にデバイス・ドライバによって処理される受信パケットの数、ソフトウェア・キュー上の送信パケットの数、または同様のものなどのパラメータを監視するとよい。続いてデバイス・ドライバは、リソース不足状況を示す、所定の高閾値の超過が生じたかどうかを、ワークロード監視により判断する（ステップ508）。ステップ508にて、ワークロードが所定の高閾値より高くなっていれば、デバイス・ドライバは、メモリ内に追加のキュー・ペアを動的に割り当てて初期化する（ステップ510）。続いてデバイス・ドライバは、ネットワーク・アダプタ内のRSSメカニズムをプログラム／再プログラムして追加の処理エンジンの動的挿入を可能にし（ステップ512）、デバイス・ドライバは、新たに割り当てられたキュー・ペアへの送信タプル・ハッシングを有効化し（ステップ514）、動作はその後、ステップ506に戻る。

30

【 0 0 4 9 】

ステップ508にて、ワークロードが所定の高閾値より高くなっていなければ、デバイス・ドライバは、ワークロードが所定の低閾値を下回ったかどうかを判断する（ステップ516）。ステップ516にて、デバイス・ドライバが、ワークロードは所定の低閾値を下回っていないと判断すると、動作はステップ506に戻る。ステップ516にて、デバイス・ドライバが、ワークロードが所定の低閾値を下回ったと判断すると、デバイス・ドライバは、割り当てられているキュー・ペアが残り1つのみあるかどうかを判断する（ステップ518）。ステップ518にて、デバイス・ドライバが、1つのみのキュー・ペアが残っていると判断すると、動作はステップ506に戻る。ステップ518にて、デバイス・ドライバが、2つ以上のキュー・ペアが残っていると判断すると、デバイス・ドライバは、ネットワーク・アダプタ内のRSSメカニズムを動的に再プログラムして、割り当てられているキュー・ペアの削除を可能にするとよい（ステップ520）。続いてデバイス・ドライバは、特定されたキュー・ペアへの送信タプル・ハッシングを無効化する（ステップ522）。デバイス・ドライバは続いて、特定されたキュー・ペアへのワークロー

40

50

ドが休止したかどうかを判断する（ステップ524）。ステップ524にて、デバイス・ドライバが、特定されたキュー・ペアへのワークロードが休止していないと判断すると、動作はステップ524に戻る。ステップ524にて、デバイス・ドライバが、特定されたキュー・ペアへのワークロードが休止したと判断すると、デバイス・ドライバは、メモリから特定されたキュー・ペアを除去し（ステップ526）、その結果、特定されたキュー・ペアによって使用されたメモリが解放される。動作は続いて、ステップ506に戻る。

【0050】

各図面のフローチャートおよびブロック図は、本発明の様々な実施形態によるシステム、方法およびコンピュータ・プログラム製品の考えられる実装のアーキテクチャ、機能性、および動作を示す。この関連で、フローチャートまたはブロック図内の各ブロックは、指定の論理機能（単数または複数）を実装する1つ以上の実行可能命令を含むモジュール、セグメント、またはコードの一部を表すこともできる。なお、さらに、いくつかの代替の実装では、ブロック内に示されている機能が、図面に示されているのとは異なる順序で生じてよい。例えば、関連する機能性次第で、連続して示されている2つのブロックが実際には事実上同時に実行されてもよく、または、各ブロックが逆順で実行されることがあってもよい。なお、さらに、ブロック図もしくはフローチャート図またはその両方の各ブロック、およびブロック図もしくはフローチャート図またはその両方の複数ブロックの組み合わせは、指定の機能もしくは動作を実行する専用ハードウェア・ベース・システム、または専用ハードウェアおよびコンピュータ命令の組み合わせにより実装することができる。

【0051】

したがって、例示実施形態は、基礎をなすハードウェアの動的再構成によってアクティブ・メモリ・シェアリング（AMS）および中央処理ユニット（CPU）利用のためのリソース・アフィニティを提供するメカニズムを提供し、パフォーマンスまたはサービスの中断なしに、変化するワークロードの要求に応じる。トラフィック・フローまたはワークロードが、所定の閾値超に増えるのに従って、デバイス・ドライバは、必要に応じて追加のキュー・ペアを動的に割り当てる。トラフィック・フローおよびワークロードが、所定の最小閾値未満に下がるのに従って、デバイス・ドライバは、キュー・ペアを除去し、通常動作の最小限のリソース利用まで再び低下するとよい。ワークロードまたはトラフィックが増加し、続いて減少すると、このサイクルが繰り返される。

【0052】

上述のとおり、当然のことながら、例示実施形態は、完全にハードウェアの実施形態、完全にソフトウェアの実施形態、またはハードウェアおよびソフトウェア両方の構成要素を含む実施形態の形態をとり得る。例示の一実施形態では、例示実施形態のメカニズムは、ファームウェア、常駐ソフトウェア、マイクロコードなどを含むがこれらに限定されないソフトウェアまたはプログラム・コードにおいて実装される。

【0053】

プログラム・コードの記憶もしくは実行またはその両方に適したデータ処理システムは、システム・バスを介してメモリ要素に直接または間接的に結合された少なくとも1つのプロセッサを含む。メモリ要素は、プログラム・コードを実際に行う間に用いられるローカル・メモリ、大容量ストレージ、および、実行中にコードが大容量ストレージから読み出されなければならない回数を減らすために少なくとも一部のプログラム・コードの一時的なストレージとなるキャッシュ・メモリを含むことができる。

【0054】

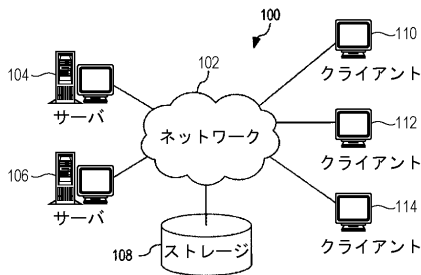
入出力、すなわちI/Oデバイス（限定はされないが、キーボード、ディスプレイ、ポインティング・デバイスなどを含む）は、直接、または介在するI/Oコントローラを介して、システムに結合されることが可能である。ネットワーク・アダプタもシステムに結合されて、データ処理システムが、他のデータ処理システムまたはリモート・プリンタまたはストレージ・デバイスに、介在するプライベートまたはパブリック・ネットワークを介して結合された状態となることを可能にしてもよい。モデム、ケーブル・モデム、およ

びイーサネット（Ｒ）・カードが、現在利用可能なタイプのネットワーク・アダプタのごく一部である。

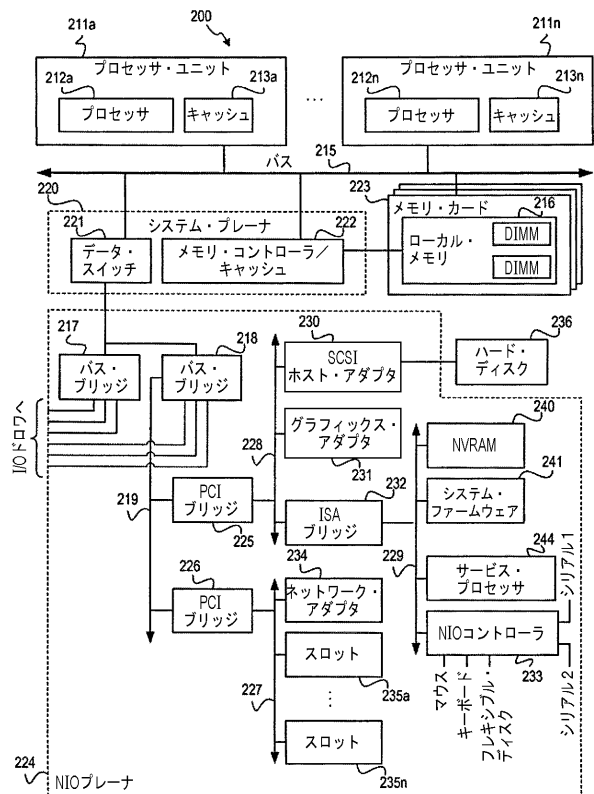
【 0 0 5 5 】

本発明の記載は、例示および説明のために示されたものであり、包括的であることも、開示された形態の発明に限定されることも目的としていない。当業者には、多数の変更および変形が明らかであろう。実施形態は、本発明の原理、実際の応用を最もよく説明して、当業者が、意図される特定の用途に適する様々な変更を用いた様々な実施形態に関して、本発明を理解できるように選ばれ、記載された。

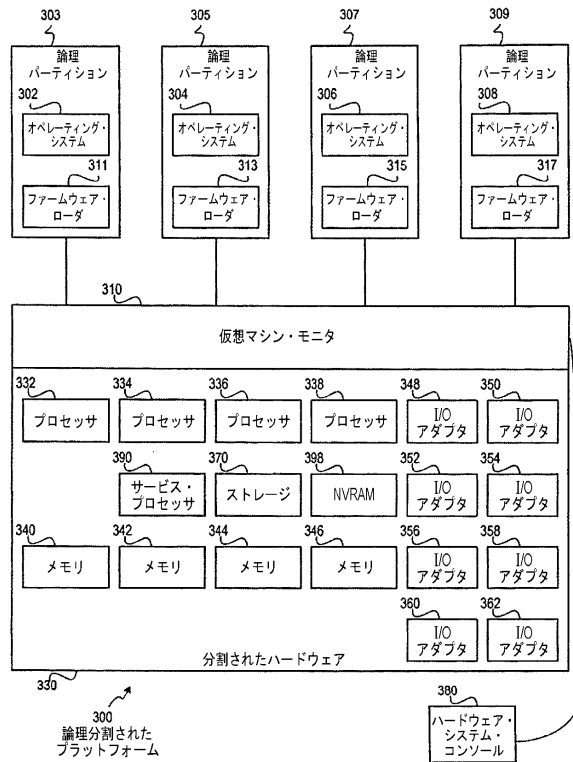
【 図 １ 】



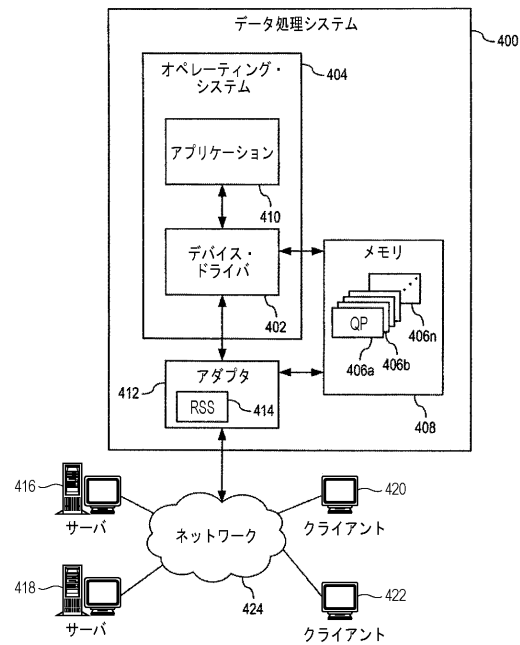
【 図 ２ 】



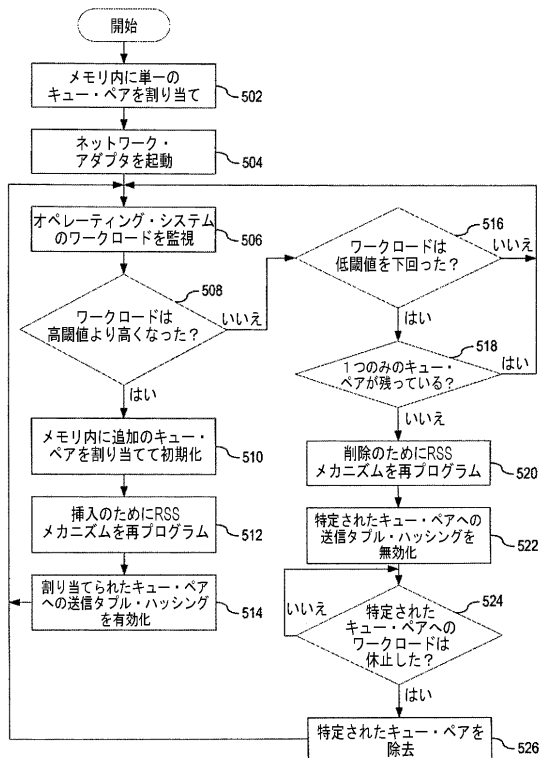
【図 3】



【図 4】



【図 5】



フロントページの続き

- (72)発明者 カルドナ、オマール
アメリカ合衆国 7 8 7 5 8 テキサス州 オースティン バーネット ロード 1 1 5 0 1 9 0
5 - 5 ビー 0 1 7
- (72)発明者 オークス、マシュー、ライアン
アメリカ合衆国 7 8 7 5 8 テキサス州 オースティン バーネット ロード 1 1 5 0 1 9 0
5 - 8 ジー 0 1 9
- (72)発明者 カニンガム、ジェームス、ブライアン
アメリカ合衆国 7 8 7 5 8 テキサス州 オースティン バーネット ロード 1 1 5 0 1 エム
ディー 9 5 5 1
- (72)発明者 シャルマ、ラケシュ
アメリカ合衆国 7 8 7 5 8 テキサス州 オースティン バーネット ロード 1 1 5 0 1 エム
ディー 9 5 5 1

審査官 森谷 哲朗

- (56)参考文献 米国特許出願公開第 2 0 1 1 / 0 1 4 2 0 6 4 (U S , A 1)
米国特許出願公開第 2 0 1 1 / 0 1 5 3 9 3 5 (U S , A 1)
国際公開第 2 0 0 9 / 0 2 7 3 0 0 (W O , A 2)
特開 2 0 0 2 - 2 0 2 9 5 9 (J P , A)
特開 2 0 0 8 - 0 6 0 7 0 0 (J P , A)
特開 2 0 0 7 - 2 5 7 0 9 7 (J P , A)