

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4410795号
(P4410795)

(45) 発行日 平成22年2月3日(2010.2.3)

(24) 登録日 平成21年11月20日(2009.11.20)

(51) Int.Cl. F I
G 0 6 F 9/52 (2006.01) G 0 6 F 9/46 4 7 2 Z

請求項の数 13 (全 18 頁)

(21) 出願番号	特願2006-506127 (P2006-506127)	(73) 特許権者	390009531
(86) (22) 出願日	平成16年4月15日 (2004. 4. 15)		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2006-524381 (P2006-524381A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43) 公表日	平成18年10月26日 (2006.10.26)		アメリカ合衆国10504 ニューヨーク州 アーモンク ニュー オーチャードロード
(86) 国際出願番号	PCT/GB2004/001624		
(87) 国際公開番号	W02004/094863	(74) 代理人	100108501
(87) 国際公開日	平成16年11月4日 (2004. 11. 4)		弁理士 上野 剛史
審査請求日	平成19年4月13日 (2007. 4. 13)	(74) 代理人	100112690
(31) 優先権主張番号	10/422, 193		弁理士 太佐 種一
(32) 優先日	平成15年4月24日 (2003. 4. 24)	(74) 代理人	100091568
(33) 優先権主張国	米国 (US)		弁理士 市位 嘉宏
早期審査対象出願			

最終頁に続く

(54) 【発明の名称】 共有リソースの同時アクセス

(57) 【特許請求の範囲】

【請求項 1】

コンピュータの共有リソースにアクセスする方法であって、前記共有リソースを使用する複数の要求を非同期的に発することが可能な複数のプロセスのうちの第1のプロセスにおいて、

(a) 前記共有リソースを使用する現在の要求に固有の識別子をアトミックに割当てするステップであって、前記固有の識別子を、前記複数のプロセスによって発せられた他の要求に対し前記現在の要求が受け取られる順番に基づいて前記現在の要求に割当てするステップと、

(b) 前記固有の識別子を用いる前記現在の要求を準備することにより当該現在の要求の処理を完了させるステップと、

(c) 前記現在の要求の処理が完了した後に、完了した要求のカウントをアトミックに増分させるステップと、

(d) 前記現在の要求の処理が完了した後に、最近完了した要求より前に発せられた他の要求の処理が完了したかどうかをアトミックに判別するステップと、

(e) 前記最近完了した要求より前に発せられた他の要求の処理が完了したという判別に応答して、処理が完了したディスパッチされていない要求の各々をディスパッチするステップと、

を含む方法。

【請求項 2】

10

20

前記ステップ (b) が、非アトミックに行われる、請求項 1 に記載の方法。

【請求項 3】

前記第 1 のプロセスにおいて、前記現在の要求の処理が完了した後に、前記現在の要求が前記最近完了した要求であるかどうかをアトミックに判別するステップをさらに含む、請求項 1 に記載の方法。

【請求項 4】

前記第 1 のプロセスにおいて、前記現在の要求の処理が完了した後に、前記現在の要求が最後のディスパッチ動作の後に発せられた最初の要求であるかどうかをアトミックに判別するステップと、

前記現在の要求が前記最後のディスパッチ動作の後に発せられた最初の要求であるという判別に応答して前記現在の要求をディスパッチするステップをさらに含む、請求項 1 に記載の方法。

10

【請求項 5】

前記ステップ (e) が、複数の要求をバッチでディスパッチするステップを含む、請求項 1 に記載の方法。

【請求項 6】

前記ステップ (a) が、要求カウンタを増分させるステップを含み、

前記ステップ (c) が、完了カウンタを増分させるステップを含み、

前記方法が、

前記現在の要求が最近完了した要求であるという判別に応答して、前記現在の要求の処理が完了した後に、前記現在の要求の前記固有の識別子と同値になるように最後の完了カウンタをアトミックに更新するステップをさらに含む、

20

前記ステップ (d) が、前記最後の完了カウンタを前記完了カウンタと比較するステップを含み、

前記ステップ (e) が、前記最後の完了カウンタと同値になるように最後のディスパッチ・カウンタを更新するステップを含む、請求項 1 に記載の方法。

【請求項 7】

前記要求カウンタ、前記完了カウンタ、前記最後の完了カウンタ、及び前記最後のディスパッチ・カウンタの各々が、単調に増加する一連の番号を用いて更新される、請求項 6 に記載の方法。

30

【請求項 8】

前記最近完了した要求より前に発せられた他の要求の処理がまだ完了されていないという前記ステップ (d) の否定的な判別に応答して、処理が完了したディスパッチされていない要求の各々のディスパッチを据え置くステップをさらに含む、

前記複数のプロセスのうちの第 2 のプロセスにおいて、

前記現在の要求より前に割当てられた固有の識別子を有する、前記現在の要求に対し先行する要求の処理を完了させるステップと、

前記先行する要求の処理が完了した後に、前記最近完了した要求より前に発せられた他の要求の処理が完了されたことをアトミックに判別するステップと、

前記最近完了した要求より前に発せられた他の要求の処理が完了されたという判別に応答して、処理が完了したディスパッチされていない要求の各々をディスパッチするステップと、

40

を含む、請求項 1 に記載の方法。

【請求項 9】

前記共有リソースが、入力/出力アダプタのフレーム・バッファの少なくとも 1 つのロットを含み、

前記ステップ (b) が、前記ロットを用いてデータを通信するステップを含み、

前記ステップ (e) が、前記通信されるデータの受信をクライアントに通知するステップを含む、請求項 1 に記載の方法。

【請求項 10】

50

前記コンピュータが、論理的に区画化されたコンピュータを含み、

前記入力/出力アダプタが、前記論理的に区画化されたコンピュータの区画間通信のために用いられる仮想ローカル・エリア・ネットワークに結合された仮想入力/出力アダプタである、請求項 9 に記載の方法。

【請求項 1 1】

前記共有リソースが、プロセッサ、ハードウェア・デバイス、入力/出力デバイス、ストレージ・デバイス、アダプタ、コントローラ、データベース、ファイルシステム、サービス、及びこれらのコンポーネントからなる群から選択される、請求項 1 に記載の方法。

【請求項 1 2】

共有リソースと、

前記共有リソースを使用する複数の要求を非同期的に発することが可能な複数のプロセスを実行するように構成された少なくとも 1 つのプロセッサと、

第 1 のプロセスにおいて、前記共有リソースを使用する現在の要求に固有の識別子をアトミックに割り当てるステップと、前記固有の識別子を用いる前記現在の要求を準備することにより当該現在の要求の処理を完了させるステップと、前記現在の要求の処理が完了した後に、完了した要求のカウントをアトミックに増分させるステップと、前記現在の要求の処理が完了した後に、最近完了した要求より前に発せられた他の要求の処理が完了したかどうかをアトミックに判別するステップと、前記最近完了した要求より前に発せられた他の要求の処理が完了されたという判別に応答して、処理が完了したディスパッチされていない要求の各々をディスパッチするステップとを行うために、前記少なくとも 1 つのプロセッサ上で実行されるように構成されたプログラム・コードと、
を含み、前記固有の識別子は、前記複数のプロセスによって発せられた他の要求に対し前記現在の要求が受け取られる順番に基づいて当該現在の要求に割り当てられる、装置。

【請求項 1 3】

請求項 1 ないし請求項 1 1 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ためのコンピュータ・プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に、コンピュータ及びコンピュータ・ソフトウェア、特に、コンピュータの共有リソースへの同時アクセスを管理することに関する。

【背景技術】

【0002】

現代社会においてはコンピュータへの依存度が高まり続けているので、コンピュータ・技術は、増加する需要についていくべく最前線で進歩していかなければならない。重要な研究及び開発努力の 1 つの特定の対象は、並列性、すなわち多数のタスクを並列に処理する性能である。

【0003】

増加した並列処理を容易にするために、多数のコンピュータ・ソフトウェア及びハードウェア技術が開発されている。ハードウェアの観点から、コンピュータは、作業負荷の受容能力の増加を与えるために、ますます多数のマイクロプロセッサに依存しつつある。さらに、多数のスレッドを並列に実行する能力をサポートして、多数のマイクロプロセッサの使用を通して達成できるのと同じ多くの性能ゲインを効果的に与える幾つかのマイクロプロセッサが開発されている。ソフトウェアの観点から、コンピュータ・プログラムが多数のスレッドを同時に実行できるようにし、それにより多数のタスクを本質的に同時に実行することができるマルチスレッド化オペレーティング・システム及びカーネルが開発されている。

【0004】

さらに、幾つかのコンピュータは、論理的区画化の概念を実装し、そこでは、単一の物理的コンピュータが、本質的に多数の独立した「仮想」コンピュータ（論理区画と呼ばれ

10

20

30

40

50

る)のように動作可能にされ、物理的コンピュータ(例えば、プロセッサ、メモリ、入力/出力デバイス)における種々のリソースは、種々の論理区画の間で割当てられる。各論理区画は、個別のオペレーティング・システムを実行し、ユーザ及び論理区画上で実行するソフトウェア・アプリケーションの見地から、完全に独立したコンピュータのように動作する。

【0005】

並列性は、多数のタスクを一度に行う能力によりシステム性能を効果的に高めるが、並列性の1つの副作用は、多数の同時処理動作を同期させる必要があるために、特に、多数のスレッドによって共有されることが可能なシステム・リソースに関するシステムの複雑さが増すことである。特定のリソースにアクセスすることができる個別のスレッド又はプロセスは、典型的には、他のスレッド又はプロセスのアクティビティを認識しない。したがって、1つのプロセスが別のプロセスに関する特定のリソースに予期しない形でアクセスして、不確定の結果及び潜在的なシステムエラーをもたらす恐れがある。

10

【0006】

しかし、ほんの一例として、イーサネット(登録商標)・ネットワーク・アダプタのような入力/出力(I/O)アダプタは、コンピュータにおいて実行される多数のスレッド又はプロセスによって共有されることが可能な有限の数のリソースを有することができる。こうしたタイプの1つのリソースはフレーム・バッファであり、それは、イーサネット(登録商標)と互換性のあるネットワークに及び該ネットワークから通信されるデータ・フレームを格納するのに用いられる。フレーム・バッファは、異なるプロセスによって順番に用いられるスロットに配置されることが多く、I/Oアダプタは、各スロットの消費を通知される。しかしながら、スロットの1つの消費が多大な時間量を占領し、その間、他のプロセスが他のスロットを利用できないことがある。さらに、2つのプロセスが同じスロットを同時に使用しようとする懸念がある。

20

【0007】

これらの懸念に対処するために、ロック又はセマフォといった逐次化機構を用いて、共有リソースへのアクセスを一度に一つのプロセスに制限することができる。ロック又はセマフォは、本質的には、特定の共有エンティティ又はリソースにアクセスするためにマルチスレッド化環境におけるプロセス又はスレッドによって排他的に得られる「トークン」である。プロセス又はスレッドがリソースにアクセスする前に、最初にシステムからトークンを取得しなければならない。別のプロセス又はスレッドがトークンを同時に所有する場合には、他のプロセス又はスレッドによってトークンが解放されるまで、前者のプロセス又はスレッドは、リソースにアクセスすることを許可されない。このようにして、リソースへのアクセスは、不確定の動作が起こるのを防止するために効果的に「逐次化」される。

30

【0008】

ロックされたリソースにアクセスする必要があるあらゆるスレッド又はプロセスは、リソースへのアクセスを進める前に、そのリソースのロックが解除されるまで待たなければならないことが多いので、共有リソースが逐次化機構によって管理されるときにはいつも、並列化の機会が減少されることが多い。したがって、共有リソースにロックがかけられる頻度及び持続時間を最小にすることへの強い要望がある。

40

【0009】

前述のI/Oアダプタの例においては、フレーム・バッファへの対立のないアクセスを保証するための最も直接的な解決策は、フレーム・バッファのあらゆる使用に対して単にロックをかけることである。しかしながら、そのようにすると、フレーム・バッファに格納されたデータを準備し、解釈し、その他の処理をするプロセスによって大量の作業が実行される必要があるため、並列化の機会が大幅に制限される。したがって、これらの全てのアクティビティが発生する間の他のプロセスによるフレーム・バッファへのアクセスのロックアウトは、並列化の観点からは逆効果である。

【0010】

50

或る状況においては、共有リソースにロックをかける前に、共有リソースの使用に関連する幾つかの事前作業を行って、共有リソースが他のプロセスに対してロックされる必要がある時間の量を減らすことが可能である。しかしながら、そのようにすることを意図された従来の試みは、リソースの利用可能性の重複チェックと、より複雑な事前作業が実行されることを要求することが見出されている。

【発明の開示】

【発明が解決しようとする課題】

【0011】

したがって、マルチスレッド化されたコンピュータにおいて実行される多数のスレッド又はプロセスによる共有リソースへのアクセスを調整する改善された方法への大きな必要性がある。

10

【課題を解決するための手段】

【0012】

本発明は、従来技術に関連するこれらの問題及び他の問題に対処するために、マルチスレッド化されたコンピュータにおける複数のプロセスによる共有リソースへのアクセスが、共有リソースを使用する要求を受け取る順番と、こうした要求を受け取った後にそれらの要求の処理が完了する順番との両方を追跡するアトミック・オペレーションの集合を介して管理されるといような装置、プログラム及び方法を提供する。要求のディスパッチは、最近完了された要求より前に受け取った、ディスパッチされていない要求の処理が完了するまで、効果的に据え置かれる。

20

【0013】

特に、本発明の1つの態様によれば、共有リソースは、共有リソースを使用する要求を非同期的に発することが可能な複数のプロセスによってアクセスされることができる。第1のプロセスにおいて、共有リソースを使用する現在の要求に固有の識別子がアトミックに割当てられる。この場合、固有の識別子は、複数のプロセスによって発せられた他の要求に対し現在の要求を受け取る順番に基づいて現在の要求に割当てられる。その後、固有の識別子を用いる現在の要求の処理が完了し、そうした処理が完了した後に、完了した要求のカウントがアトミックに増分され、最近完了した要求より前に発せられた他の要求の処理が完了したかどうかについてのアトミックな判別がなされる。完了していた場合、処理が完了した、ディスパッチされていない要求の各々がディスパッチされる。

30

【0014】

本発明に係る幾つかの実施形態においては、要求の処理の完了は、非アトミックに行うことができ、したがって、要求の処理と関連して行われるアトミック・オペレーションの持続時間を制限し、それにより共有リソースへの競合を減少させる。さらに、幾つかの実施形態においては、ディスパッチされていない要求のディスパッチの結果として、多数の要求がバッチで対処され、それにより個々のディスパッチのオーバーヘッドが減少する。

【0015】

ここで、本発明を、添付の図面を参照しながら、一例として説明する。

【発明を実施するための最良の形態】

【0016】

40

以下に示される実施形態は、複数のクライアントによって発せられた、共有リソースを使用する要求を受け取る順番と、共有リソースを使用する要求のディスパッチと関連してこうした要求の処理が完了する順番との両方を追跡するものである。要求は、非同期的に発せられることができ、さらに、順番どおりではなく完了させることができる。結果として、要求を受け取る際に、要求を受け取る順番に基づいて、要求に固有の識別子が割当てられる。さらに、完了した要求の番号が追跡され、最近完了した要求より前に発せられた他の要求の処理が完了したかどうかを判別するために、各要求の取り扱いと関連して判別がなされる。この条件が満たされたと判断されると、処理が完了した、ディスパッチされていない要求の各々がディスパッチされる。

【0017】

50

以下で一層明らかとなるように、要求を受け取る順番、及びこうした要求が完了する順番の追跡は、共有リソースに関連する一連のカウンタの使用を通じて行われ、これらのカウンタは、セマフォのような排他的ロック機構を介してアトミックにアクセスされることが可能である。要求カウンタは、要求を受け取る際に、要求に固有の識別子を割当てるために用いることができる。完了した要求（すなわち、準備及びサブミッションといった処理が完了した要求）の番号を追跡するために、完了カウンタを用いることができる。最後の完了カウンタは、最近完了した要求を追跡するのに用いることができ、一方、最後のディスパッチ・カウンタは、最近ディスパッチされた要求を追跡するのに用いることができる。

【 0 0 1 8 】

本発明に係る多くの実施形態においては、上述のカウンタへのアクセスのみがアトミックである必要がある。要求の取り扱いに関係する他の処理は、非アトミックに取り扱われることが多く、したがって、要求の取り扱いに関連して行われるアトミック・オペレーションの持続時間が制限され、それにより、そうでなければ共有リソースへの同時アクセスを妨げる競合が減少される。

【 0 0 1 9 】

さらに、本発明に係る実施形態は、要求のバッチでのディスパッチを与える能力をサポートする。要求のディスパッチは、本発明に係る異なる形で行われる。例えば、共有リソースが、ネットワーク・アダプタのような入力/出力アダプタのフレーム・バッファにおけるスロットである場合には、ディスパッチオペレーションは、入力/出力アダプタを介してそのクライアントに関するパケットが受け取られるときにはいつでも、クライアントへの割り込みの生成として実装することができる。共有リソースが、SCSIネットワークのようなネットワークにおける大容量ストレージ・デバイスである場合には、ディスパッチ・オペレーションは、オペレーティング・システムへのI/O要求の完了通知とすることができる。

【 0 0 2 0 】

要求の性質は、典型的には、当該共有リソースの特定のタイプに応じて変わることになる。例えば、データベース・リソースにおいては、要求はデータベース・トランザクションを表すことができ、一方、ファイル管理システム・リソースにおいては、要求はファイルアクセス要求を表すことができる。ストレージ・デバイスにおいては、要求は入力/出力(I/O)要求を表すことができ、一方、プロセッサ・リソースにおいては、要求はプロセッサに論理区画を割当てるための要求を表すことができる。本発明の開示内容の利点を有する他の共有リソース及び要求の実装が、当業者には明らかであろう。

【 0 0 2 1 】

ここで、幾つかの視点を通じて同じ番号が同じ部分を表す図面に移ると、図1は、本発明に係る共有リソースへの同時及び非同期アクセスをサポートするために用いることができる汎用コンピュータ・アーキテクチャ50の例示的な構成を示す。特に、アーキテクチャ50は、汎用共有リソース56への複数のクライアント54によるアクセスを管理するリソース・マネージャ52を組み込んでいることが示されている。

【 0 0 2 2 】

共有リソース56は、コンピュータ、例えば、種々のタイプの入力/出力、イーサネット(登録商標)・アダプタ、SCSIアダプタその他のネットワーク・アダプタ、ストレージ・アダプタ及びコントローラ、ワークステーション・アダプタ及びコントローラ、データベース・コントローラ、プロセッサといったその他のハードウェア・リソース、その他の入力/出力アダプタ及びコントローラなどによってアクセスされることができる事実上あらゆるタイプのコンピュータ・リソースを表すと考えられる。さらに、共有リソース56は、例えば、カーネル又はオペレーティング・システム・サービス、ファイル管理システム、データベースなどといった種々のタイプのソフトウェア・リソースを含むと考えられる。共有リソースはまた、例えば、バッファ又は入力/出力アダプタに常駐する他のメモリ・リソースといった上記のタイプのリソースのサブコンポーネントを含むことがで

10

20

30

40

50

きる。

【0023】

クライアント54は、それぞれ、事実上あらゆるプログラム・コード、又はその他のクライアントに対して非同期的な形でリソース56にアクセスすることができる外部プログラマブル・デバイスも表すと考えられる。図示された実施形態においては、典型的には、各クライアント54は、例えば、シングルスレッド化及び/又はマルチスレッド化プロセッサを含む1つ又はそれ以上のプロセッサ、ソフトウェアベースの又はオペレーティング・システム支援マルチスレッド化、論理的区画化、及び/又は多数の実行インスタンスが共有リソースに並列に非同期的にアクセスできるようにする他の技術を使用して、コンピュータに常駐する複数の並列処理又はスレッドの1つにおいて実行される、コンピュータ上のプログラム・コードを用いて実装される。

10

【0024】

共有リソース56へのアクセスを管理するために、リソース・マネージャ52は、種々のクライアント54によって生成されたリソース・アクセス要求のための個別のサーバとして働く複数のサーバ・プロセス58を含む。各サーバ・プロセス58は、所与のクライアント54に静的に又は動的に割当てられ、サーバ・プロセスとクライアントとの間のマッピングは、必ずしも1:1である必要はなく、時間の経過と共に変化してもよいことがわかるであろう。

【0025】

リソース・マネージャ52はまた、共有リソース56へのアクセスを調整するためにサーバ・プロセス58と対話するディスパッチ・プロセス60を含む。以下でより明らかとなるように、ディスパッチ・プロセス60は、一度に多数のアクセス要求をディスパッチするためにバッチでのディスパッチを実行することができる。さらに、ディスパッチ・プロセス60はまた、共有リソース56からのデータのリターンを管理ことができ、入来データをしかるべきクライアントに送付するために、サーバ・プロセス58と、あるいは直接各クライアントと対話することができる。しかしながら、他の実施形態においては、リターン・データは、ディスパッチ・プロセス60とは別に取り扱われてもよい。さらに、リソース・マネージャ52は、多数の共有リソースを管理ことができ、又は多数の共有リソースを管理するために多数のリソース・マネージャを用いることができる。

20

【0026】

本発明に係る共有リソース56への同時アクセスをサポートするために、リソース・マネージャ52は、4つのアトミック・カウンタ62、64、66及び68を含む。カウンタ62は、ここでは要求カウンタと呼ばれ、クライアント54の1つによって生成された次のアクセス要求に割当てられるべき識別子を与えるために用いられる。生成された各要求のために、要求カウンタ62が次の識別子に増分され、識別子は、要求を固有に識別するために、サーバ・プロセスにおけるMY__REQUEST変数70に格納される。種々の実施形態においては、要求カウンタ62は、カウンタの現在値が、最後に受け取られた要求を表すか、又は要求が受け取られたときに次の要求に割当てられることになる値を表すように、事前に又は事後に増分される。

30

【0027】

カウンタ64は、ここでは完了カウンタと呼ばれ、完了したアクセス要求の数を追うために用いられる。カウンタ66は、ここでは最後の完了(「LAST__DONE」)カウンタと呼ばれ、完了した最後の(すなわち最新の)アクセス要求の識別子を維持するために用いられる。典型的には、要求が順番どおりではなく完了することが許される、すなわち、後続する要求が先行する要求より先に完了することが許されるとすれば、カウンタ64及び66は両方とも必要とされる。

40

【0028】

カウンタ68は、最後のディスパッチ・カウンタと呼ばれ、ディスパッチ・プロセス60によってディスパッチされた共有リソースへの最後の(又は最新の)アクセス要求の識別子を格納するために用いられ、典型的にはディスパッチ・プロセスが最後に通知された

50

最後の完了カウンタ66の値に設定される。カウンタ62、64、66及び68を識別するために用いた特定の名称は、単に便宜的に選定されたのであって、本発明の範囲を制限するものではない。

【0029】

図示された実施形態においては、各カウンタ62、64、66及び68がゼロに初期化され、カウンタ62及び64が1だけ増分され、その結果、各カウンタは、常に増加する整数で常に更新される。しかしながら、本発明に係る要求を固有に識別するために、他の単調な順番番号又は他の識別子を用いることもできることが分かるであろう。さらに、特定の数の要求が処理された後にカウンタがロール・オーバーすることを可能にするカウンタ・サイズを用いることができるが、実行可能な時間内でカウンタがロール・オーバーすることを可能にしないカウンタ・サイズ（例えば、8バイトカウンタ）を用いることによって、処理が実質的に簡単化されることが見出された。また、システムの始動の間に各カウンタを再び初期化して、実際問題として、コンピュータ・アーキテクチャ50の動作の通常の過程ではカウンタのロール・オーバーが起こらないようにすることが望ましい。

10

【0030】

ここで図2に移ると、本発明に係るクライアント要求を処理するのに用いられる例示的なクライアント要求処理ルーチン72が、かなり詳しく示されている。ルーチン72は、各要求に応答してサーバ・プロセス58によって実行され、したがって、多数のクライアントが同時に要求をサブミットするときには、ルーチン72の多数のインスタンスがあらゆる所与の時点で並列に実行される。

20

【0031】

ルーチン72は、ブロック74において、カウンタ・ロックを得ることによって、すなわち、カウンタ62 - 68に関連するロック又はセマフォーを得ることによって始まる。図示された実施形態においては、こうしたセマフォーは、全てのカウンタ62 - 68に割当てることができ、又は代替的に、こうしたカウンタの各々に別のロックを与えることができ、或いはカウンタのサブセットに所与のロックを割当てることができる。

【0032】

ロックが多数のカウンタによって共有されるか、又はカウンタの1つ又はサブセットにのみ割当てられるかどうかに関係なく、ブロック74においてロックが得られると、制御がブロック76に移って要求カウンタ62を増分する。次に、ブロック78が、要求に関連付けられることになるMY__REQUEST変数70を要求カウンタ62の現在値に設定する。その後、ブロック80がカウンタ・ロックを解除する。その結果、ブロック74 - 80が、要求カウンタにおけるアトミック増分動作を有効に実行し、また、要求カウンタの増分後の値に基づいて要求に固有の要求識別子を割当てて。他の実施形態においては、要求カウンタの増分前の値を用いることができる。さらに、ロックの下で要求カウンタの有効なコピーが得られる限り、ブロック78のMY__REQUEST変数への値の割当ては、ロックの外部で行われても良い。しかしながら、用いられる方法に関係なく、要求は、複数のクライアントによって発せられた他の要求に対して要求が受け取られる順番に基づいて、固有の識別子を有効に割当てられる。

30

【0033】

次に、ブロック82が、クライアントに代わって要求を準備し、サブミットすることによって要求の処理を完了する。しかしながら、注目すべきは、ブロック82は、ロック又はセマフォーの下では実行されず、したがって、共有リソースにアクセスする際の対立の結果として現在のプロセス又は別のプロセスを停止させることがある、あらゆる競合の問題を伴わない。要求を準備し、サブミットすることに関連するステップは、典型的には、特定の要求のタイプ及びアクセスされるリソースのタイプに応じて変化することになる。例えば、ブロック82において、入力/出力アダプタへの要求をサブミットする目的のために、直接メモリ・アクセス(DMA)動作を介した伝送のためにネットワーク・フレームを準備するステップといった種々のステップを行うことができる。本発明に係るブロック82においては、要求を準備し、要求をサブミットし、及び/又は、要求の処理を、例

40

50

えば、メモリ割当て、データのコピー又は伝送、プールされたリソースの取得などといった他の方法で完了することに関連する他の動作を実行することもできる。

【0034】

次に、ブロック84はまた、カウンタ・ロックが得られるのに必要なだけ待って、カウンタ・ロックを得る。カウンタ・ロックが得られると、制御がブロック86に移って完了カウンタ64を増分し、それにより現在の要求の処理の完了を確立する。

【0035】

次に、ブロック88で、現在の要求の識別子が最後の完了カウンタ66に格納された値より大きいかどうかを判別する。大きい場合には、現在の要求は、完了した最新の又は最後の要求であることが分かる。したがって、制御がブロック90に移って、最後の完了カウンタ66の値が現在の要求の識別子と同値に設定される。次に、制御がブロック92に移る。さらに、ブロック88に戻り、現在の要求の識別子が最後の完了カウンタ66に格納された値より大きくない場合には、ブロック88は、制御を直接ブロック92に移し、それによりブロック90をバイパスする。

【0036】

ブロック92は、典型的には、完了カウンタ64に格納された値が最後の完了カウンタ66に格納された値に等しいかどうかを判別することによって、最近完了した要求より前に発せられた他の要求の処理が完了したかどうかを判別する。完了カウンタ64の値がより小さいときには、前の要求の少なくとも1つがまだ完了していないことを示すので、この条件が満たされているときには、最近完了した要求(カウンタ66において識別された)より時間が先行する全ての要求が完了したことがわかる。

【0037】

ブロック92における条件が満たされた場合には、制御がブロック94に移って、最後のディスパッチ・カウンタ68の値が最後の完了カウンタ66に格納された値と同値に設定され、それは、最後の完了カウンタにおいて識別された要求とそれまでの要求を含む全ての要求が完了したことを示す。次に、制御がブロック96に移って、全ての今のところまだディスパッチされていない要求(典型的には、最後の完了カウンタ66の値と、ブロック94を実行する前の最後のディスパッチ・カウンタ68に格納された前の値との間の要求)をディスパッチするようにディスパッチ・プロセスに通知する。次に、制御がブロック98に移って、ブロック84において得られたカウンタ・ロックを解除し、それによりルーチン72が完了する。本発明に係る他の実施形態においては、ブロック96の後にブロック94が実行されることを認識されたい。

【0038】

ブロック92に戻り、全ての先行の要求が完了していない場合には、ブロック92が制御をブロック100に移し、典型的には現在の要求の識別子が最後のディスパッチ・カウンタ68に格納された値より大きいかどうかを判別することによって、現在の要求が、ディスパッチ・プロセスによる要求の最後のディスパッチの後に始められた最初の要求であったかどうかを判別する。

【0039】

そうである場合には、制御がブロック102に移って、カウンタ68の値が現在の要求識別子と同値に設定される。次に、制御がブロック96に移って、ディスパッチ・プロセスに現在の要求をディスパッチするように通知する。そうでなければ、ブロック100に戻り、ブロックにおいて示された条件が満たされない場合には、後で実行される他の作業が要求を最終的にはディスパッチさせるので、このルーチン72のインスタンスを通して実行されるべき更なる作業は要求されない。したがって、ブロック100は、制御を直接ブロック98に移してカウンタ・ロックを解除し、ルーチンを終える。

【0040】

ブロック100において行われる決定、及び結果としてのブロック102の流れは、幾つかの実施形態においては省略されても良い。上記の機能は、バッチでディスパッチした後に始められた最初の要求の待ち時間を減らすという観点での利点を与える。幾つかの環

10

20

30

40

50

境においては、最初の要求の待ち時間は大きな懸念とはならず、そのため上記の機能は、こうした実施形態においては省略されても良い。こうした事例においては、ブロック 9 2 は、そこで示された条件が満たされていない場合には、制御をブロック 9 8 に直接移すことができる。

【 0 0 4 1 】

上記の要求を取り扱う技術は、多数のクライアントによる共有リソースのアクセスを管理する観点で多くの固有の利点を与える。多数のアクセス要求が、要求の生成の観点から同期化又は逐次化を必要とすることなく生成され、保持されることができる。例えば、要求が、制限された数のブロック又はスロットに区画化されたバッファに組み込まれる場合には、クライアントは、こうしたスロットの完全な消費を同期させる必要なしに特定のスロットに同時に割当てられる。さらに、1つの共有リソースに対する多数の要求を、1つずつではなくバッチでディスパッチすることができ、それにより各ディスパッチに関連するオーバーヘッド（例えば、割り込みを介してクライアントに通知することに関連するオーバーヘッド）の多くが排除される。さらに、共有リソースへのアクセスを保護し、管理するのに用いられるロックにおける競合が減少される。ロックは、前述のカウンタの比較的短い更新の間の特定の要求のためにのみかけられ、何らかの同期化機構が得られることを必要とすることなく、そのディスパッチが行われる前の要求の作成及び準備に関連するオーバーヘッドの多くをイネーブルにする。

【 0 0 4 2 】

上述のように、共有リソースへの同時アクセスを管理するための上記の技術は、多様な用途に用いることができる。図 3 - 図 5 と関連して後で詳しく説明される 1 つの特定の用途は、図 3 のコンピュータ 1 0 のような論理的に区画化されたコンピュータにおける物理的又は仮想イーサネット（登録商標）・アダプタへのアクセスを管理することである。特に、イーサネット（登録商標）・アダプタ、トークン・リング・アダプタなどのような入力/出力アダプタは、この事例においてはフレーム・バッファのスロットである有限の数のリソースのみを有することが良く知られている。典型的には、多数のクライアントがスロットを使用することを許される。さらに、スロットの消費は、長時間にわたって該スロットを拘束して、その間、他のクライアントが該スロットを使用できなくなることがある。多数のクライアントが、他のクライアントを妨害することなくこうしたスロットを同時に消費できるようにすることにより、コンピュータ・システムにおけるより有効な並列実行が可能になり、それによりシステム性能が向上する。後述するように、コンピュータ 1 0 は、本発明に係るこうしたスロットの同時消費をサポートする。

【 0 0 4 3 】

コンピュータ 1 0 は、通常は、例えば、ネットワーク・サーバ、中間価格帯のコンピュータ、メインフレーム・コンピュータなど、例えば、IBM eServer（eServer は IBM の商標）コンピュータといったあらゆる数のマルチユーザ・コンピュータを表す。しかしながら、本発明は、他のコンピュータ及びデータ処理システム、例えば、ワークステーション、デスクトップ・コンピュータ、ポータブル・コンピュータなどといった単一ユーザ・コンピュータ、又は他のプログラム可能電子デバイス（例えば、内蔵コントローラなどを組み入れる）において実装されてもよいことを認識されたい。さらに、本発明はまた、論理的に区画化されていないマルチスレッド化コンピュータと組み合わせ

【 0 0 4 4 】

図 3 に最も良く示されるように、コンピュータ 1 0 は、一般に、バス 1 6 を介してメモリ 1 4 に結合された 1 つ又はそれ以上のプロセッサ 1 2 を含む。各プロセッサ 1 2 は、シングルスレッド化プロセッサとして、又は、複数のハードウェア・スレッド 1 8 を組み込んでいることが示されたプロセッサ 1 2 a のようなマルチスレッド化プロセッサとして実装することができる。大部分は、マルチスレッド化プロセッサ 1 2 a における各ハードウェア・スレッド 1 8 は、コンピュータに常駐するソフトウェアによって、独立プロセッサと同様に処理される。これに関しては、本開示の目的上、シングルスレッド化プロセッサ

10

20

30

40

50

は、シングル・ハードウェア・スレッド、すなわち単一の独立型実行ユニットを組み入れると考えられる。しかしながら、コンピュータにおける多数のタスクの並列性能をさらに支援するために、ソフトウェアベースのマルチスレッド化又はマルチタスク化を、シングルスレッド化プロセッサ及びマルチスレッド化プロセッサの両方と組み合わせて用いることができる。

【0045】

さらにまた、図3に示されるように、1つ又はそれ以上のプロセッサ12（例えば、プロセッサ12b）は、システムの初期プログラム・ロード（IPL）を管理するために、及びシステム・ハードウェアを監視し、診断し、構成するために特殊ファームウェア・コードを実行するのに用いられるサービス・プロセッサとして実装することができる。一般に、コンピュータ10は、オペレーティング・システム及びコンピュータに常駐するアプリケーションを実行するのに用いられる1つのサービス・プロセッサと多数のシステム・プロセッサを含むが、本発明はこの特定の実装に限定されない。幾つかの実装においては、サービス・プロセッサは、バス16を通す以外の形で、コンピュータにおける種々の他のハードウェア・コンポーネントに結合することができる。

10

【0046】

メモリ14は、1つ又はそれ以上のレベルのデータ、命令及び/又は組み合わせキャッシュだけではなく、例えば、DRAMベースの主ストレージといった1つ又はそれ以上のレベルのメモリ・デバイスを含むことができ、特定のキャッシュは、当該技術分野では周知のように個々のプロセッサ又は多数のプロセッサと関連して動作する。さらに、メモリ14は、バス20、例えば1つ又はそれ以上のネットワーク・アダプタ22（コンピュータとネットワーク24をインターフェースするための）、1つ又はそれ以上のストレージ・コントローラ26（コンピュータと1つ又はそれ以上のストレージ・デバイス28をインターフェースするための）、及び1つ又はそれ以上のワークステーション・コントローラ30（複数のワークステーション・アダプタを介して1つ又はそれ以上の端末又はワークステーション32とインターフェースするための）を介して多くのタイプの外部デバイスに結合される。

20

【0047】

図3はまた、区画マネージャ又はハイパーバイザー36によって管理される複数の論理区画34を含む、コンピュータ10上の論理的に区画化されたコンピューティング環境を実装するのに用いられる主ソフトウェア・コンポーネント及びリソースをかなり詳細に示す。当該技術分野では公知のあらゆる数の論理区画をサポートすることができ、あらゆる時点でコンピュータに常駐する論理区画の数は、区画がコンピュータに加えられ、又はコンピュータから除去される際に動的に変化することができる。

30

【0048】

図示されたIBM eServerベースの実装においては、区画マネージャ36は、2つのプログラム・コードの層からなる。ここではディスパッチ不可能部分38と呼ばれる第1の層は、コンピュータ10のファームウェア又はライセンス内部コード（LIC）内で実装され、それは、ハードウェア・アクセスの詳細から例えばオペレーティング・システムのようなより高次の層を分離しながら、種々のハードウェア・コンポーネントとの低レベルのインターフェースを与えるために用いられる。ファームウェアはまた、サービス・プロセッサ12bのようなサービス・プロセッサと通信することができる。ディスパッチ不可能部分38は、コンピュータ10のための多くの低レベルの区画管理機能、例えば、ページ・テーブル管理などを与える。ディスパッチ不可能部分38はまた、タスクの概念をもたず、主に、ソフトウェアのより高次の層からの機能コールを介してアクセス可能である。

40

【0049】

区画マネージャ36におけるプログラム・コードの第2の層は、ここではディスパッチ可能部分40と呼ばれる。再配置オフの状態で行われ、ソフトウェアのより高次の層からの機能コールを介してアクセス可能な、タスクの概念をもたないディスパッチ不可能部

50

分38とは対照的に、ディスパッチ可能部分40は、(あらゆるオペレーティング・システムと同様に)タスクの概念を有し、再配置オンの状態で実行される。ディスパッチ可能部分は、典型的には、それがユーザから隠れていること以外は、区画とほぼ同じような形で実行される。ディスパッチ可能部分は、一般に、区画の作成及び削除、同時I/Oメンテナンス、プロセッサ、メモリ及び他のハードウェア・リソースの種々の区画34への割当てといったより高いレベルの区画管理動作を管理する。

【0050】

各論理区画34は、典型的には、コンピュータ10における利用可能なリソースの一部に静的に及び/又は動的に割当てられる。例えば、各論理区画は、1つ又はそれ以上のプロセッサ12、及び/又は、1つ又はそれ以上のハードウェア・スレッド18、並びに、
10
利用可能なメモリスペースの一部に割当てられる。論理区画は、プロセッサのような特定のハードウェア・リソースを共有して、所与の1つのプロセッサが、1つより多い論理区画によって用いられるようにすることができる。別のハードウェアにおいては、リソースは、一度に1つだけの論理区画に割当てられることができる。

【0051】

付加的なリソース、例えば、大容量ストレージ、バックアップ・ストレージ、ユーザ入力、ネットワーク接続、したがってI/Oアダプタは、典型的には、当該技術分野では周知の形で1つ又はそれ以上の論理区画に割当てられる。リソースは、例えば、1つ1つのバスに基づいて、又は1つ1つのリソースに基づいて、といった多くの形で割当てられ、
20
多数の論理区画は、同じバスにおけるリソースを共有する。幾つかのリソースは、一度に多数の論理区画に割当てられることもできる。

【0052】

各論理区画34は、区画化されていないコンピュータのオペレーティング・システムと同じ形で論理区画の主動作を制御するオペレーティング・システム42を用いる。例えば、各オペレーティング・システム42は、インターナショナル・ビジネス・マシーズ社から入手可能なOS/400オペレーティング・システムを用いて実装することができる。

【0053】

各論理区画34は、別々の、すなわち独立したメモリスペースにおいて実行され、それにより各論理区画は、そうした各論理区画において実行される各ユーザ・アプリケーション(ユーザアプリ)44の見地から、独立した区画化されていないコンピュータとほぼ同様に働く。したがって、ユーザ・アプリケーションは、典型的には、区画化環境において用いるためのあらゆる特別な構成を必要としない。

【0054】

別々の仮想コンピュータとしての論理区画34の性質から、論理区画があたかも別の物理マシン上にあるようにして別の論理区画と通信することを可能にする区画間通信をサポートすることが望ましい。したがって、幾つかの実装においては、ディスパッチ不可能部分38において、論理区画34がイーサネット(登録商標)・プロトコルのようなネットワーク・プロトコルを介して別の論理区画と通信することを可能にする仮想ローカル・エリア・ネットワーク(LAN)46をサポートすることが望ましい。区画間の通信を
40
サポートする別の方法も、本発明に従ってサポートすることができる。

【0055】

本発明に従って他の論理的に区画化された環境を用いることもできることが認識されるであろう。例えば、あらゆる区画34から分離されたディスパッチ可能部分40を用いるのではなく、ディスパッチ可能部分の機能を、1つ又はそれ以上の論理区画に代替的に組み入れてもよい。

【0056】

一般に、本発明の実施形態を実装するために実行されたルーチンは、オペレーティング・システム又は特定のアプリケーション、コンポーネント、プログラム、オブジェクト、モジュール又は命令のシーケンス、或いはそのサブセットの一部のいずれとして実装され
50

るかに関係なしに、ここでは「コンピュータ・プログラム・コード」又は単に「プログラム・コード」と呼ばれる。プログラム・コードは典型的に、コンピュータの種々のメモリ及びストレージ・デバイスに種々の時点で常駐する1つ又はそれ以上の命令を含み、それは、コンピュータの1つ又はそれ以上のプロセッサによって読み取られ、実行されたときに、本発明の種々の態様を具体化するステップ又は要素を実行するのに必要なステップをコンピュータに実行させる。さらに、本発明は、満足に機能するコンピュータ及びコンピュータ・システムとの関連で以下に説明されるが、本発明の種々の実施形態は、種々の形態のプログラム製品として配布されることができ、そして、本発明は、実際に配布を行うのに用いられる信号伝達媒体の特定のタイプに関係なく等しく適用されることを当業者であれば認識するであろう。信号伝達媒体の例は、この限りではないが、揮発性及び不揮発性メモリ・デバイス、ディスクその他のリムーバブル・ディスク、ハードディスクドライブ、磁気テープ、光ディスク（例えば、CD-ROM、DVDなど）などのような書き込み可能型媒体、並びに、デジタル及びアナログ通信リンクのような伝送型媒体を含む。

10

【0057】

さらに、以下に説明される種々のプログラム・コードは、本発明の特定の実施形態において実装されるアプリケーション又はソフトウェア・コンポーネントに基づいて識別することができる。しかしながら、つけられたあらゆる特定のプログラム名は、単に便宜のために用いられたのであり、したがって本発明は、そうした名称によって識別され及び/又は示される何らかの特定のアプリケーションにおける使用にのみ限定されるものではないと認識されたい。さらに、コンピュータ・プログラムを、ルーチン、プロシージャ、メソッド、モジュール、オブジェクトなどに構成する無数の典型的な方法、並びに、典型的なコンピュータ内に常駐する種々のソフトウェア層（例えば、オペレーティング・システム、ライブラリ、API、アプリケーション、アプレットなど）の間にプログラム機能を割当てる種々の方法があり、本発明は、ここで説明されたプログラム機能の特定の構成及び割当てに限定されるものではないと認識されたい。

20

【0058】

図3に示された例示的な環境は、本発明を制限することを意図されたものではないことを当業者は理解するであろう。実際は、本発明の範囲から逸脱することなく、他の代替的なハードウェア及び/又はソフトウェア環境を用いることができることを当業者は理解するであろう。具体的には、本発明は、論理的に区画化された環境における使用に限定されるものではなく、共有リソースへの同時及び非同期アクセスをサポートする多くの他の環境において用いることができる。

30

【0059】

ここで図4に移ると、ここでは仮想入力/出力アダプタにおけるフレーム・バッファである特定のタイプの共有リソースへの同時アクセスを実行することに関連する主ソフトウェア・コンポーネントが示されている。この実施においては、複数のクライアント120は、仮想ローカル・エリア・ネットワーク（VLAN）を介して区画34間の区画間通信をサポートする目的のために、種々の区画34におかれる。図示された実施形態においては、各クライアントは、例えば、区画に常駐するデバイス・ドライバとして実装することができる。

40

【0060】

区画マネージャ36内には、複数の仮想入力/出力アダプタ（IOA）122が配置され、それらは、VLAN上での通信のためにクライアント120によってアクセスされることができる。こうした構成の下では、1つの区画のクライアントが別の区画のクライアントにデータを通信する必要があるときにはいつも、データは、非仮想IOAにデータが転送されるのと同様のように、区画マネージャ34において送信側クライアントから送信側クライアントに割当てられた仮想IOA122に転送される。データは、仮想IOAによって従来のネットワーク・パケット（例えば、イーサネット（登録商標）互換パケット）と同様にフォーマットされ、異なる区画34の受信側クライアント120に割当てら

50

れる別の仮想 I O A 1 2 2 に通信される。このデータ通信は、物理的 L A N において用いられるような物理的ケーブル上ではなくソフトウェアにおいて行われるが、その他の点では、使用される通信プロトコルは同じである。次に、受信側クライアントに割当てられた仮想 I O A 1 2 2 が、受信したデータを、非仮想 I O A とほぼ同じようにして受信側クライアントに渡す。したがって、各クライアント 1 2 0 の見地から、データが或るクライアントによって送信され、別のクライアントによって受信される方法は、データが物理的 L A N 上で通信されるのと同じである。

【 0 0 6 1 】

この方法で仮想 L A N を実装することによって、あたかも他の区画は異なる物理的コンピュータ上に常駐するかのようにして、或る区画において実行されるアプリケーションが別の区画の別のアプリケーションにデータを通信するので、論理的に独立したコンピュータのような区画の論理構成が維持される。多くの事例においては、これは、論理的に区画化されていないコンピュータ上で実行されて、論理的に区画化されたコンピュータとまさに同じように動作するように構成されたアプリケーションを可能にする。

【 0 0 6 2 】

各仮想 I O A 1 2 2 は、1 つ又はそれ以上のフレーム・バッファ 1 2 4 と関連し、それらは、非仮想 I O A と同様にイーサネット（登録商標）と互換性のあるデータ・フレームを送受信するのに用いられる。各フレーム・バッファ 1 2 4 は、バッファにおける別個のデータエリアを表す複数のスロット 1 2 6 を含む。各スロット 1 2 6 はまた、関連するクライアント 1 2 0 によって用いられるために、区画におけるローカル・バッファ（例えばバッファ 1 2 8 ）へのポインタを含む。

【 0 0 6 3 】

図示された実施形態においては、フレーム・バッファ 1 2 4 におけるスロット 1 2 6 は、種々のクライアント 1 2 0 によって順番に用いられることができる共有リソースを表す。多数のクライアント 1 2 0 がスロットを消費することができ、スロットは要求の処理の間にクライアントによって消費されるので、これらのリソースが消費される際にクライアントに通知する必要がある。図示された実施形態においては、例えば、仮想 I O A からクライアントに割り込みを発行することによってクライアントに通知される。以下でより明らかとなるように、本発明によれば、これらの試みは、要求ごとにではなくバッチで取り扱って、ディスパッチ・プロセスの効率を改善することができる。

【 0 0 6 4 】

図 4 に示された実施形態においては、各仮想 I O A 1 2 2 は、図 1 のリソース・マネージャ 5 2 と類似したリソース・マネージャとして動作する。したがって、各仮想 I O A 1 2 2 は、複数のサーバ・プロセス 1 3 0、ディスパッチ・プロセス 1 3 2、要求カウンタ 1 3 4、完了カウンタ 1 3 6、最後の完了カウンタ 1 3 8、及び最後のディスパッチ・カウンタ 1 4 0 を含む。さらに、各サーバ・プロセス 1 3 0 は、関連するサーバ・プロセスによって発せられた各要求を固有に識別するために、M Y _ R E Q U E S T 変数 1 4 2 を格納する。本発明に係るリソース管理は、典型的には、仮想 L A N からのデータの受信と関連して起こり、結果として、クライアント 1 2 0 は、データが関連する仮想 I O A によって受信される際の割り込みを介して通知される。

【 0 0 6 5 】

ここで図 5 に移ると、図 3 及び図 4 の環境のクライアント要求を取り扱うのに特に適したクライアント要求処理ルーチン 1 7 2 が、かなり詳しく示されている。ルーチン 1 7 2 は、多くの要求においてルーチン 7 2 と類似しているが、特に、仮想入力/出力アダプタのフレーム・バッファのスロットを消費するためのクライアント要求を取り扱うように構成されている。一般に、各クライアント要求は、フレーム・バッファにおける利用可能なスロットをアトミックに取得し、新しい固有の識別子が、要求カウンタをアトミックに増分することによってこうした要求に割当てられる。フレーム・バッファにおけるスロットの数を法とする要求識別子は、使用されるべき特定のスロットを得るために用いられる。この実装のためには、ここでは使用済みスロットカウンタと呼ばれる付加的なカウンタが

10

20

30

40

50

、消費されたスロットの総数を維持するために用いられる。このカウンタは、クライアントがリソースを解放するときにはいつも、典型的にはルーチン 172 以外のルーチンに実装された機能を介して減分される。

【0066】

ルーチン 172 は、使用済みスロットの総数が利用可能なスロットの最大数（総スロット定数又は変数によって表される）に等しいかどうかを判別することによって、ブロック 174 において始まる。等しい場合には、クライアントの要求が不履行となり、ブロック 176 に示されるようにエラーが戻される。そうでなければ、ブロック 174 は、制御をブロック 178 に移して、使用済みスロットの数を増分させる。次に、ブロック 180、182、184 及び 186 は、要求カウンタをアトミックに増分し、ルーチン 72 のブロック 74 - 80 と同様にして要求識別子を取得する。

10

【0067】

次に、制御はブロック 188 に移って、フレーム・バッファにおけるスロットの総数を法とする現在の要求識別子をとることによって、現在の要求のための現在のスロット（MY_SLOT）を設定する。次に、ブロック 190 において、現在のスロットへの要求が準備され、サブミットされ、それにより要求の処理が完了する。したがって、制御がブロック 192 に移って、カウンタ・ロックが取得され、次いで、ブロック 194 に移って、ルーチン 72 のブロック 84 及び 86 と同じようにして完了カウンタを増分させる。特に、ブロック 196、198、200、202、204、206、208 及び 210 に示されるようなルーチン 172 における残りの動作は、ルーチン 72 のブロック 88、90、92、94、96、98、100 及び 102 と関連して上記で開示されたのと同じであり、ブロック 204 における要求のディスパッチは、ディスパッチ・プロセスへの割り込みを介して起こり、通知される必要がある区画における種々のクライアントに付加的な割り込みを生成する。さらに、本発明に係る幾つかの実施形態においては、ルーチン 72 のブロック 100 及び 102 と同様に、ブロック 208 及び 210 は、ルーチン 172 から省略することができる。

20

【0068】

上記の実装における本発明の使用の実例として、4つの論理区画が仮想 LAN 上で互いに結合され、仮想 IOA が各論理区画に割り当てられ、3つの論理区画が第4の区画に送信される必要があるデータを有するケースを考える。第1の区画が送信されるべき64KBのデータを有し、第2の区画が送信されるべきたったの2KBのデータを有し、第3の区画が送信されるべき64KBのデータを有すると考える。このケースでは、第1、第2及び第3の区画は、それぞれの仮想 IOA に要求を発し、そして、これらの仮想 IOA が、第4の区画の仮想 IOA に送信されるべきデータのフレームをフォーマットするとみなす。さらに、仮想 IOA が、番号順に、すなわち第1、第2及び第3の区画からの要求が第4の区画の仮想 IOA によって順番に受信されるように、第4の区画の仮想 IOA に要求を発する。したがって、ルーチン 172 を用いて、要求に順番に固有の識別子が割り当てられ、第4の区画の仮想 IOA のフレーム・バッファのスロットが順番に消費される。

30

【0069】

しかしながら、要求に識別子が割り当てられると、仮想 IOA における各サーバ・プロセスは、例えば、そのプロセスによって消費されるスロットによって示されるクライアント・バッファへのそれぞれのフレームの DMA のために、その要求の取り扱いを並列に進める。しかしながら、図示された例においては、第2の区画からのデータ・フレームは第1の区画のデータ・フレームよりかなり小さい（2KB vs . 64KB）ので、第1の区画の要求が完了する前に第2の区画の要求が完了することが起こりうる。そうした場合、ルーチン 172 は、完了カウンタを増分させ、最後の完了カウンタを、第2の区画からの要求の固有の識別子に設定する。しかしながら、ルーチンはまた、ブロック 200 において、全ての先行の要求が完了したわけではないことに気付き、その結果、第1の区画からの要求が完了するまで第4の区画のクライアントへの通知を据え置くことになり、この時点

40

50

でバッチでディスパッチが行われて、両方のデータ・フレームの受信が第4の区画に通知されることになる。さらに、第3の区画からの要求がかなり遅れて完了する場合には、別のディスパッチが起こって、第4区画に第3データ・フレームの受信が通知される。

【図面の簡単な説明】

【0070】

【図1】本発明に係る形で複数のクライアントによる共有リソースへの同時アクセスをサポートする汎用コンピュータ・アーキテクチャのブロック図である。

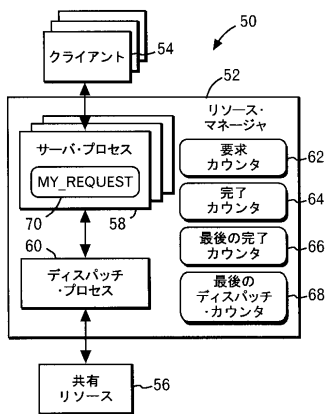
【図2】クライアント要求を処理するために図1のコンピュータ・アーキテクチャのサーバ・プロセスによって実行されるクライアント要求処理ルーチンのプログラム・フローを示すフローチャートである。

【図3】本発明に係る共有リソースの同時アクセスを実行する論理的に区画化されたコンピュータの基本ハードウェア・コンポーネントのブロック図である。

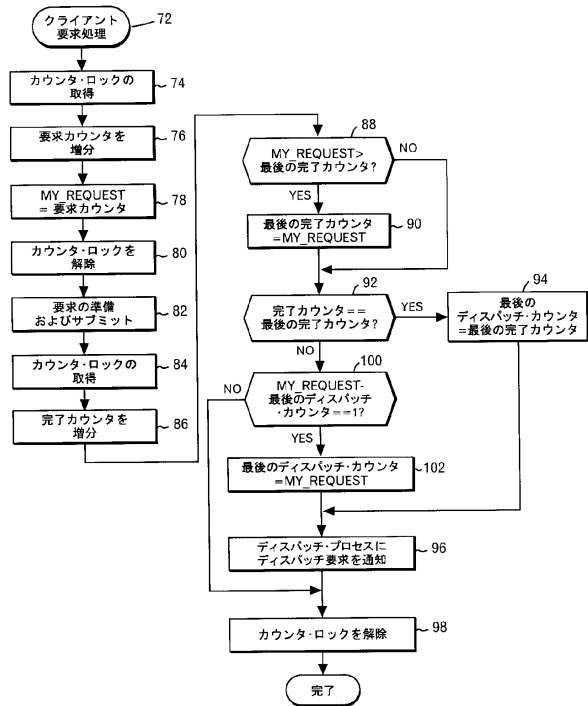
【図4】コンピュータの共有リソースへの同時アクセスを管理することに関係する図3のコンピュータの基本コンポーネントのブロック図である。

【図5】図4に関連するサーバ・プロセスによって実行されるクライアント要求処理ルーチンのプログラム・フローを示すフローチャートである。

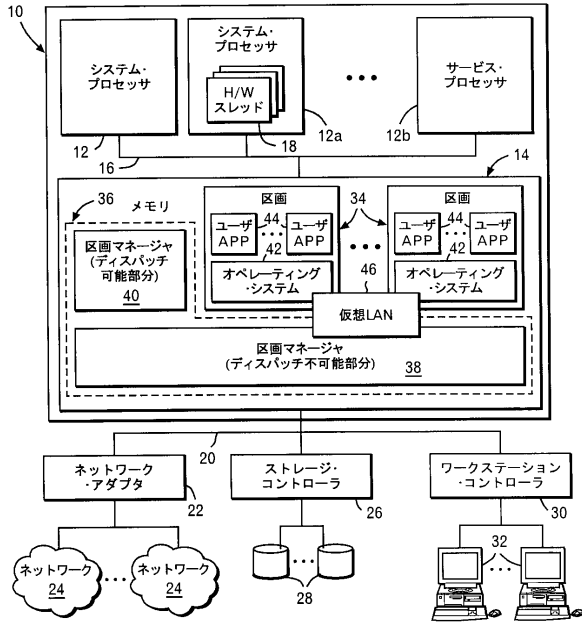
【図1】



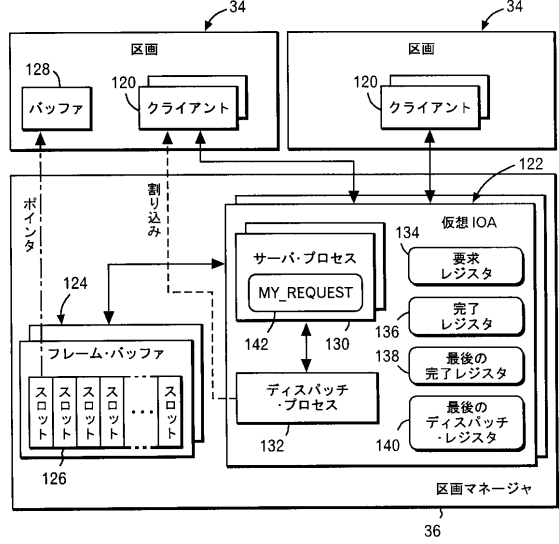
【図2】



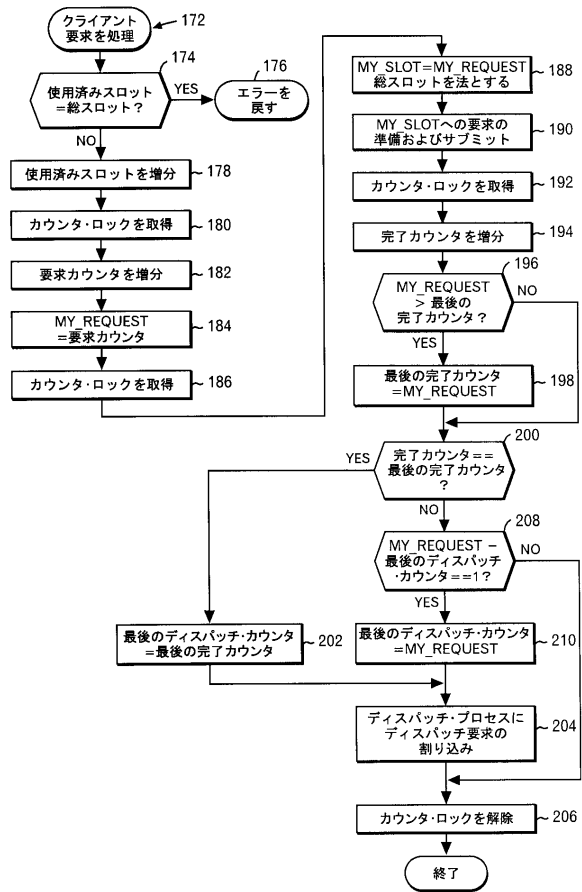
【図3】



【図4】



【図5】



フロントページの続き

(74)代理人 100086243

弁理士 坂口 博

(72)発明者 アームストロング、トロイ、デビッド

アメリカ合衆国 55902 ミネソタ州 ロチェスター アシュレイ・レーン・サウスウェスト
1310

(72)発明者 リュッケ、カイル、アラン

アメリカ合衆国 55906 ミネソタ州 ロチェスター キャシディ・ドライブ・ノースイース
ト 3023

審査官 鈴木 修治

(56)参考文献 特開2000-090057(JP,A)

特開2001-014178(JP,A)

特開平11-327931(JP,A)

特開平09-198265(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46-9/54