

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第3980487号

(P3980487)

(45) 発行日 平成19年9月26日(2007.9.26)

(24) 登録日 平成19年7月6日(2007.7.6)

(51) Int. Cl.		F I			
<b>G06F</b>	<b>9/46</b>	<b>(2006.01)</b>	G06F	9/46	350
<b>G06F</b>	<b>9/50</b>	<b>(2006.01)</b>	G06F	9/46	462Z

請求項の数 16 (全 26 頁)

(21) 出願番号	特願2002-571989 (P2002-571989)	(73) 特許権者	390009531
(86) (22) 出願日	平成14年2月1日(2002.2.1)		インターナショナル・ビジネス・マシー ズ・コーポレーション
(65) 公表番号	特表2004-530196 (P2004-530196A)		INTERNATIONAL BUSIN ESS MACHINES CORPO RATION
(43) 公表日	平成16年9月30日(2004.9.30)		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/GB2002/000456	(74) 代理人	100086243
(87) 国際公開番号	W02002/073397		弁理士 坂口 博
(87) 国際公開日	平成14年9月19日(2002.9.19)	(74) 代理人	100091568
審査請求日	平成15年10月3日(2003.10.3)		弁理士 市位 嘉宏
(31) 優先権主張番号	09/801,993	(74) 代理人	100108501
(32) 優先日	平成13年3月8日(2001.3.8)		弁理士 上野 剛史
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 区分処理環境におけるリソース平衡化

(57) 【特許請求の範囲】

【請求項1】

第1のオペレーティング・システムを含む第1のパーティションと、第2のオペレーティング・システムを含む第2のパーティションとを有するコンピューティング・システムにおける方法であって、

a) 前記第1のパーティションにおいて、前記第1のパーティションのネットワーク・トラフィック及び前記第1のパーティションのプロセッサの利用率を含む、前記第1のパーティションの第1スループット情報を収集するステップと、

b) 前記第1のパーティションにおいて、前記第1のパーティションの前記第1スループット情報に基づいて、前記第1のパーティションのプロセッサの累積利用率の変化に対する前記第1のパーティションの累積トラフィックの変化の割合である速度基準を計算するステップと、

c) 前記第1のパーティションから前記第2のパーティションのパーティション・マネージャに前記速度基準を伝達するステップと、

d) 前記パーティション・マネージャが、前記速度基準の値が、前記プロセッサの利用率が高く、かつ前記トラフィックが低いことを示す場合に、前記第1のパーティションに、少なくともプロセッサ、メモリ及びI/Oデバイスのいずれか1つを含むリソースをより多く割り振るよう指示するリソース平衡化指示を作成するステップと、

e) 前記リソース平衡化指示に従って、前記パーティション・マネージャにより前記第1のパーティションにリソースを割り振るステップと

10

20

を含む方法。

【請求項 2】

前記パーティション・マネージャは、前記第 2 のパーティションで実行されるワークロード・マネージャ、およびハイパーバイザを備える、請求項 1 に記載の方法。

【請求項 3】

パーティション間の通信は、パーティション間のメモリ共有を含む、請求項 1 に記載の方法。

【請求項 4】

前記コンピューティング・システムは、ネットワーク・アドレス及びオフセットを物理アドレスに変換する共有 I/O アダプタを更に有し、前記第 1 のパーティションは、ソース・レジスタにより指定される物理アドレスからデスティネーション・レジスタにより指定される物理アドレスヘデータを移動させるデバイス・ドライバを更に有し、

パーティション間の通信は、前記共有 I/O アダプタと前記デバイス・ドライバによる メモリからメモリへの直接のメッセージパッシングを含む、請求項 1 に記載の方法。

【請求項 5】

前記コンピューティング・システムは、ネットワーク・アドレス及びオフセットを物理アドレスに変換し、かつ、ソース・レジスタにより指定される物理アドレスからデスティネーション・レジスタにより指定される物理アドレスヘデータを移動させる共有 I/O アダプタを更に有し、

パーティション間の通信は、前記共有 I/O アダプタによる メモリからメモリへの直接のメッセージ・パッシングを含む、請求項 1 に記載の方法。

【請求項 6】

前記第 1 のパーティションのネットワーク・トラフィックについての情報はパケット・アクティビティ・カウンタを利用して取得される、請求項 1 に記載の方法。

【請求項 7】

前記第 1 のパーティションのネットワーク・トラフィックについての情報は、前記第 1 のパーティションに関連するネットワーク・パケットを計数することによって取得される、請求項 1 に記載の方法。

【請求項 8】

前記第 1 のパーティションに関連するネットワーク・パケットとは、前記第 1 のパーティションにおいて受信されるパケットである、請求項 7 に記載の方法。

【請求項 9】

前記第 1 のパーティションに関連するネットワーク・パケットとは、前記第 1 のパーティションにおいて送信されるパケットである、請求項 7 に記載の方法。

【請求項 10】

前記プロセッサの利用率はシステム・アクティビティ・カウンタを利用して取得される、請求項 1 に記載の方法。

【請求項 11】

第 1 のオペレーティング・システムを含む第 1 のパーティションと、第 2 のオペレーティング・システムを含む第 2 のパーティションとを備えるコンピューティング・システムであって、さらに、

a) 前記第 1 のパーティションのネットワーク・トラフィック及び前記第 1 のパーティションのプロセッサの利用率を含む、前記第 1 のパーティションの第 1 スループット情報を収集する手段と、

b) 前記第 1 のパーティションの前記第 1 スループット情報に基づいて、前記第 1 のパーティションのプロセッサの累積利用率の変化に対する前記第 1 のパーティションの累積トラフィックの変化の割合である速度基準を計算する手段と、

c) 前記第 1 のパーティションから 前記第 2 のパーティションのパーティション・マネージャに前記速度基準を伝達する手段と、

d) 前記パーティション・マネージャ中で、前記速度基準の値が、前記プロセッサの利用

10

20

30

40

50

率が高く、かつ前記トラフィックが低いことを示す場合に、前記第1のパーティションに、少なくともプロセッサ、メモリ及びI/Oデバイスのいずれか1つを含むリソースをより多く割り振ることを指示するリソース平衡化指示を作成する手段と、

e) 前記リソース平衡化指示に従って、前記パーティション・マネージャにより前記第1のパーティションにリソースを割り振る手段と  
を備えるシステム。

【請求項12】

前記パーティション・マネージャは、前記第2のパーティションで実行されるワークロード・マネージャ、およびハイパーバイザを備える、請求項11に記載のシステム。

【請求項13】

パーティション間の通信は、パーティション間のメモリ共有を含む、請求項11に記載のシステム。

【請求項14】

前記コンピューティング・システムは、ネットワーク・アドレス及びオフセットを物理アドレスに変換する共有I/Oアダプタを更に有し、前記第1のパーティションは、ソース・レジスタにより指定される物理アドレスからデスティネーション・レジスタにより指定される物理アドレスヘデータを移動させるデバイス・ドライバを更に有し、

パーティション間の通信は、前記共有I/Oアダプタと前記デバイス・ドライバによるメモリからメモリへの直接のメッセージ・パッシングを含む、請求項11に記載のシステム。

【請求項15】

前記コンピューティング・システムは、ネットワーク・アドレス及びオフセットを物理アドレスに変換し、かつ、ソース・レジスタにより指定される物理アドレスからデスティネーション・レジスタにより指定される物理アドレスヘデータを移動させる共有I/Oアダプタを更に有し、

パーティション間の通信は、前記共有I/Oアダプタによるメモリからメモリへの直接のメッセージ・パッシングを含む、請求項11に記載のシステム。

【請求項16】

第1のオペレーティング・システムを含む第1のパーティションと、第2のオペレーティング・システムを含む第2のパーティションとを有するコンピューティング・システム内で実行されるコンピュータ可読プログラムであって、前記プログラムは、

前記第1のパーティションを、

a) 前記第1のパーティションのネットワーク・トラフィック及び前記第1のパーティションのプロセッサの利用率を含む、前記第1のパーティションの第1スループット情報を収集する手段と、

b) 前記第1のパーティションの前記第1スループット情報に基づいて、前記第1のパーティションのプロセッサの累積利用率の変化に対する前記第1のパーティションの累積トラフィックの変化の割合である速度基準を計算する手段と、

c) 前記第1のパーティションから前記第2のパーティションのパーティション・マネージャに前記速度基準を伝達する手段として機能させ、

前記パーティション・マネージャを、

d) 前記速度基準の値が、前記プロセッサの利用率が高く、かつ前記トラフィックが低いことを示す場合に、前記第1のパーティションに、少なくともプロセッサ、メモリ及びI/Oデバイスのいずれか1つを含むリソースをより多く割り振ることを指示するリソース平衡化指示を作成する手段と、

e) 前記リソース平衡化指示に従って、前記第1のパーティションにリソースを割り振る手段として機能させる、プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

10

20

30

40

50

( 関連出願の相互参照 )

本出願は、以下の同時係属米国特許出願に関連する。

Baskey他による米国特許出願番号 09 / 801407 「INTER-PARTITIONMESSAGE PASSING METHOD, SYSTEM AND PROGRAM PRODUCT FOR THROUGHPUT MEASUREMENT IN A PARTITIONED PROCESSING ENVIRONMENT」( 整理番号 P O U 9 2 0 0 0 - 0 2 0 0 U S 1 )、

Baskey他による米国特許出願番号 09 / 802185 「INTER-PARTITIONMESSAGE PASSING METHOD, SYSTEM AND PROGRAM PRODUCT FOR A SHARED I/O DRIVER」( 整理番号 P O U 9 2 0 0 0 - 0 2 0 2 U S 1 )、および

Baskey他による米国特許出願番号 09 / 801492 「INTER-PARTITIONMESSAGE PASSING METHOD, SYSTEM AND PROGRAM PRODUCT FOR A SECURITY SERVER IN A PARTITIONED PROCESSING ENVIRONMENT」( 整理番号 P O U 9 2 0 0 1 - 0 0 1 2 U S 1 )。 10

【 0 0 0 2 】

本発明は、一般には、パーティションに区分されたデータ処理システムに関し、詳細には、システムのパーティション内で複数のオペレーティング・システム・イメージを実行することが可能なユニプロセッサ・システムおよびマルチプロセッサ・システムに関する。

【 背景技術 】

【 0 0 0 3 】

今日の中企業、大企業の大半は、かつては集中化された「温室」であったデータ・センターの到達範囲を自組織の範囲一杯まで、そして実際には範囲を越えて拡張するために各自のITインフラストラクチャを発展させている。このような発展の推進力の元となるのは、部分的には、従来は本質的に異なっていた部門別の業務を相互に結びつけ、納入業者や顧客とリアルタイムで意思の疎通を図りたいという要望であり、そして、電子商取引を行う媒体としてのインターネットの急速な成長と、それに伴い、そのような接続性を得るための相互接続とビジネス間のソリューションを次第に多く利用できるようになったことによりさらに加速されている。 20

【 0 0 0 4 】

この近年の発展に伴うのは、今日の企業が、多くの異なる作動プラットフォームを動的にリンクして相互に接続されたシームレスなシステムを作り出す必要性である。諸企業には、集中化していない仕入れ業務、適用分野に応じた要件、および合併に関連する活動から生じる本質的に異なる技術プラットフォームの作成などの要因による、異種の要素を含む情報システム・インフラストラクチャを備えることを特徴とすることが多い。さらに、納入業者、提携業者、および顧客間を結ぶリアルタイムの企業外部の接続性を促進したいという要望から、異種の環境に接続性を提供しようとするさらに強い動機が生じる。 30

【 0 0 0 5 】

急速に増大する顧客の様々な要求に対して、情報技術の提供者は、企業のデータ・センターへのより高い接続性を求めるニーズに対応するデータ処理ソリューションの考案に着手している。

【 0 0 0 6 】

本明細書の主題に関連する背景情報には以下が含まれる。コンピュータ・システムのアクティビティ分析について記載する、Ruffin他による米国特許出願番号 09 / 183961 「COMPUTATIONAL WORKLOAD-BASED HARDWARE SIZER METHOD, SYSTEM AND PROGRAM PRODUCT」; 論理パーティション間の共有メモリについて記載する、Temple他による米国特許出願番号 09 / 584276 「INTER-PARTITIONSHARED MEMORY METHOD, SYSTEM AND PROGRAM PRODUCT FOR A PARTITIONED PROCESSING ENVIRONMENT」; 高帯域幅の内蔵アダプタについて記載する、Baskey他による米国特許出願番号 09 / 253246 「A METHOD OF PROVIDING DIRECT DATA PROCESSING ACCESS USING QUEUED DIRECT INPUT-OUTPUT DEVICE」; システム中の2つの異なるクライアント・サーバのパーティショニングについて記載する、Temple他による米国特許出願番号 09 / 583501 「Heterogeneous Client Server Method, System and Program Product For A Partitioned Processing Environment」; 複数 40 50

パーティションのワークロード管理について記載する、IBM文書SG24-5326-00「OS/390 WorkloadManager Implementation and Exploitation」ISBN: 0738413070; および、ESA/390命令セット・アーキテクチャについて記載するIBM文書SA22-7201-06「ESA/390Principles of Operation」。これらの文書は、参照により本明細書に組み込む。

【0007】

当初、動作上の依存性を有する可能性がある各種のアプリケーションに対する処理サポートを同時に提供する統合システムを供給する必要性により、区分多重処理システムの市場が拡大した。単一の物理的なコンピューティング・システムで複数のオペレーティング・システム・イメージをサポートする能力を提供するこの区分システムは、かつてはメインフレーム・コンピュータ（IBM S/390システムなど）の分野でのみ用いられていたが、次第に広範囲の供給者から販売されるようになってきている。例えば、Sun Microsystems社は最近、Ultra Enterprise 10000ハイエンド・サーバでシステム・パーティショニングの一形態を提供し始めた。これについては、1996年12月12日に出願され、1999年8月3日に発行され、SunMicrosystems社に譲渡された、Drogichen他による米国特許第5,931,938号「Multiprocessor Computer HavingConfigurable Hardware System Domains」に詳細に記載される。他社もこの種のシステムに対する関心を示す方向性を表明している。

【0008】

この業界による採用は、1企業内の各種の計算ワークロードを1台（あるいは数台）の物理的なサーバ・コンピュータに統合する際に、そして動的に再構成することが可能なハードウェア環境でテスト・レベルおよび製造レベルのコードを同時に実装するために、システムの区分から得られる「systems within a system」の利点を明確に示すものである。さらに、前述の相互参照される特許出願で挙げたIBM S/390コンピュータ・システムなど、特定の区分多重処理システムでは、そのシステムで行う作業に割り当てられた優先度に応じて、リソース（プロセッサ、メモリ、およびI/O）を論理パーティション内および論理パーティション間で動的に割り振ることができる（IBMおよびS/390は、InternationalBusiness Machines社の登録商標）。このワークロードの優先度に応じて動的にリソースを割り振る能力により、長年にわたる容量計画の問題が対処される。従来は、データ・センターの管理者が、一時的なワークロードの急増に対処するために予想される計算ワークロードに過剰な量のリソースを意図的に指定することによりこの問題に対応していた。

【特許文献1】米国特許出願番号09/801407

【特許文献2】米国特許出願番号09/802185

【特許文献3】米国特許出願番号09/801492

【特許文献4】米国特許出願番号09/183961

【特許文献5】米国特許出願番号09/584276

【特許文献6】米国特許出願番号09/253246

【特許文献7】米国特許出願番号09/583501

【特許文献8】米国特許第5,931,938号

【特許文献9】米国特許第5,442,802号

【特許文献10】米国特許出願番号09/411417

【特許文献11】米国特許第5,414,851号

【特許文献12】米国特許第5,269,009号

【特許文献13】米国特許出願番号09/677338

【非特許文献1】IBM文書SG24-5326-00「OS/390 WorkloadManager Implementation and Exploitation」ISBN: 0738413070

【非特許文献2】IBM文書SA22-7201-06「ESA/390Principles of Operation」

【非特許文献3】IBM文書SC28-1855-06「OS/390V2R7.0 OSA/SF User's

10

20

30

40

50

Guide」

【非特許文献4】G 3 2 1 - 5 6 4 0 - 0 0 「S/390 clustertechnology: Parallel Sysplex」

【非特許文献5】IBM文献SC 3 4 - 5 3 4 9 - 0 1 「MQSeries QueueManager Clusters」

【非特許文献6】「Coupling Facility ConfigurationOptions: A Positioning Paper」(GF 2 2 - 5 0 4 2 - 0 0, IBM社)

【非特許文献7】IBM Redbook Document SG24-4810-01 「UnderstandingRS/6000 Performance and Sizing」

【発明の開示】

10

【発明が解決しようとする課題】

【0009】

こうした区分システムは、企業内の異種のシステムを取り込むようにデータ・センターを拡張することを容易にするが、現在、こうしたソリューションは、異種または同質の区分プラットフォームを相互運用性がある単一の区分システムに機能的に統合する単純な機構は提供しない。実際、これらの新しいサーバは、単一の物理的ハードウェア・プラットフォーム内でのオペレーティング・システム・イメージの統合は可能にするが、サーバの複数のパーティションに存在するオペレーティング・システム間の相互運用性に対するニーズには充分に対応していない。この相互運用性の問題はさらに、各種のパーティション内に異種のオペレーティング・システムを有する異種要素を含むシステムにおいて難しくなる。また、こうしたシステムは通例、そのような異種のプラットフォーム間でのパーティション間のリソース共有のタイプにも対応していない。パーティション間のリソース共有が実現すると、パーティション間の高帯域幅、低待ち時間の相互接続が可能になる。これらの問題に対するシステム統合のソリューションにより、別個のパーティションで実行されるプロセス間の通信を行うためのより堅固な機能を可能にして、そのようなアプリケーションが個別のオペレーティング・システムで実行されているとき、それらのアプリケーションは実際には互いに対してローカルであるという事実を活用できることから、この相互運用性の問題に対処することが重要である。

20

【0010】

上述のTemple他による「INTER-PARTITION SHARED MEMORYMETHOD, SYSTEM AND PROGRAM PRODUCT FOR A PARTITIONED PROCESSING ENVIRONMENT」という名称の米国特許出願番号09/584276では、数個のオペレーティング・システムの「カーネル」を拡張することにより、共有ストレージの使用を容易にしてパーティションを越えたメモリ共有を実現する。「カーネル」とは、オペレーティング・システム中の中核的なシステム・サービス・コードである。ネットワーク・メッセージ・パッセージのプロトコルは、そのようにして作成したインタフェース上で実装することができるが、しばしば、1つまたは複数のオペレーティング・システムの変更に頼ることなく効率的なプロセス間通信を可能にすることが求められる。また、上述のTemple他による米国特許出願番号09/584276、または米国特許第5,931,938号に記載されるSunMicrosystems社のUltra Enterprise 10000ハイエンド・サーバのように、メモリ領域を共有するためにパーティションの分離を制限せずに済むことも望まれる。同時に、ネットワーク速度ではなくメモリ速度によりパーティション間で情報をやり取りすることも望まれる。したがって、アドレスを共有せずにパーティション・メモリ間でメモリを移動する方法が求められる。

30

40

【0011】

IBM S/390 Gbit Ethernet (R) (1995年8月15日に発行され、IBM社に譲渡された米国特許第5442802号「Asynchronous Coprocessor Data Mover Method and Means」)のI/Oアダプタを使用して、あるパーティションのカーネル・メモリから別のパーティションのカーネル・メモリにデータを移動することができるが、データは、第1のカーネル・メモリからアダプタのキュー・バッファに移動され、次いでアダプタの第2のキュー・バッファに転送してから第2のカーネル・メモリに

50

転送される。すなわち、メモリからメモリに転送する際に合計3回のデータ移動がある。どのメッセージ・パッシング通信方式でも、データ・アクセスの待ち時間が共有ストレージへの一度の格納とストレージからの取り出しの待ち時間に近づくように、データ移動操作の回数を最少に抑えることが望ましい。移動機能は、転送するデータ・ブロックごとに3回のデータ移動操作を有する。この操作のうち1つまたは2つの操作を除去する方法が望まれる。

#### 【0012】

同様に、IBM S/390 Parallel Sysplex Coupling Facilityマシンは、パーティション間のメッセージ・パッシングを容易にするために使用することができ、また実際に使用されている。しかし、この場合、データの転送は、第1のカーネル・メモリから結合機構へ、そして結合機構から第2のメモリへと行われる。これには、求められる1回の移動ではなく2回のデータ操作が必要となる。

10

#### 【0013】

多くのコンピュータ・システムでは、無許可あるいは不当なアクセスを通じたマシン上のデータおよびアプリケーションの不正使用を防止するために、ユーザの識別を確認することが望まれる。各種のオペレーティング・システムおよびアプリケーション・システムは、この目的のためにユーザ認証および他のセキュリティ・サービスを備える。区分システム、あるいは実際には任意のシステムのクラスタまたはネットワークに入ろうとするユーザを、エントリ時に1回のみ、あるいは重要なリソースに対する要求時や重要なシステム保守機能の実行時などの重要なチェックポイントで確認することが望ましい。この要望は「シングル・サインオン」の必要性として知られる。このために、各種パーティションのセキュリティ・サーバは相互に対話するか、または統合しなくてはならない。この例には、OS/390におけるウェブから受信した「デジタル証明書」を処理し、それを従来のユーザIDおよびパスワード検証および権利付与にマッピングするOS/390 SAF(RACF)インタフェースの強化や、Kerberosセキュリティ・サーバ、および台頭しつつあるディレクトリ・サービス向けのLDAP規格がある。

20

#### 【0014】

さらに、電子商取引は競争的な性質を持つため、ユーザ認証および権利付与のパフォーマンスは、従来のシステムに比べて重要性が高い。労働者は、業務の開始時に認証されるまで待つことを予期しているかもしれないが、顧客は、認証に時間がかかり過ぎる場合には単によそに行ってしまうことが考えられる。ウェブは公衆性があるので、暗号化の使用はこの問題を難しくする。また、しばしば、他のオペレーティング・システムのためには書かれていないデバイス・ドライバがあるオペレーティング・システムに存在することがある。そのような場合には、別のパーティションからあるパーティションのデバイス・ドライバに効率的な形でインタフェースを取ることが望ましい。現在このタイプの動作に利用できるのはネットワーク接続だけである。

30

#### 【0015】

分散システムに伴う問題の1つは、他のシステムが過度に利用されている一方での、「ホワイト・スペース」、すなわちあるシステムで十分に利用されていないリソースの管理である。システム間あるいはシステム・イメージ間でワークロードを移動するIBMのLoad LevelerやOS/390オペレーティング・システムのワークロード・マネージャが備えるParallel Sysplex機能などのワークロード・バランスがある。区分コンピューティング・システムでは、パーティション間でワークロードではなくリソースを移すことが可能であり、また望ましい。これは、機能の変更に伴う大量のコンテキストの切り替えとデータ移動が回避されることから望ましい。

40

#### 【0016】

外部のSysplexのクラスタ化した接続を使用してUNIX(R)オペレーティング・システムのソケット間接続を実装するIBM S/390の「Sysplex Sockets」は、従来技術の一部の一例である。この技術では、あるサービスが、利用可能なセキュリティのレベルを示し、アプリケーションが示す要求されるセキュリティ・レ

50

ベルに基づいて接続をセットアップする。しかし、この場合は、より高いレベルのセキュリティのために暗号化が提供され、また S y s p l e x 接続自体が、本発明によって実装されるメモリ接続に比べてはるかに深い物理トランスポート層を有する。

【 0 0 1 7 】

同様に、SSLによる認証を提供し、証明書情報を(プロキシとして)ウェブ・アプリケーション・サーバに提供するウェブ・サーバは、本発明のメモリ共有またはメモリからメモリへの直接のメッセージが効果的に使用される別の例と考えることができる。ここでは、プロキシは、セキュリティ・サーバに渡すデータを再度暗号化する必要がなく、さらに、管理すべき深い接続インタフェースを持たない。実際、当業者には、本発明のこの実施形態では、プロキシ・サーバは、基本的に、セキュリティ・サーバと同じオペレーティング・システム下で実行されるプロキシ・サーバと実質的に同じプロセスを通じてセキュリティ・サーバと通信することが理解されよう。Baskey他による米国特許出願番号09/411417「Methods, Systems and Computer Program Products for Enhanced Security Identity Utilizing an SSL Proxy」には、プロキシ・サーバを使用してセキュアHTTPプロトコルでセキュア・ソケット・レイヤ(SSL)を行うことについて述べられる。

10

【課題を解決するための手段】

【 0 0 1 8 】

本発明の一態様によれば、処理システムは、少なくとも第1のパーティションおよび第2のパーティションを含む。このシステムは、すべてのパーティションと通信するパーティション・リソース・マネージャを備える。第1のパーティションで実行されるプログラムは、第1のパーティションのネットワーク・トラフィック及び第1のパーティションのプロセッサの利用率を含む第1のパーティションの第1スループット情報を収集し、第1スループット情報に基づいて、第1のパーティションのプロセッサの累積利用率の変化に対する第1のパーティションの累積トラフィックの変化の割合である速度基準を計算する。パーティション・リソース・マネージャは、第1のパーティション、すなわち第1のパーティションで実行されるプログラムから、第1のパーティションの速度基準を受け取り、第1のパーティションのリソースに対してリソース平衡化の指示を決定する。ここで、リソース平衡化の指示は、速度基準の値に基づいて作成される。リソース・マネージャは、リソース平衡化の指示に従って第1のパーティションにリソースを割り振る。

20

30

【 0 0 1 9 】

好ましい実施形態では、パーティション・リソース・マネージャ機能は、ハイパーバイザと併せてワークロード・マネージャを使用して実装すると有利である。このシステムは、すべてのパーティションと通信するハイパーバイザを備える。第2のパーティションにはワークロード・マネージャが提供され、第1のパーティション、すなわち第1のパーティションで実行されるプログラムから速度基準についての情報を受け取り、リソース平衡化の指示を決定する。コミュニケータは、ワークロード・マネージャからハイパーバイザにリソース平衡化の指示を伝達する。第1のパーティションのカーネルは、ハイパーバイザから受け取るリソース平衡化の指示に従って第1のパーティションにリソースを割り振る。

【 0 0 2 0 】

共有されるメモリ・リソースは、複数のパーティションで実行される相互動作する複数プロセスのための指定メモリ・リソースに独立してマッピングすることが好ましい。このようにして、メモリ・リソースを共有する各パーティション中のプロセスによって共通の共有メモリ空間をマッピングすると、共有メモリは、パーティション内でそのプロセスに割り当てられ、通常のプロセス実行の過程にデータの読み書きに使用できるメモリ・リソースのように見える。

40

【 0 0 2 1 】

さらなる実施形態では、プロセスは相互に依存しており、共有メモリ・リソースは、いずれかまたは両方のプロセスによる後のアクセスのために、いずれかまたは両方のプロセスから格納を行うことができる。

50

## 【 0 0 2 2 】

さらなる実施形態では、システムは、パーティション中の各種プロセスを共有メモリ空間に接続するプロトコルを含む。

## 【 0 0 2 3 】

別の実施形態では、区分に関係なくすべての物理メモリへの物理的アクセスを備える I/O アダプタにより、あるパーティションのカーネル空間から別のパーティションのカーネル空間への直接のデータ移動が可能になる。I/O アダプタのすべてのメモリにアクセスする能力は、パーティション間の I/O リソース共有を可能にする区分コンピュータ・システム中の機能から当然得られる結果である。このような共有については、1995 年 5 月 9 日に発行された米国特許第 5,414,851 号「METHOD AND MEANS FOR SHARING I/O RESOURCES BY A PLURALITY OF OPERATING SYSTEMS」に記載される。ただし、アダプタは、データ・ムーバを使用することにより、あるパーティションのメモリから別のパーティションのメモリに直接データを移動する能力を有する。

10

## 【 0 0 2 4 】

さらなる実施形態では、カーネル・メモリ間のデータ移動のための機能は、ネットワーク通信アダプタのハードウェアおよびデバイス・ドライバ内に実装される。

## 【 0 0 2 5 】

さらなる実施形態では、メモリ間インタフェースを通じたローカルであるが異種の要素を含むセキュアな接続のためにそれぞれが最適化された TCP/IP スタックから、ネットワーク・アダプタを駆動する。

20

## 【 0 0 2 6 】

別の実施形態では、データ・ムーバ自体を区分処理システムの通信ファブリック中に実装し、I/O アダプタによって制御して、さらに直接的なメモリからメモリへの転送を容易にする。

## 【 0 0 2 7 】

さらに別の実施形態では、オペランドとして供給されるネットワーク・アドレスおよびオフセットを物理アドレスに変換することができる、特権 CISC 命令のマイクロコードによってデータ・ムーバを制御し、それにより、データ・ムーバは、実際のアドレスと仮想アドレスを 2 つのパーティションに有する物理アドレス間で、move character long 命令 (IBM S/390 の MVCL 命令。IBM 文書 SA22-7201-06 "ESA/390 Principles of Operation" を参照) に相当する操作を行う。

30

## 【 0 0 2 8 】

さらに別の実施形態では、すべての物理メモリへの仮想メモリ・アクセスおよび実際のメモリ・アクセスを行い、オペランドとして供給されるネットワーク・アドレスおよびオフセットを物理アドレスに変換することができるハイパーバイザで実行されるルーチンによってデータ・ムーバを制御し、それによりデータ・ムーバは、実際のアドレスと仮想アドレスを 2 つのパーティションに有するアドレス間で move character long 命令 (IBM S/390 MVCL) に相当する操作を行う。

## 【 0 0 2 9 】

パーティションの 1 つでサーバ・プロセスを実装し、他のパーティションでクライアント・プロセスを実装することにより、区分システムは、異種の要素からなる単一システムのクライアント・サーバ・ネットワークを実装することができる。既存のクライアント/サーバ・プロセスは通例ネットワーク・プロトコル接続によって相互動作するので、そうしたプロセスは、本発明のメッセージ・パッシングの実施形態で容易に実装され、インタフェースの変更に頼ることなくパフォーマンスとセキュリティの利点が得られる。ただし、共有メモリの実施形態でのクライアント/サーバ・プロセスの実装は、パフォーマンスまたは配置の速度、あるいはその両方で有利である可能性がある。

40

## 【 0 0 3 0 】

さらなる実施形態では、共有メモリまたはメモリ間のメッセージ・パッシングを利用するアプリケーション・サーバに対して、信頼性のある/保護されたサーバ環境が提供され

50

る。この環境では、現行技術のような追加的な暗号化または認証を必要とせずに、外部化を伴う認証および認証データのセキュリティの露出が回避される。

【0031】

特定の実施形態では、ウェブ・サーバは、OS/390、Z/OS、またはVM/390の下で実行される「SAF」セキュリティ・インタフェースと、メモリ・インタフェースを通じて通信するOS/390向けのLinuxで実行されるLinux Apacheである。この実施形態では、Linuxの「プラグ可能認証モジュール」を変更して、メモリ接続を通じてSAFインタフェースを駆動する。

【0032】

さらなる実施形態では、セキュリティ・クリデンシャル/コンテキストを共有メモリに格納するか、メモリ間の転送を介して複製するように、Policy DirectorあるいはRACFのようなセキュリティ・サーバを変更する。

10

【0033】

このように、本発明の好ましい実施形態は、複数の異種のオペレーティング・システム・イメージをサポートすることが可能な区分されたコンピュータ・システムを含み、それらのオペレーティング・システム・イメージは、メモリ位置を共有せずに、それらのメモリ位置間でメモリ速度により同時にメッセージをやり取りすることができる。これは、あるパーティションのあるカーネル・メモリ空間から第2のパーティションのカーネル・メモリ空間への直接のデータ移動を協働して容易にする特殊なデバイス・ドライバを備えるI/Oアダプタを使用して行う。

20

【0034】

次いで添付図面を参照して単なる例として本発明の一実施形態を説明する。

【発明を実施するための最良の形態】

【0035】

本発明の好ましい実施形態の特定の態様について述べる前に、区分処理システムの基本的な構成要素を概説すると有用であろう。これを背景知識として使用することにより、本発明の特定の有利な特徴をどのように区分システムで用いてそのパフォーマンスを向上できるかをより深く理解することができる。IBM文書SC28-1855-06「OS/390 V2R7.0 OSA/SF User's Guide」を参照されたい。この書籍は、OS/390オペレーティング・システムの一要素であるOpenSystems Adapter Support Facility (OSA/AF)の使用法について記載する。この書籍には、OSA/SFをセットアップし、OS/2インタフェースまたはOSA/SFコマンドのいずれかを使用してOSAをカスタマイズし、管理するための指示が提供される。G321-5640-00「S/390cluster technology: Parallel Sysplex」は、汎用目的の大規模な商用マーケットプレース向けに開発されたクラスタ化マルチプロセッサ・システムについて述べる。S/390 Parallel Sysplexシステムは、高度にスケーラブルなクラスタ化コンピューティング環境で完全なデータ共有と並列処理を組み合わせるために設計されたアーキテクチャに基づく。Parallel Sysplexシステムは、コスト、パフォーマンス範囲、および可用性の面で著しい利点をもたらす。IBM文献SC34-5349-01「MQSeriesQueue Manager Clusters」は、MQSeriesのキュー・マネージャ・クラスタについて記述し、クラスタの概念、用語、および利点を説明する。この文献は、新規のコマンドと変更されたコマンドの構文を要約し、キュー・マネージャのクラスタをセットアップし、維持するタスクの例をいくつか挙げる。IBM出版SA22-7201-06「ESA/390Principles of Operation」は、参照のためにESA/390アーキテクチャの詳細な定義を含む。この書籍は主にアセンブラ言語のプログラマが使用する参考文献として書かれたものであり、各機能に依存するアセンブラ言語プログラムを作成するのに必要とされる詳細のレベルで各機能を説明しているが、ESA/390の機能面の詳細に関心を持つ者にとっては有用であろう。

30

40

【0036】

上述の文献は、当技術分野の現行状態の事例を提供し、本発明の背景を理解する助けと

50

なろう。これらの参考文献は、参照により本明細書に組み込まれる。

【0037】

図1に、区分処理システム100を構成する基本要素を示す。システム100は、ブロックAおよびBとして図示するブロックに区分することが可能な物理的なメモリ・リソースからなるメモリ・リソース・ブロック101と、区分されたメモリ・リソース101と一致するように論理的または物理的に区分することができる1つまたは複数のプロセッサからなることが可能なプロセッサ・リソース・ブロック102と、同様に区分することが可能な入出力(I/O)リソース・ブロック103とからなる。これらの区分されたリソース・ブロックは、スイッチ・マトリクスなどからなる相互接続ファブリック104を介して相互に接続される。相互接続ファブリック104は、プロセッサ102Bをメモリ101Bに接続するなど、パーティション内でリソース同士を相互接続する機能を果たすことができ、また、プロセッサ102Aをメモリ101Bに接続するなどパーティション間のリソースを相互に接続する役割も果たすことができることは理解されよう。本明細書で使用する用語「ファブリック」は、システムの諸要素を相互接続するための当技術分野で知られる一般的な方法を意味する。ファブリックは、単純な二地点間バスであっても高度なルーティング機構であってもよい。この図には2つのパーティション(AおよびB)を有するシステムを示しているが、この表現は説明を簡略にするために選択したものであり、さらに、本発明は、使用可能なリソースとパーティション技術が可能にする数のパーティションを実装するように構成可能なシステムを包含することは容易に認識されよう。

10

【0038】

図を考察すると、個別に見たパーティションAおよびBはそれぞれ、単独のデータ処理システムの構成要素、すなわちプロセッサ、メモリ、およびI/Oからなることが容易に理解されよう。この事実は、区分された処理システムに、独自の「systems within a system」の利点を与える特徴である。事実、そして本明細書で例証するように、現在利用可能なパーティション処理システムを分ける主要な差異は、システム・リソースを区分することが可能な境界と、その境界を越えてパーティション間でリソースを移動できる容易さである。

20

【0039】

パーティションを分ける境界が物理的な境界である第1のケースの最適な例となるのはSunMicrosystemsのUltra Enterprise 10000システムである。Ultra Enterprise 10000システムでは、パーティションは物理的な境界に従って区画され、具体的には、ドメインまたはパーティションは、それぞれがいくつかのプロセッサ、メモリ、およびI/Oデバイスからなる1つまたは複数の物理的なシステム・ボードからなる。ドメインは、そのシステム・ボードの1つまたは複数、およびそのボードに装着されたI/Oアダプタとして定義される。ドメインは、独自開発のバスおよびスイッチ・アーキテクチャによって相互に接続される。

30

【0040】

図2に、物理的に区分された処理システム200を構成する要素の高レベル表現を示す。図2の参照から分かるように、システム200は2つのドメインあるいはパーティションAおよびBを含む。パーティションAは、2つのシステム・ボード202A1および201A2からなる。パーティションAの各システム・ボードは、メモリ201A、プロセッサ202A、I/O203A、および相互接続媒体204Aを含む。相互接続媒体204Aは、システム・ボード201A1上の構成要素が相互と通信することを可能にする。同様に、単一のシステム・ボードからなるパーティションBも同様の構成処理要素、すなわちメモリ201B、プロセッサ202B、I/O203B、および相互接続204Bを含む。パーティションに分けられたシステム・ボードに加えて、各システム・ボードに結合され、パーティション内のシステム・ボード間の相互接続と、異なるパーティションのシステム・ボードの相互接続を可能にする相互接続ファブリック205が存在する。

40

【0041】

次のタイプのシステム・パーティションを論理的パーティショニングと呼ぶ。このシス

50

テムでは、各種のパーティションへのリソースの割り当てを制約する物理的な境界がなく、システムは、その物理的な場所に関係なく任意のパーティションに割り当てることができる使用可能なリソースのプールを備えるものと考えることができる。この点が、例えば所与のシステム・ボード（システム・ボード 201A1 など）上のすべてのプロセッサが必然的に同じパーティションに割り当てられる、物理的に区分されたシステムとの違いである。IBM AS/400 システムは、論理的に区分された専用のリソース処理システムの例である。AS/400 システムでは、ユーザは、プロセッサ、メモリおよび I/O をそれらの物理的な場所に関係なく所与のパーティションに含めることができる。したがって、例えば、物理的に同じカード上に位置する 2 つのプロセッサを 2 つの異なるパーティションのリソースとして指定することができる。同様に、カードなど所与の物理的なパッケージ中のメモリ・リソースは、そのアドレス空間の一部を論理的に 1 パーティションのために確保し、残りのアドレス空間を別のパーティションのために確保することができる。

10

#### 【0042】

AS/400 システムなど論理的に区分された専用リソース・システムの特徴は、パーティションへのリソースの論理的マッピングが、手動によるシステムの再構成によるみ変更することができる、静的に行われる割り当てである点である。図 3 を参照すると、プロセッサ 302A1 は、物理的にはシステム中のどこに位置してもよく、論理的にはパーティション A のために確保されたプロセッサを表す。ユーザが、プロセッサ 302A1 をパーティション B にマッピングし直したい場合は、プロセッサをオフライン状態にし、変更に対応するように手動で再度マッピングしなければならない。論理的に区分されたシステムは、例えば固定数のプロセッサをサポートするシステム・ボードなどの物理的な区分の境界による制限を受けないので、リソースの区分により高い精度を提供する。ただし、そのような論理的に区分された専用のリソース・システムの再構成を行うには、パーティションの再マッピングを行うリソースの動作を中断しなくてはならない。したがって、そのようなシステムでは、物理的に区分されたシステムに伴う制限の一部は回避されるが、なおパーティション間のリソースの静的なマッピングに伴う再構成の制限があることが分かる。

20

#### 【0043】

この結果、我々は、論理的に区分された共有リソース・システムの検討に至った。そのようなシステムの一例は、IBM S/390 コンピュータ・システムである。論理的に区分された共有リソース・システムの特徴の 1 つは、プロセッサなど論理的に区分されたリソースを複数のパーティションで共有できる点である。この特性により、論理的に区分された専用リソース・システムの再構成の制約が効果的に克服される。

30

#### 【0044】

図 4 に、論理的に区分されたリソース共有システム 400 の概略的な構成を示す。論理的に区分された専用のリソース・システム 300 と同様に、システム 400 は、システム中の物理的な場所に関係なく任意のパーティション（この例では A または B）に論理的に割り当てることが可能なメモリ 401、プロセッサ 402、および I/O リソース 403 を含む。ただし、システム 400 から分かるように、特定のプロセッサ 402 または I/O 403 の論理パーティションへの割り当ては、「ハイパーバイザ」（408）で実行されるスケジューラに従って仮想プロセッサ（406）および I/O ドライバ（407）を入れ替えることにより動的に変更することができる。（ハイパーバイザは、仮想マシンのリソースをスケジューリングし、割り振る管理プログラムである。）プロセッサおよび I/O を仮想化することにより、適切な優先順位付けを行って全オペレーティング・システム・イメージを動作中および非動作時に入れ替えることが可能になり、パーティションがそれらのリソースを動的に共有できるようになる。

40

#### 【0045】

論理的に区分された共有リソース・システム 400 は、プロセッサと I/O リソースを共有する機構を提供するが、パーティション間のメッセージ・パッシングは、既存のシス

50

テムでは完全には対応されていなかった。これは、既存の区分システムではパーティション間の通信を可能にすることができないということではない。実際、そのような通信はここに記載する各タイプの区分システムで行われている。しかし、それらの実装で、ハイパーバイザ、共有メモリの実装、あるいはアダプタやチャネル通信デバイスの標準的なセット、あるいはパーティション間を接続するネットワークの介在なしに、カーネル・メモリからカーネル・メモリにデータを移動する手段を提供するものはない。

**【 0 0 4 6 】**

米国特許第 5 , 9 3 1 , 9 3 8 号に記載される Sun Microsystems Ultra Enterprise 100 00 システムを代表とする、物理的に区分された多重処理システムでは、マスク・レジスタを適切に設定することにより、システム・メモリの一領域を、ハードウェア・レベルの複数のパーティションからアクセス可能にすることができる。この Sun 特許は、パーティション間ネットワークのためのバッファリング機構および通信手段として使用できることを指摘する以外には、この機能をどのように活用するかを教示しない。上述の Temple 他による米国特許出願番号 0 9 / 5 8 4 2 7 6 は、異種の要素からなる区分システムにおいて共有メモリ機構をどのように構築し、活用するかを教示する。

**【 0 0 4 7 】**

「Coupling Facility Configuration Options: A Positioning Paper」( G F 2 2 - 5 0 4 2 - 0 0 , I B M 社 ) に詳細に記述されるように、I B M S / 3 9 0 システムでは、共通にアドレス指定される物理メモリを「統合された結合機能」として使用するのための同様の内部のクラスタ化機能について記載される。ここでは、共有ストレージは実際にリポジトリであるが、そこへの接続は、X C F と称される I / O 様のデバイス・ドライバを通じて行われる。ここでは、共有メモリは結合機能中に実装されるが、S / 3 9 0 でないオペレーティング・システムがそれを使用するには拡張を作成することが必要である。さらに、この実装では、1 パーティションのカーネル・メモリから結合機能のメモリにデータを移動させ、そして第 2 のパーティションのカーネル・メモリにデータを移動させる。

**【 0 0 4 8 】**

カーネルとは、ハードウェア・リソースの割り振りなどの基本的機能を行うオペレーティング・システムの部分である。カーネル・メモリとは、カーネルが自身の機能を実行するために使用する、カーネルが利用できるメモリ空間である。

**【 0 0 4 9 】**

これに対し、この実施形態は、どちらかのパーティションまたはハードウェア中でオペレーティング・システムに共有ストレージの拡張を提供することなく、新規な I / O アダプタとそのデバイス・ドライバのイネープリング機能を使用して、あるパーティションのカーネル・メモリから別のパーティションのカーネル・メモリに 1 回の操作でデータを移動する手段を提供する。

**【 0 0 5 0 】**

ここで説明する実施形態の動作を理解する助けとして、オペレーティング・システム中のプロセス間通信を理解すると有用である。図 5 を参照すると、プロセス A ( 5 0 1 ) および B ( 5 0 3 ) は、それぞれアドレス空間メモリ A ( 5 0 2 ) およびメモリ B ( 5 0 4 ) を有する。これらのアドレス空間には、カーネル ( 5 0 5 ) によるシステム・コールの実行により実際のメモリが割り振られている。カーネルは、独自のアドレス空間であるメモリ K ( 5 0 6 ) を有する。1 つの通信形態では、プロセス A および B は、バッファ 5 1 0 を作成し、そこに接続し、アクセスする適切なシステム・コールを行うことにより、メモリ K にバッファ 5 1 0 を作成することで通信する。これらの呼び出しの意味は、システムごとに異なるが効果は同じである。第 2 の通信形態では、メモリ S ( 5 0 7 ) のセグメント 5 1 1 をメモリ A ( 5 0 2 ) およびメモリ B ( 5 0 4 ) のアドレス空間にマッピングする。このマッピングが完了すると、プロセス A ( 5 0 1 ) および B ( 5 0 3 ) は、両プロセスが理解する任意のプロトコルに従って、メモリ S ( 5 0 7 ) の共有セグメントを自由に使用することができる。

**【 0 0 5 1 】**

米国特許出願番号 09 / 583501 「Heterogeneous ClientServer Method, System and Program Product For A Partitioned ProcessingEnvironment」を図 6 に表し、ここではプロセス A ( 601 ) および B ( 603 ) が、異なるオペレーティング・システム・ドメイン、イメージ、あるいはパーティション ( パーティション 1 ( 614 ) およびパーティション 2 ( 615 ) ) に存在する。この場合には、メモリ K1 ( 606 ) およびメモリ K2 ( 608 ) を各自のカーネル・メモリとして有するカーネル 1 ( 605 ) およびカーネル 2 ( 607 ) がある。メモリ S ( 609 ) は、ここではパーティション 1 およびパーティション 2 の両方からアクセス可能な物理メモリの空間である。このような共有は、これらに限定しないが、UE10000 のメモリ・マッピング実装、または S / 390 のハイパーバイザ実装、あるいは区分によって生じるアクセスの障壁を制限する他の任意の手段を含む任意の実装によって可能にすることができる。これに代わる例として、共有メモリを最も高い物理メモリ・アドレスにマッピングし、コンフィギュレーション・レジスタ中の先頭アドレスによって共有空間が定義される。

#### 【 0052 】

慣例として、メモリ S ( 609 ) は、メモリ K1 およびメモリ K2 にマッピングされた、カーネル 1 およびカーネル 2 の拡張によって使用される共有セグメント ( 610 ) を有する。セグメント 610 は、メモリ K1 ( 606 ) およびメモリ K2 ( 608 ) にマッピングされて上述の第 1 の形態によるパーティション間通信を可能にするメモリ ( 609 ) のセグメントについての定義および割り振りのテーブルを保持するのに使用されるか、または上記で図 5 を参照して説明した第 2 の通信形態に従ってメモリ A ( 602 ) およびメモリ B ( 604 ) にマッピングされたセグメント S2 ( 611 ) を定義するために使用される。本発明の一実施形態では、メモリ S は、制限されたサイズであり、実際のストレージ中で固定されている。ただし、メモリを固定する必要はなく、付随するページ管理タスクが効率的に管理されるのであれば、より大きな共有ストレージ空間を可能にすることを企図する。

#### 【 0053 】

ここで参照する本発明の第 1 の実施形態では、共有メモリ構成データ・セット ( S M C D S ) ( 613 ) からデータを読み出し、メモリ S ( 609 ) のセグメント S1 ( 610 ) にテーブルを構築する共有メモリ構成プログラム ( S M C P ) ( 612 ) と呼ばれるスタンドアロンのユーティリティ・プログラムにより、共有ストレージの定義および割り振りのテーブルをメモリ中にセットアップする。したがって、どのカーネルがどのストレージ・セグメントを共有するかについての割り振りと定義は、このユーティリティが作成する構成によって固定され、あらかじめ決定される。各種のカーネル拡張は次いでその共有ストレージを使用して、パイプ、メッセージ・キュー、ソケットなどの、各種のイメージ間、プロセス間の通信構造を実施し、さらには、各自の規定および規則に従っていくつかのセグメントを共有メモリ・セグメントとしてユーザ・プロセスに割り振る。これらのプロセス間の通信は、I P C A P I 618 および 619 を通じて可能になる。

#### 【 0054 】

共有ストレージの割り振りテーブルは、イメージ識別子、セグメント番号、g i d、u i d、「スティッキー・ビット」、および許可ビットからなるエントリを含む。スティッキー・ビットは、関連付けられたストアがページング可能でないことを示す。この例の実施形態では、スティッキー・ビットは確保され、1 と想定する ( すなわちデータはメモリ中のその場所に固定される、すなわち「固着されて」いる )。セグメントを使用する各グループ、ユーザ、およびイメージは、テーブルにエントリを有する。慣例的に、すべてのカーネルは、テーブルを読むことができるが、書き込むことはできない。初期化時に、カーネル拡張は、構成テーブルを読み取り、イメージを越えるプロセス間通信が他のプロセスから要求された際に使用する独自の割り振りテーブルを作成する。割り振られた空間の一部またはすべては、プロセス間通信を要求する他のプロセスからの要求によりカーネルが作成する「パイプ」、ファイル、およびメッセージ・キューを実装するためにカーネルが使用する。パイプとは、カーネル機能を通じて第 2 のプロセスに導かれる 1 プロセスか

10

20

30

40

50

らのデータである。パイプ、ファイル、およびメッセージ・キューは、Linux、OS/390、UNIX、および大半のUNIX(R)オペレーティング・システムで使用される、標準的なUNIX(R)オペレーティング・システムのプロセス間通信のAPIおよびデータ構造である。共有空間の一部は、システムを越えた直接的なメモリ共有のために、さらなるカーネル拡張により他のプロセスのアドレス空間にマッピングすることができる。

#### 【0055】

共有メモリの割り振り、使用、および仮想アドレス空間へのマッピングは、各自の規定および変換プロセスに従って各カーネルによって行われるが、基本的なハードウェア・ロックおよびメモリ共有のプロトコルは、システムの残りの部分の基礎となる共通のハードウェア設計アーキテクチャによって駆動される。

10

#### 【0056】

通信を行うために、上層レベルのプロトコルは共通でなければならない。好ましい実施形態では、これは、各種のオペレーティング・システム・イメージそれぞれに、UNIX(R)オペレーティング・システムとともに使用するIPC(プロセス間通信)APIを実装させることによって行われ、拡張がその要求をイメージを越えた要求であると識別する。この拡張は、パラメータか、または別個の新しい識別子/コマンド名によることができる。

#### 【0057】

図4および7を参照すると、本発明では、チャンネルまたはネットワーク接続を通じたデータ転送と、オペレーティング・システムに対する共有メモリの拡張の使用がともに回避されることが分かる。パーティション714のアプリケーション・プロセス(701)は、カーネル1(705)を呼び出すソケット・インタフェース718にアクセスする。ソケット・インタフェースとは、TCP/IPスタックの特定のポートを、リスン中のユーザ・プロセスに関連付ける構造である。カーネルはデバイス・ドライバ(716)にアクセスし、デバイス・ドライバ(716)は、I/Oアダプタ(720)のハードウェアを通じて、メモリ(401)にとってはメモリからメモリへの移動に見えるような形で、パーティション714および715のプロセッサ(402)またはファブリック(404)あるいはその両方に実装されるキャッシュ・メモリを回避して、カーネル・メモリ1(706)からカーネル・メモリ2(708)にデータを転送させる。データを移動すると、I/Oアダプタは次いでパーティション715のデバイス・ドライバ(717)にアクセスし、データを移動したことを通知する。デバイス・ドライバ717は次いで、カーネル2を待っているデータがソケット(719)にあることをカーネル2(707)に通知する。するとソケット(719)は、データをアプリケーション・プロセス(703)に提示する。このようにメモリからメモリへの直接の移動が行われ、同時に外部のインタフェースでデータを移動することが回避され、またメモリ共有のためのどちらかのオペレーティング・システムの拡張も回避される。

20

30

#### 【0058】

これに対し、図8に示す従来技術によるシステムでは、個別のメモリ移動操作を使用してカーネル・メモリ1(706)からアダプタ・メモリ・バッファ1(721)に移動を行う。第2のメモリ移動操作で、アダプタ・メモリ・バッファ1(721)からアダプタ・メモリ・バッファ2(722)にデータを移動する。そして第3のメモリ移動操作で、アダプタ・メモリ・バッファ2(722)からカーネル・メモリ2(708)にデータを移動する。これは、2つのカーネル・メモリ間でデータを移動するのに3回の別個のメモリ移動操作を使用することを意味し、これに対して図7の本発明では、1回のメモリ移動操作でカーネル・メモリ1(706)からカーネル・メモリ2(708)に直接データを移動する。これには、ユーザ・プロセスから見た待ち時間を低減する効果がある。

40

#### 【0059】

さらなる実施形態を図4および9によって例証する。ここでは、実際のデータ・ムーバ・ハードウェア(821)をファブリック(404)中に実施する。この実施形態の動作

50

は、I/Oアダプタ820内の管理状態(822)に従って、ファブリック(404)内のムーバ・ハードウェアによって実際にデータを移動する点を除いては上記の説明と同様に進行する。

【0060】

このようなファブリックに位置するデータ・ムーバの一例が、1993年12月7日発行のRobert D. Herzi他による「Processor System with Improved Memory Transfer Means」という名称の米国特許第5,269,009号に記載される。この文献は参照によりその全体を本明細書に含める。参照する同特許に記載される機構を、パーティションの主要な記憶位置間のデータ転送を含むように拡張する。

【0061】

本発明の実施形態は以下の要素を含む。CPUの設計によって定義される、基礎となる共通のデータ移動プロトコル、I/Oアダプタまたはファブリック・ハードウェアあるいはその両方、I/Oアダプタとのインタフェースを実装するデバイス・ドライバの異種のセット、好ましい実施形態ではソケット・インタフェースとして示す共通の高レベルのネットワーク・プロトコル、物理メモリ・アドレスへのネットワーク・アドレスのマッピング、およびI/Oアダプタ(820)が各パーティションのカーネル・メモリおよびデバイス・ドライバと通信するのに使用するI/O割り込みベクトルまたはポインタ。

【0062】

データ・ムーバは、ハードウェアの状態機械として、またはマイクロコードおよびマイクロプロセッサを用いてI/Oアダプタ中に実装することができる。あるいは、I/Oアダプタによって制御される、マシンの通信ファブリック中のデータ・ムーバを使用するように実装してもよい。このようなデータ・ムーバの一例が、1993年12月7日発行のHerzi他による米国特許第5,269,009号「PROCESSOR SYSTEM WITH IMPROVED MEMORYTRANSFER MEANS」に記載される。

【0063】

図10を参照すると、実装に関係なくデータ・ムーバは以下の要素を有する。メモリからのデータはソース・レジスタ(901)に保持され、データは、データ・アライナ(902および904)を通じてデスティネーション・レジスタ(903)に渡され、次いでメモリに戻される。したがって、連続した動作の一部としてメモリからの取り出しそしてメモリへの格納がある。すなわち、位置合わせプロセスは、メモリ・ラインから複数のワードが取り出されると行われる。位置合わせしたデータは、メモリ・ストアが開始されるまで、デスティネーション・レジスタ(903)でバッファリングする。ソース・レジスタ(901)およびデスティネーション・レジスタ(903)は、移動操作中に取り出しと格納の間でどれほどの重複が許されるかに応じて、単一ラインまたは複数ライン分のメモリ・データを保持するために使用することができる。メモリのアドレス指定は、移動中に取り出しと格納が行われるアドレスを追跡するカウンタ(905および906)から行われる。制御およびバイト数の要素(908)は、アライナ(902および904)を通じたデータ・フローを制御し、メモリ・アドレスに合わせてソース・カウンタ(905)または宛先カウンタ(906)を選択させる。コントローラ(908)は、アドレス・カウンタ(905および906)の更新も制御する。

【0064】

図11を参照すると、データ・ムーバは、デバイス・ドライバによって実装される、特権CISC命令(1000)として実装することもできる。MVXLは、カウント・レジスタによって指定されるバイト数を、ソース・レジスタで指定される物理アドレスから、デスティネーション・レジスタで指定される物理アドレスに移動する。この命令は特権命令である。(仮想アドレス間ではMVC Lが同じ機能を行う。)ここで、デバイス・ドライバは、仮想アドレスではなく物理アドレスをレジスタにロードし、パーティションを越えたデータ移動を可能にする。このようなCISC命令は、S/390のMove Page、Move Character Longなどのパーティション間のデータ移動に適したハードウェア機能を利用するが、ネットワーク・アドレスおよびオフセットを物理メモリのアドレスにマッピング

10

20

30

40

50

するテーブルに従って物理的にメモリをアドレス指定する特権も有する。最後に、データ・ムーバおよびアダプタは、仮想アダプタとして機能するハイパーバイザ・コードによって実装することができる。

#### 【 0 0 6 5 】

図 1 2 に、次のステップからなる、データ・ムーバがアダプタにある場合のデータ・ムーバの動作を示す。

1 1 0 1 ユーザがデバイス・ドライバを呼び出し、ソース・ネットワークの ID、ソース・オフセット、宛先ネットワークの ID を提供する。

1 1 0 2 デバイス・ドライバがアドレスをアダプタに転送する。

1 1 0 3 アダプタがアドレスを転送する。

ID から物理的な基底アドレスを検索する ( テーブル・ルックアップ )  
ロックおよび現在の宛先オフセットを得る

オフセットを加算する

境界を確認する

1 1 0 4 アダプタがカウントとアドレスをレジスタにロードする

1 1 0 5 アダプタがデータ移動を実行する

1 1 0 6 アダプタがロックを解除する

1 1 0 7 アダプタがユーザに「戻る」デバイス・ドライバに通知する

10

#### 【 0 0 6 6 】

図 1 3 に、以下のメソッドを含む、プロセッサ通信ファブリックで実装されるデータ・ムーバ・メソッドを示す。

1 2 0 1 ユーザがデバイス・ドライバを呼び出し、ソース・ネットワークの ID、ソース・オフセット、宛先ネットワークの ID を提供する

1 2 0 2 デバイス・ドライバがアドレスをアダプタに送る

1 2 0 3 アダプタがアドレスを変換する

ID から物理的な基底アドレスを検索する ( テーブル・ルックアップ )  
ロックおよび現在の宛先オフセットを得る

オフセットを加算する

境界を確認する

アダプタが、ロックおよび物理アドレスをデバイス・ドライバに戻す

1 2 0 4 デバイス・ドライバがデータ移動を実行する

1 2 0 5 デバイス・ドライバがロックを解除する

1 2 0 6 デバイス・ドライバが戻る

20

30

#### 【 0 0 6 7 】

以上、区分されたコンピューティング・システム内で異種の要素を含む相互動作を実施する 2 つの方式を説明した。方式の 1 つでは、共有メモリ機能とオペレーティング・システム・カーネルの拡張を使用してパーティションを越えたプロセス間の通信プロトコルを可能にし、もう一方の方式では、共有 I / O アダプタの能力を使用してすべての物理メモリをアドレス指定して、メモリからメモリへのメッセージ・パッシングを 1 回の操作で実施する。

40

#### 【 0 0 6 8 】

前述の構造により、単一システムのクライアント - サーバ・モデルを活用するいくつかの本発明の実装が得られる。この構造を実装する一方式は、サーバの作業キューを共有される記憶空間に入れて、各種のクライアントが要求を追加できるようにするものである。そして、クライアントがそこに入れられた情報にアクセスできるように、「リモート」クライアントのためのリターン・バッファも共有メモリ空間になければならない。あるいは、上述のメッセージ・パッシング方式を使用して、既存のネットワーク指向のクライアント / サーバを迅速かつ容易に展開することができる。これらの実装は例証として提供しており、限定的と見なすべきではない。実際、当業者は、この構造を各種の方式でさらに発展させ、単一システムの枠組みの中で様々なタイプの異種のクライアント - サーバ・シス

50

テムを実装できることが容易に理解される。

【0069】

パーティション・クラスタのワークロード管理

本発明のパーティション・リソース・マネージャを用いたパーティションのワークロード管理について、OS/390の実装により説明する。パーティション・リソース・マネージャのタスクは、あるパーティションで実行されるワークロード・マネージャの連携からなり、ハイパーバイザと通信する別のパーティションの他のワークロード・マネージャ（あるいは特定のオペレーティング・システムに適したスループット情報ジェネレータ）からスループット情報を得、パーティション・リソースの割り振りを調整する役割を担う。当業者には、本発明は、システム・アーキテクチャの必要性に従ってリソース・マネージャ機能を実装することにより、他のシステム・アーキテクチャでも実装できることは容易に理解されよう。

10

【0070】

図14を参照すると、OS/390オペレーティング・システムのワークロード・マネージャ(WLM)(1308)は、S/390のパーティション・ハイパーバイザと通信して、各パーティションに割り振られたリソースを調整する機能を持つ。これは、LPARクラスタリングとして知られる。ただし、OS/390でないパーティション(1301)については、WLMは、そのパーティションのオペレーティング・システムまたはアプリケーションに基づくのではなく、ハイパーバイザから供給できる利用率およびその他の情報のみに基づいて割り振りを行わなければならない。上述の待ち時間が低いパーティションを越えた通信(1305)を使用してパーティションからWLM(1308)に情報をパイピングすることは、システムを越えてリソースを割り振るジョブをより適切に行うためにWLMが必要とする情報をWLM(1308)に与える、オーバーヘッドが非常に低い手段となる。これは、アプリケーションがワークロード管理のために装備されない場合であっても効果的でありうる。この理由は、通例、制御されるシステムは、システムで送受信されるIPパケットを数えるTCP/IPスタック中のパケット・アクティビティ・カウンタにアクセスするコマンドであるUNIX(R)オペレーティング・システムの「NETSTAT」(UNIX(R)オペレーティング・システムの標準的なコマンド・ライブラリの一部)を実装することができ、また、ビジー・サイクルとアイドル・サイクルを数え、利用率データ(1302)を生成するカーネル中のシステム・アクティビティ・カウンタにアクセスする標準的なUNIX(R)オペレーティング・システム・コマンドであるUNIX(R)オペレーティング・システムの「VMSTAT」(UNIX(R)オペレーティング・システムの標準コマンド・ライブラリの一部)を実行することができるからである。既存のNETSTATコマンドおよびVMSTATコマンドを使用することは必須ではなく、パケット数と利用率をそれらに供給する基礎的な機構を使用してリソースおよび経路長のコストを最小に抑えることが最良であることは理解されよう。このデータを組み合わせると「速度」の基準(1303)とし、それをワークロード・マネージャ(WLM)のパーティション(1307)に送ることにより、WLM(1308)は、ハイパーバイザにリソースの調整を行わせることができる。CPUの利用率が高く、パケット・トラフィックが低い場合、パーティションはより多くのリソースを必要とする。接続(1304および1306)は、相互接続(1305)の実施形態に応じて異なる。共有メモリの実施形態では、それらの接続はUNIX(R)オペレーティング・システムのPIPE、Message Q、SHMEM、あるいはソケット構造である。データ・ムーバの実施形態では、通例ソケット接続である。

20

30

40

【0071】

一実施形態では、「速度」基準は以下のように得る(IBM Redbook Document SG24-4810-01「Understanding RS/6000 Performance and Sizing」に記載されるUNIX(R)オペレーティング・システム・コマンドNETSTATおよびVMSTATを参照されたい)。

(NETSTAT)の全パケットについての時間間隔データを使用してスループットの

50

概略を得る。

時間間隔についてのCPUデータ (VMSTAT) を使用して、CPUの利用率の概略を得る。

ピークを1としてトラフィックを正規化してそれらの概略データをプロットし、表示する (1401)。

トラフィック対CPUについて累積的な相関分析を行う (1402)。

トラフィックの関係は、関数  $T(C)$  に沿った曲線になる。

ここで使用する例 (1402) では、 $T(C) = y(x) = 0.802 + 1.13x$  となる。

$S = dT/dC = dy/dx$  が速度基準であり、この例では  $S = 1.13$  となる。

$S$  が傾向線よりも小さい場合には、より多くのリソースが必要とされる。

10

#### 【0072】

図15の例ではこれは2度発生している (1403および1404)。管理図は、諸業界で監視プロセスを作成するための標準的な方法である。1405では  $S$  を管理図として動的にプロットしている。この例でパケット・トラフィックとCPUとの間に見られるような関係を与えられると、統計的な制御理論に基づいて収集データを様々な形で監視し、調整することができる。これらの方法は、通例、アクションをトリガする制御変数の閾値に依存する。あらゆるフィードバック・システムでそうであるように、制御不能状態に近いと判断された時に即座にアクションを行わせることが必要である。さもないとシステムが不安定になる可能性がある。この実施形態では、これは、内部通信が提供する低待ち時間の接続によって行われる。

20

#### 【0073】

静的な環境では、 $S$  を使用して、どの利用率でより多くのリソースが必要とされるのかを明確にすることができる。これは平均にわたって有効であるが、 $S$  はワークロードと時間の関数でもある。図15を参照すると、まずそれは50~60%の間のどこかであることが分かり、次いで  $S$  の谷があるとその後少なくとも1回の時間間隔で利用率のピークがあることが分かる。したがって、WLMは利用率ではなく  $S$  を与えられた場合により好適なジョブを行う。これは、 $S$  が、よりタイムリーなリソース調整を可能にする「先行する指標」であるためである。区分されたマシンのリソースはパーティションによって共有されるので、ワークロード・マネージャは  $S$  のデータを複数のパーティションから得なければならない。データの転送は、非常に低いオーバーヘッド、かつ高い速度で行う必要がある。この実施形態はこれら両方の条件を可能にする。図14を参照すると、ワークロード・マネージャがないパーティション (1301) では、プログラム・ステップ (1303) でパラメータ (この例では「 $S$ 」) を評価するために使用する利用率およびパケットのデータ (1302) をモニタが収集する。プログラムは次いで、低待ち時間のパーティション間通信機能 (1305) への接続 (1304) を使用し、パーティション間通信機能 (1305) は、ワークロード・マネージャを備えるパーティション (1307) の接続 (1306) にそのデータを渡し、その接続により「論理パーティション・クラスタ・マネージャ」 (1308) に入力が提供される。「論理パーティション・クラスタ・マネージャ」については、2000年10月2日出願の米国特許出願番号09/677338「METHOD AND APPARATUS FOR ENFORCING CAPACITY LIMITATIONS IN A LOGICALLY PARTITIONED SYSTEM」に記載される。

30

40

#### 【0074】

この場合、パーティション・データをワークロード・マネージャに通信する最も効率的な方法は、メモリ共有を通じて行うものであるが、ソケットの待ち時間が十分に低くデータの時間伝達が可能である場合には、内部のソケット接続も機能する。これは、ワークロードと、求められる制御の精度の両方に応じて決まる。

#### 【0075】

上述は、ワークロード・マネージャがリソースを割り振るための情報を供給する1方式であるが、これは決して限定的なものと解釈すべきではない。この例を選択したのは、これ

50

が、大量の新たなコードを用いずに、すべてではなくとも大半のオペレーティング・システムから収集することが可能な基準であるためである。クライアント・システムは、応答時間やユーザ数など、WLMサーバに渡す任意の基準の任意の手段を実装することができる。

#### 【0076】

##### 非間接的I/O

時に、デバイス・ドライバは、ハードウェアによってサポートされる可能なオペレーティング・システムのうち1つのみでしか利用できない場合がある。共有メモリにデバイス・ドライバのメモリ・インタフェースを提示し、すべての付属システムがドライバ・プロトコルに準拠することにより、デバイスを複数のシステムで共有することができる。実際には、1つのパーティションが他のパーティションのIOPになることができる。デバイスに過負荷を課すと、単一システムからデバイスに過負荷を課す場合と同様の否定的な結果になることを理解した上で、デバイスへのアクセスは単一システムのレベルに近いものになる。図16を参照すると、デバイス・ドライバ(1501)は、共有メモリ(1511)を通じた、アプリケーションおよびアクセス・メソッド(1503)からのI/Oサービスに対する要求に応答する。

10

#### 【0077】

このメッセージ・パッシングの実施形態をいくつかのデバイスに使用することは可能であるが、ソケット、スタックおよびデータ移動の待ち時間を許容しなければならない。これは、ネイティブのデバイスとネットワークに接続されたデバイスとの中間と考えることができる。

20

#### 【0078】

デバイス・ドライバを実行するシステム・イメージに割り振られたプロセッサ・リソースを、アプリケーションを実行するシステム・イメージに割り振られたプロセッサ・リソースから分離した場合にはさらなる機能の増大が得られる。この分離を行うと、I/O割り込みの対象とならないプロセッサでは、I/O割り込みとそれに伴うコンテキストの切り替えによるキャッシュおよびプログラム・フローの中断が回避される。

#### 【0079】

##### 共通のセキュリティ・サーバ

アプリケーションがウェブ対応型になり、統合されるのに従い、ユーザの確認と権利付与の確認は、従来のシステムの場合と比べて一般的な課題になる。この課題をさらに増大するのは、アプリケーションを統合するために異種のシステムをまとめる必要性である。この結果、LDAP、Kerberos、RACF、およびその他のセキュリティ機能を統合された形で使用するには、通常、セキュリティ機能を行う共通のセキュリティ・サーバへのネットワーク接続が必要とされる。これは、パフォーマンスに影響を与える。また、ネットワーク盗聴者(sniffer)によるセキュリティの露出もある。共有メモリ接続またはメモリ・ムバ接続を介して共通のセキュリティ・サーバをウェブ・サーバに接続した場合、このアクティビティは著しく速度が上がり、接続は内部化され、セキュリティが向上する。さらに、そのような環境では、一部の顧客は、特にLinuxの場合には、他のUNIX(R)オペレーティング・システムに基づくパスワード保護よりも、S/390「RACF」、あるいは他のOS/390「SAF」インタフェースのユーザ認証によるより高いセキュリティを選ぶ可能性がある。Linuxシステムでは、適合し、カスタマイズすることが意図される「プラグ可能認証モジュール」によってそこでユーザ認証を行うので、そのような共有サーバに対するクライアント・サイドを比較的容易に構築することができる。ここでは、共有メモリ・インタフェースまたはメモリからメモリへのデータ・ムバ・インタフェースを介してセキュリティ・サーバにアクセスし、ウェブ・サーバはそのアクセスを求めて競争する。その結果生じる作業キューはセキュリティ・サーバによって実行し、必要に応じて共有メモリ・インタフェースを通じて応答を返す。この結果、ウェブ・アプリケーションに対する強化されたセキュリティとパフォーマンスが発揮される。図17を参照すると、セキュリティ・サーバ(1601)は、共有メモリ(16

30

40

50

11)を通じたユーザ・プロセス(1603)からのアクセス要求に応答する。ユーザ・プロセスは、カーネル2(1607)のセキュリティ・クライアント・プロセス(LINUXの場合にはこれがPAMとなる)との標準的なプロセス間通信(IPC)インタフェースを使用し、セキュリティ・クライアント・プロセスは、共有メモリ(1610)を通じてカーネル1(1605)のカーネル・プロセスと通信し、カーネル・プロセスは、セキュリティ・サーバ・インタフェース(OS/390またはZ/OSの場合はSAF)をユーザ・プロセス(1603)のプロキシとして駆動し、共有メモリ(1610)を通じてカーネル2(1607)のセキュリティ・クライアントに認証を返す。

【0080】

別の実施形態では、共有メモリに入れられたデータは、単一操作のデータ・ムーバを介してカーネル・メモリ1(1606)からカーネル・メモリ2(1608)に移動し、共有メモリの開発を回避するとともに、ネットワーク接続も回避する。

10

【0081】

以下で、共通のセキュリティ・サーバ(1601)を第1のパーティション(1614)で実行し、少なくとも1つのセキュリティ・クライアント(またはプロキシ)(1603)を少なくとも1つの第2のパーティション(1615)で実行する、区分された処理システムにセキュリティを提供するセキュリティ・サーバにおける通信ステップの実装の一例を挙げる。

【0082】

ユーザが認証を要求する。セキュリティ・クライアント(1603)がユーザからパスワードを受け取る。セキュリティ・クライアントは、セキュリティ・サーバ(1610)からのアクセスが可能なメモリ位置に要求を入れ、要求を入れたことを通知する。第1のパーティション(1614)の「セキュリティ・デーモン」が通知を認識し、第1のパーティション(1614)の「プロキシ」クライアント(1616)を開始する。プロキシ(1616)クライアントは、セキュリティ・サーバ(1601)にとってネイティブのインタフェースを使用して、その要求をもってセキュリティ・サーバを呼び出す。セキュリティ・サーバ(1601)は、要求を処理し、プロキシ・クライアント(1616)にサーバ応答を返す。プロキシ・クライアントは、セキュリティ・サーバの応答を第2のパーティションのセキュリティ・クライアントによるアクセスが可能なメモリに入れ、要求を入れたことを通知する。この通知により、その権限付与をポイントするセキュリティ・クライアント(1603)が起動する。セキュリティ・クライアント(1603)は、ユーザに応答を戻す。一実施形態では、第2のパーティション(1615)のセキュリティ・クライアント(1603)は、共有メモリ・インタフェース(1609)を用いて第1のパーティション(1614)のセキュリティ・サーバ(1601)と通信し、したがってネットワーク接続のセキュリティの露出を防止し、パフォーマンスを高める。別の実施形態では、第2のパーティションのセキュリティ・クライアントは、図9に示すデータ・ムーバ(821)を使用した内部のメモリからメモリへの移動により、第1のパーティションのセキュリティ・サーバと通信する。図9を参照すると、この第2の実施形態は、セキュリティ・クライアントをプロセスA(803)として実装し、セキュリティ・プロキシはプロセスB(801)として実装され、したがって外部のネットワーク接続が回避され、共有メモリの実装が回避される。

20

30

40

【図面の簡単な説明】

【0083】

【図1】区分データ処理システムの概要の図である。

【図2】1つまたは複数のシステム・ボードからなるパーティションを有する、物理的に区分された処理システムの図である。

【図3】論理的に区分されたリソースが各自の個々のパーティションのために確保された、論理的に区分された処理システムの図である。

【図4】論理的に区分されたリソースをいくつかのパーティション間で動的に共有することが可能な、論理的に区分された処理システムの図である。

50

【図5】UNIX(R)オペレーティング・システムの「プロセス間通信」の構造の図である。

【図6】スタンドアロン・ユーティリティによってロードされた構成テーブルに従って実際のメモリが共有される実施形態の図である。

【図7】I/Oアダプタおよびそのドライバの機能を使用してパーティション間のデータ転送を容易にする実施形態の図である。

【図8】従来技術によるシステムの図である。

【図9】区分データ処理システムの通信ファブリック中に実装されたデータ・ムーバによってパーティション間の実際のデータ転送を行う実施形態の図である。

【図10】データ・ムーバ例の構成要素の図である。

10

【図11】IBM S/390移動命令の例示的フォーマットの図である。

【図12】アダプタ・データ・ムーブを行う例示的ステップの図である。

【図13】プロセッサ・データ・ムーブを行う例示的ステップの図である。

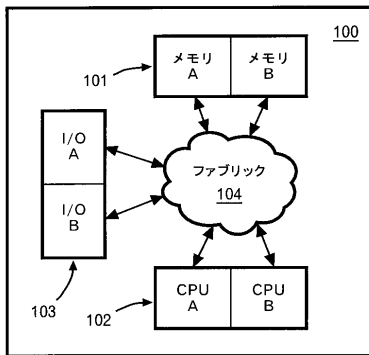
【図14】ワークロード・マネージャ(WLM)の高レベル図である。

【図15】典型的なワークロード管理データを示す図である。

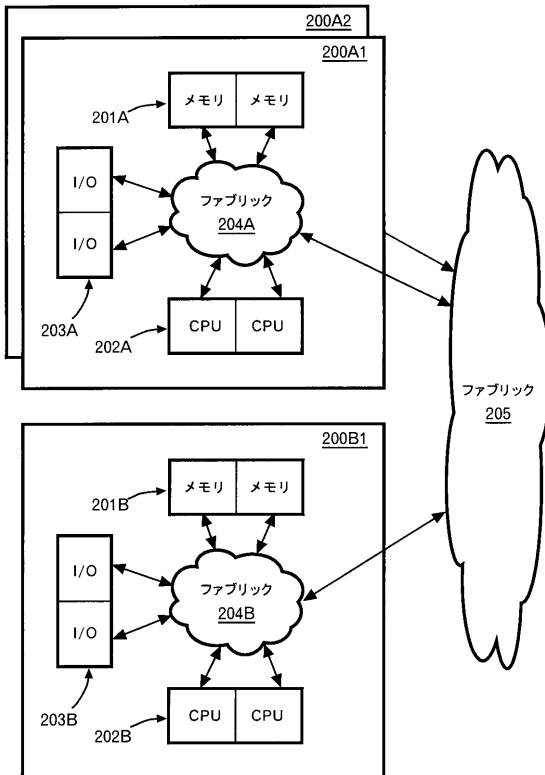
【図16】間接的なI/Oを使用したクライアント/サーバのクラスタ化の図である。

【図17】クライアント/サーバのサーバ・クラスタ化の図である。

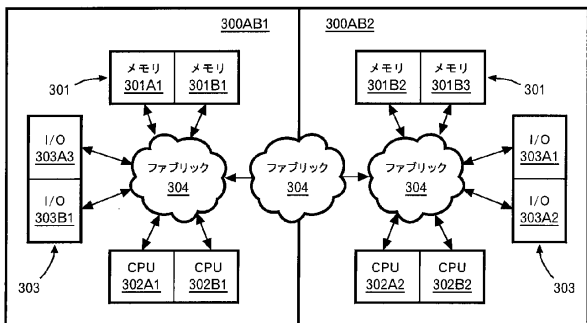
【図1】



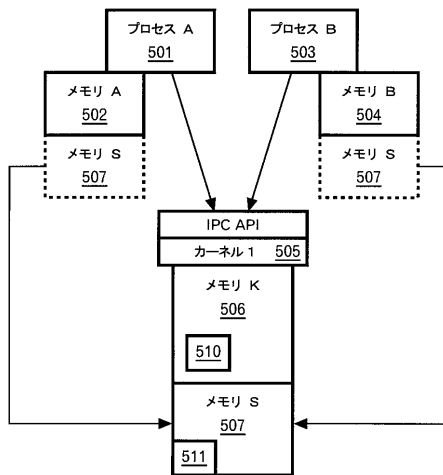
【図2】



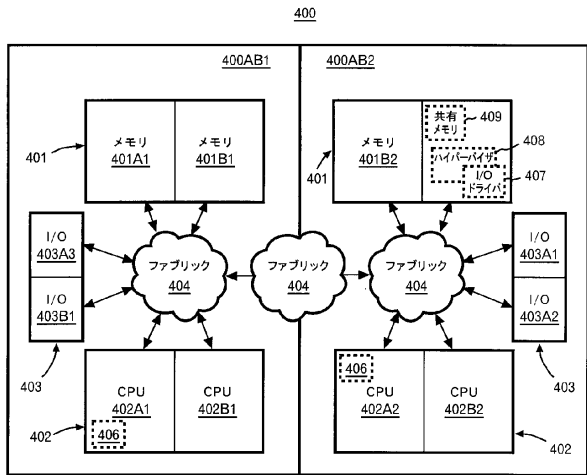
【 図 3 】



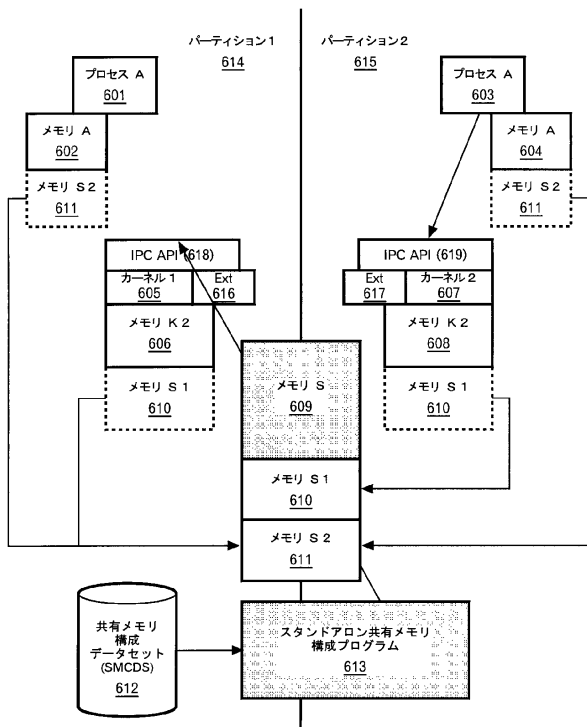
【 図 5 】



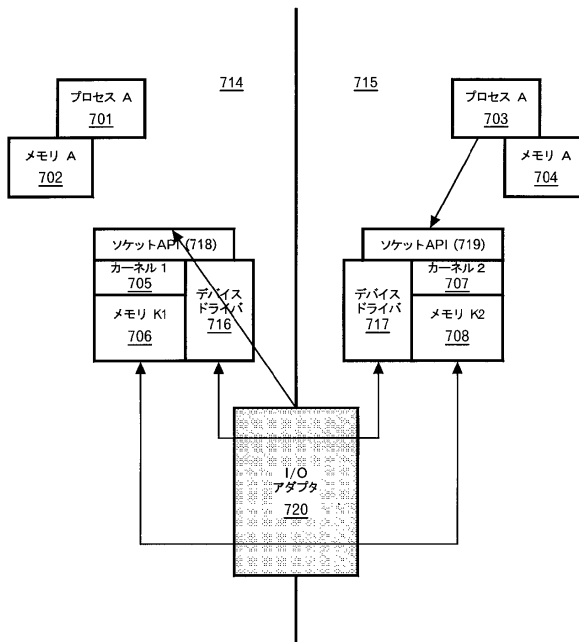
【 図 4 】



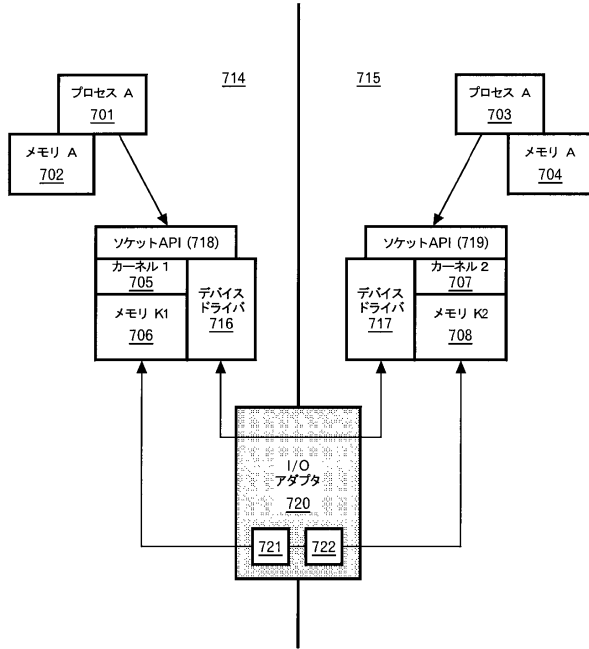
【 図 6 】



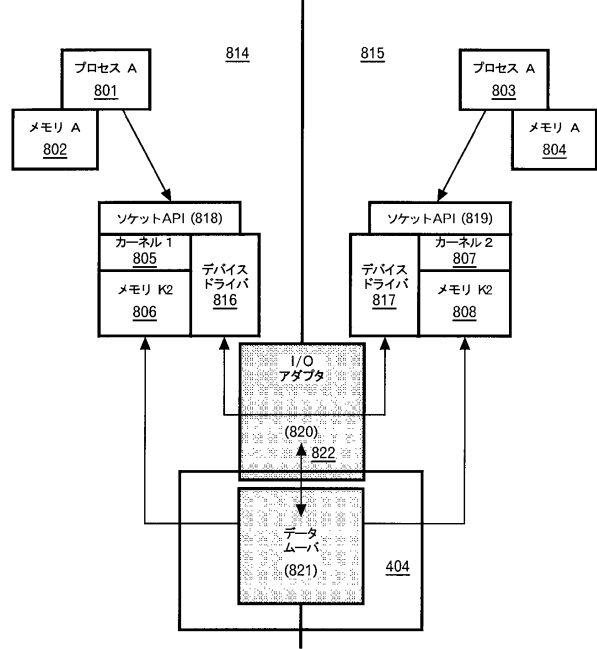
【 図 7 】



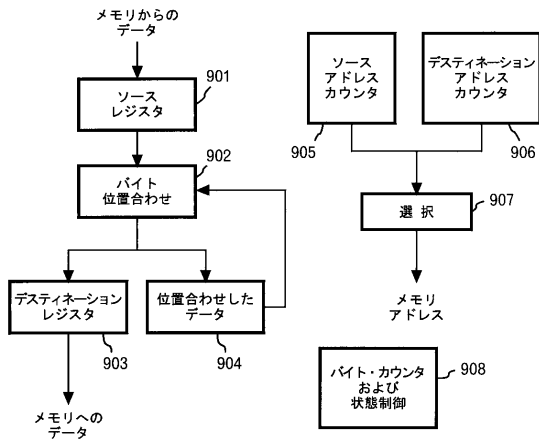
【 図 8 】



【 図 9 】



【 図 10 】



【 図 11 】

1000

MVXL	バイト数	デスティネーション	ソース

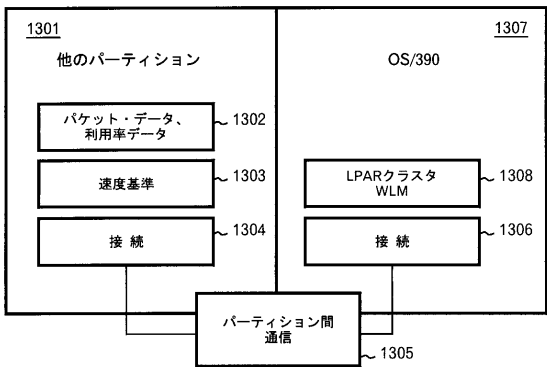
【 図 12 】

1101. ユーザがデバイス・ドライバを呼び出し、  
 - ソース・ネットワークID、  
 - ソース・オフセット、  
 宛先ネットワークのIDを提供する
1102. デバイス・ドライバがアドレスをアダプタに転送する
1103. アダプタがアドレスを変換する  
 - IDから物理的な規定アドレスを検索する(テーブル・ルックアップ)  
 - ロックと現在の宛先オフセットを得る  
 - オフセットを加算する  
 - 境界を確認する
1104. アダプタがレジスタにカウントとアドレスをロードする
1105. アダプタがデータ移動を実行する
1106. アダプタがロックを解除する
1107. アダプタが、ユーザに「戻る」デバイス・ドライバを通知する

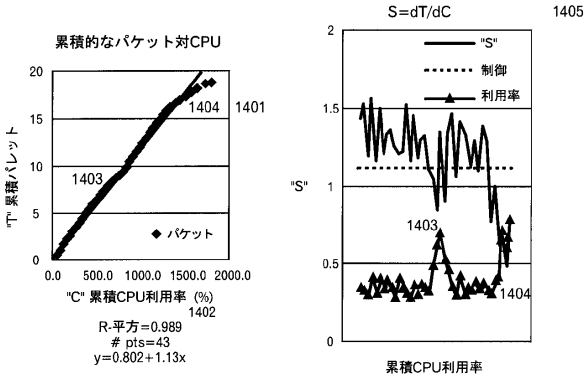
【 図 13 】

1201. ユーザがデバイス・ドライバを呼び出し、  
 - ソース・ネットワークID、  
 - ソース・オフセット、  
 宛先ネットワークのIDを提供する
1202. デバイス・ドライバがアドレスをアダプタに送る
1203. アダプタが変換を行う  
 - IDから物理的な規定アドレスを検索する(テーブル・ルックアップ)  
 - ロックと現在の宛先オフセットを得る  
 - オフセットを加算する  
 - 境界を確認する  
 - ロックと物理アドレスをデバイス・ドライバに戻す
1204. デバイス・ドライバがデータ移動を実行する
1205. デバイス・ドライバがロックを解除する
1206. デバイス・ドライバが戻る

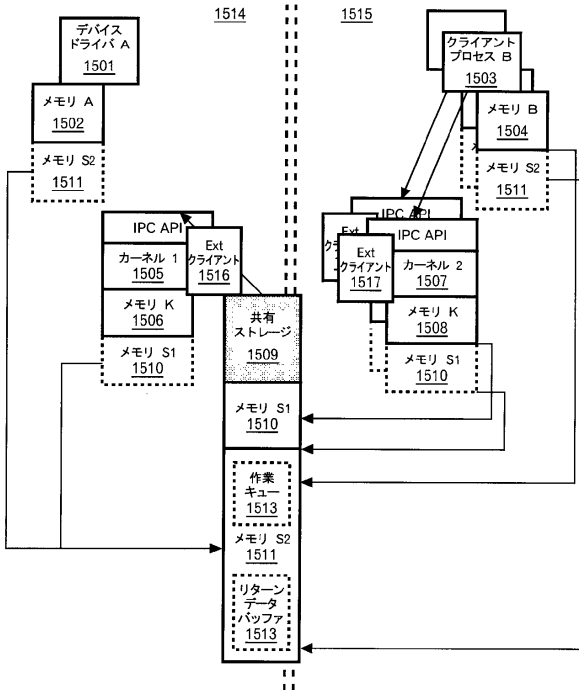
【 図 1 4 】



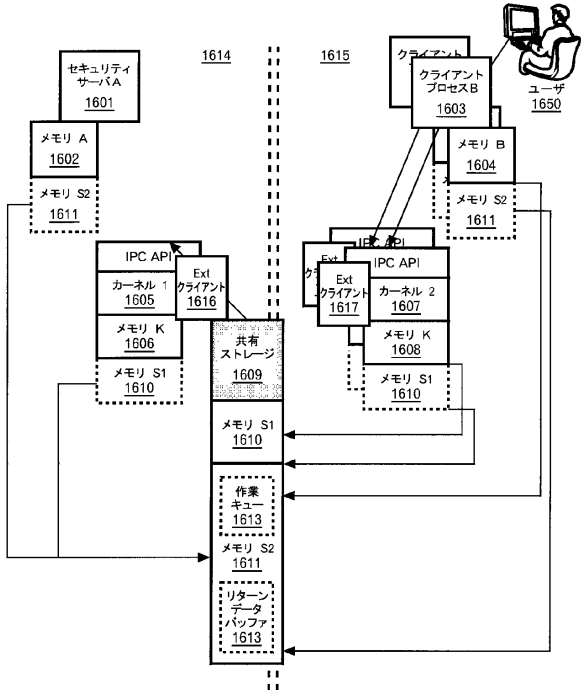
【 図 1 5 】



【 図 1 6 】



【 図 1 7 】



## フロントページの続き

- (72)発明者 クバラ、ジェフリー  
アメリカ合衆国12570 ニューヨーク州プーカグ モーガン・レーン 10
- (72)発明者 ニック、ジェフリー  
アメリカ合衆国12493 ニューヨーク州ウエスト・パーク ルート 9 ダブリュー 195  
7
- (72)発明者 テンプル、ジョゼフ、サード  
アメリカ合衆国12443 ニューヨーク州ハーレー フック・ストリート 312 私書箱50  
7
- (72)発明者 ヨコム、ピーター  
アメリカ合衆国12540 ニューヨーク州ラグランジェビル ユニット 29 ストリングム・  
ロード 129

審査官 殿川 雅也

- (56)参考文献 特開2000-132530(JP,A)  
特表2002-532806(JP,A)  
特開平10-334057(JP,A)  
特開平10-027167(JP,A)  
特開平07-295841(JP,A)  
特開昭64-002145(JP,A)  
国際公開第00/036509(WO,A1)  
特許第2638065(JP,B2)  
Kinshuk Govil et al., Cellular disco: resource management using virtual clusters on shared-memory multiprocessors, ACM Transactions on Computer Systems (TOCS), 2000年  
8月, Volume 18, Issue 3, pp.229-262
- (58)調査した分野(Int.Cl., DB名)  
G06F 9/46 - 9/54