

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4288295号
(P4288295)

(45) 発行日 平成21年7月1日(2009.7.1)

(24) 登録日 平成21年4月3日(2009.4.3)

(51) Int. Cl. F I
G 0 6 F 9/50 (2006.01) G O 6 F 9/46 4 6 2 Z
G 0 6 F 9/46 (2006.01) G O 6 F 9/46 3 5 0

請求項の数 11 (全 19 頁)

(21) 出願番号	特願2007-513913 (P2007-513913)	(73) 特許権者	390009531
(86) (22) 出願日	平成17年5月18日 (2005.5.18)		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2008-500619 (P2008-500619A)		I N T E R N A T I O N A L B U S I N E S S M A S C H I N E S C O R P O R A T I O N
(43) 公表日	平成20年1月10日 (2008.1.10)		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/EP2005/052279		
(87) 国際公開番号	W02005/119444	(74) 代理人	100108501
(87) 国際公開日	平成17年12月15日 (2005.12.15)		弁理士 上野 剛史
審査請求日	平成20年4月9日 (2008.4.9)	(74) 代理人	100112690
(31) 優先権主張番号	10/855,726		弁理士 太佐 種一
(32) 優先日	平成16年5月27日 (2004.5.27)	(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】 論理パーティションにおけるメモリへのアクセス

(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、第1の論理パーティションにおける第1のプロセスが第2の論理パーティションにおける遠隔メモリにアクセスするための方法であって、

前記遠隔メモリのための、前記遠隔メモリの記述を与える拡張メモリ間記述子を検索するステップであって、前記拡張メモリ間記述子が、異なる論理パーティションにおける遠隔メモリを記述するものとして前記拡張メモリ間記述子を示す第1のフィールドを含む、ステップと、

前記拡張メモリ間記述子に基づいて、前記第1の論理パーティションにおける前記第1のプロセスが、前記第2の論理パーティションにおける前記遠隔メモリにアクセスするステップと、

を含む、方法。

【請求項2】

前記拡張メモリ間記述子が、前記第2の論理パーティションにおける前記遠隔メモリを識別する第2のフィールドと、前記遠隔メモリのサイズを示す第3のフィールドと、前記遠隔メモリにおける開始アドレスを示す第4のフィールドと、を更に含む、請求項1に記載の方法。

【請求項3】

前記拡張メモリ間記述子が、前記第1の論理パーティションにおけるオペレーティング・システムが初期化されるときに発生される、請求項1に記載の方法。

10

20

【請求項 4】

前記拡張メモリ間記述子が、他の論理パーティションの遠隔メモリを前記第 1 の論理パーティションのオペレーティング・システムにアタッチするオペレーティング・システム・カーネル・サービスによって発生される、請求項 1 に記載の方法。

【請求項 5】

前記オペレーティング・システム・カーネル・サービスが、前記第 1 の論理パーティションにおける前記オペレーティング・システムが初期化されるときに前記オペレーティング・システム・カーネル・サービスによってパースされる前記第 2 の論理パーティションのためのデバイス・ツリーに基づいて前記拡張メモリ間記述子を発生する、請求項 4 に記載の方法。

10

【請求項 6】

前記デバイス・ツリーがハイパーバイザによって維持される、請求項 5 に記載の方法。

【請求項 7】

前記第 1 の論理パーティションがサーバ・コンピューティング・デバイスにあり、前記第 2 の論理パーティションがクライアント・コンピューティング・デバイスにある、請求項 1 に記載の方法。

【請求項 8】

前記第 2 の論理パーティションにおける前記遠隔メモリにアクセスするステップが、前記拡張メモリ間記述子を用いて直接メモリ・アクセス動作を実行することを含む、請求項 1 に記載の方法。

20

【請求項 9】

前記拡張メモリ間記述子に基づいて前記第 2 の論理パーティションにおける前記遠隔メモリにアクセスするステップが、

前記拡張メモリ間記述子を入出力 (I / O) スタックを通して物理デバイス・ドライバに渡すステップと、

前記拡張メモリ間記述子に基づいて直接メモリ・アクセス動作を発生するステップと、

前記直接メモリ・アクセス動作を I / O アダプタに提示するステップと、

を含み、前記 I / O アダプタは前記直接メモリ・アクセス動作を前記第 2 の論理パーティションに送信する、請求項 1 に記載の方法。

【請求項 10】

30

第 1 の論理パーティションにおける第 1 のプロセスが第 2 の論理パーティションにおける遠隔メモリにアクセスするためのプログラムであって、

前記遠隔メモリのための、前記遠隔メモリの記述を与える拡張メモリ間記述子を検索するステップであって、前記拡張メモリ間記述子が、異なる論理パーティションにおける遠隔メモリを記述するものとして前記拡張メモリ間記述子を示す第 1 のフィールドを含む、ステップと、

前記拡張メモリ間記述子に基づいて、前記第 1 の論理パーティションにおける前記第 1 のプロセスが、前記第 2 の論理パーティションにおける前記遠隔メモリにアクセスするステップと、

をコンピュータに実行させるプログラム。

40

【請求項 11】

第 1 の論理パーティションにおける第 1 のプロセスが第 2 の論理パーティションにおける遠隔メモリにアクセスするためのシステムであって、

前記遠隔メモリのための、前記遠隔メモリの記述を与える拡張メモリ間記述子を検索するための手段であって、前記拡張メモリ間記述子が、異なる論理パーティションにおける遠隔メモリを記述するものとして前記拡張メモリ間記述子を示す第 1 のフィールドを含む、手段と、

前記拡張メモリ間記述子に基づいて、前記第 1 の論理パーティションにおける前記第 1 のプロセスが、前記第 2 の論理パーティションにおける前記遠隔メモリにアクセスするための手段と、

50

を含む、システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に、改良したデータ処理システムを対象とする。更に具体的には、本発明は、メモリ間記述子を拡張して、論理的にパーティションに分割したコンピューティング環境において、別のパーティションからメモリを記述するために使えるようにするためのシステムおよび方法を対象とする。

【背景技術】

【0002】

AIX(R)オペレーティング・システムは、IBM社のpSeries(R)マシンにおいて、論理分割(LPAR)をサポートする。LPARは、単一プロセッサ・コンプレックス上での多数のシステム・イメージの実行を可能とする。かかる各イメージは、十分なCPU(専用または共用)、中央ストレージ、拡張ストレージおよびチャネルを有する。LPARでは、パーティション間でリソースは明確に分離され、あるパーティションは他のいかなるパーティションのシステム環境も不安定にすることができないようになっている。IBM社のマシンでは、パーティションのためのリソースのこの明らかな分離を可能とする機構は、ハイパーバイザである。

【0003】

各パーティションはそれ自身の分離したリソースの組を有するので、各パーティションは、周囲のパーティションから隠れたそれ自身の「ローカル」メモリを有する。仮想入出力(VIO)を生成するまで、あるパーティションが別のパーティションのメモリに直接アクセスする正当な理由はなかった。

【0004】

VIOによって、ポイント・ツー・ポイントの直接メモリ・アクセス(DMA)動作およびパーティション間のメモリ・コピー動作の概念が導入された。DMA動作によって、あるメモリから別のメモリへの直接データ転送が可能となり、このデータ転送を実行するために中央演算処理装置(CPU)を用いずに済む。DMA動作では、データ転送にCPUを用いなければならない場合に比べて高速のデータ転送が可能となる。メモリ・コピー動作によって、メモリの一部を、あるメモリから別のものに直接コピーすることができる。

【0005】

このため、VIOによって、例えばI/Oアダプタのような物理リソースをパーティション間で共有することが可能となる。リソースは1つのみのパーティション(「サーバ」パーティションと呼ぶ)に属することができるので、他のパーティションは、そのリソースを用いるためにサーバ・パーティションを通過しなければならない(他のパーティションを「クライアント」パーティションと呼ぶ)。

【発明の開示】

【発明が解決しようとする課題】

【0006】

共有リソースがI/Oアダプタである場合、サーバ・パーティションのI/Oアダプタを用いるために、サーバ・パーティションはI/OアダプタからクライアントのメモリへのDMA動作をセットアップする必要がある。しかしながら、DMA動作をセットアップするために、サーバ・パーティションは、クライアント・パーティションのメモリの構造に関する情報を知らなければならない。更に、サーバ・パーティションは、クライアント・パーティションのメモリ構造に関する情報を、そのI/Oスタックに渡すことができない。このI/Oスタックは、規定のソフトウェア・レイヤの階層であり、使用可能なフォーマットでI/O動作を実行するためには、この階層を介してデータを処理しなければならない(例えば論理ボリューム・マネージャ(LVM)、ヘッド・ドライバ、アダプタ・ドライバ、バス・ドライバ、カーネル等)。従って、サーバ・パーティショ

10

20

30

40

50

ンの環境においてクライアント・パーティションのメモリを記述し、この記述をサーバ・パーティションのI/Oスタックに容易に渡すことができるようなシステムおよび方法を有することは有利であろう。

【課題を解決するための手段】

【0007】

本発明は、サーバ・パーティション環境においてメモリ間記述子の使用を拡張して、例えばクライアント・パーティションのメモリ（以降、「遠隔」メモリと称する）のような別のパーティションのメモリを記述するために使えるようにするシステムおよび方法を提供する。このシステムおよび方法によって、論理パーティションの立ち上げの間等に、論理的にパーティションに分けたコンピューティング・システムにおけるオペレーティング・システムを初期化すると、デバイス・ドライバのローディング等の起動動作を実行するために、オペレーティング・システム・カーネル・サービスを呼び出す。これらのオペレーティング・システム・カーネル・サービスの一部として、あるカーネル・サービスを呼び出して、他の論理パーティションの遠隔メモリを本オペレーティング・システムにアタッチし、遠隔メモリ・コピーおよびDMA動作が他の論理パーティションの遠隔メモリによって実行可能であるようにする。

10

【0008】

アタッチ・カーネル・サービスは、コンピューティング・システムのデバイス・ツリーを調査する。これは、オープン・ファームウェア等によって生成可能であり、サーバ・パーティションにおけるオペレーティング・システムがアクセスできるようにサーバ・パーティションのメモリに維持することができる。論理パーティション当たり1つのデバイス・ツリー・インスタンスがある。デバイス・ツリーは、そのパーティションにおけるリソースのためのノードを含む。これらのリソースは、パーティションごとのローカル・メモリ、PCI/ISA I/Oスロット、プロセッサ、サポートされたファームウェア・コール等を含む。

20

【0009】

コンピューティング環境においてVIOをサポートする場合、典型的な実施形態において、サーバ・パーティションのデバイス・ツリーにおけるノードまたはその子ノードは、クライアント・パーティションのローカル・メモリの記述を含む。この記述は、クライアントのローカル・メモリの論理識別子、ローカル・メモリの開始アドレス、およびローカル・メモリの長さを含む。この情報から、各論理パーティションのローカル・メモリごとの拡張メモリ間記述子を生成することができる。次いで、これらの拡張メモリ間記述子は、遠隔メモリ・コピーおよびDMA動作を実行する際に用いるため、現行パーティションのオペレーティング・システムによって維持される。

30

【0010】

動作を実行するために遠隔メモリに対するアクセスが必要である場合、サーバ・パーティションのオペレーティング・システムは、遠隔メモリ動作の対象であるクライアント・パーティションの「遠隔」メモリを記述する拡張メモリ間記述子を用いる。本発明の拡張メモリ間記述子は、メモリ間記述子を用いるパーティションのローカル・メモリのみを記述するために通常用いられるデータ構造である既存のメモリ間記述子の拡張である。

40

【0011】

本発明においては、メモリ間記述子のフィールドを拡張して、別のパーティションのローカル・メモリすなわち「遠隔」メモリに関する情報をストアするために用いる。特に、拡張メモリ間記述子のフィールドは、アドレス空間識別子、遠隔メモリのサイズ、遠隔メモリを識別するための識別子（いくつかのクライアント・パーティションがある場合）、および、遠隔メモリ内の開始アドレスを含む。アドレス空間識別子は、メモリ間記述子を用いて「遠隔」メモリを表す場合、遠隔メモリ記述としてメモリ間記述子を識別する。このため、メモリ間記述子を調べるI/Oスタックまたはメモリ管理サービスにおけるいずれのインタフェースも、ここで、アドレス空間識別子を見て、これが遠隔メモリ記述であることを認識することができる。この結果、インタフェースは、DMA動作をセットアッ

50

プして「遠隔」メモリにアクセスするための適切なアクションを取ることができる。

【0012】

本発明のこれらおよび他の特徴および利点は、好適な実施形態の以下の詳細な説明において記載され、または、それを考慮することによって当業者に明らかとなる。

【0013】

ここで、本発明の好適な実施形態について、添付図面を参照して、一例としてのみ説明する。

【発明を実施するための最良の形態】

【0014】

提供されるのは、第1のパーティション環境においてメモリ間記述子を拡張することによって、別のパーティションのローカル・メモリを記述するようにして、直接メモリ・アクセスまたは遠隔メモリ・コピー動作を必要とする動作を容易にするための機構である。このため、本発明の実施形態は、好ましくは、論理分割機構を用いてリソースがパーティションに分けられているコンピューティング・デバイスにおいて実装される。

【0015】

ここで図面、特に図1を参照すると、本発明の実施形態を実装可能なデータ処理システムのブロック図が示されている。データ処理システム100は、システム・バス106に接続された複数のプロセッサ101、102、103、および104を含む対称マルチプロセッサ(SMP)・システムとすることができる。例えば、データ処理システム100は、IBM社(ニューヨーク州アーモンク)の製品であるIBM(R) eServer(R)とし、ネットワーク内のサーバとして実装することができる。あるいは、単一プロセッサ・システムを用いることも可能である。また、システム・バス106には、複数のローカル・メモリ160~163に対するインタフェースを提供するメモリ・コントローラ/キャッシュ108も接続されている。I/Oバス・ブリッジ110が、システム・バス106に接続され、I/Oバス112に対するインタフェースを提供する。メモリ・コントローラ/キャッシュ108およびI/Oバス・ブリッジ110は、図示するように統合することができる。

【0016】

データ処理システム100は、論理パーティション化(LPAR)データ処理システムである。このため、データ処理システム100は、同時に動作する複数の異種オペレーティング・システム(または単一オペレーティング・システムの複数インスタンス)を有することができる。これらのオペレーティング・システムの各々は、その内部で実行する任意の数のソフトウェア・プログラムを有することができる。データ処理システム100は、論理的にパーティション化して、異なるPCI I/Oアダプタ120~121、128~129、および136、グラフィック・アダプタ148、ならびにハード・ディスク・ドライブ149を、異なる論理パーティションに割り当てることができるようになっている。この場合、グラフィック・アダプタ148は、ディスプレイ・デバイス(図示せず)に対する接続を与え、ハード・ディスク・アダプタ149は、ハード・ディスク150を制御するための接続を与える。

【0017】

このため、例えば、データ処理システム100が3つの論理パーティションすなわちP1、P2、およびP3に分割されると仮定する。PCI I/Oアダプタ120~121、128~129、136の各々、グラフィック・アダプタ148、ハード・ディスク・アダプタ149、ホスト・プロセッサ101~104の各々、およびローカル・メモリ160~163からのメモリが、3つのパーティションの各々に割り当てられる。これらの例では、メモリ160~163は、デュアル・インライン・メモリ・モジュール(DIMM)の形態を取ることができる。DIMMは、通常、DIMMごとにパーティションに割り当てられない。その代わりに、パーティションは、プラットフォームから見たメモリ全体の一部分を得る。例えば、プロセッサ101、ローカル・メモリ160~163からのメモリの一部、およびI/Oアダプタ120、128、および129は、論理パーティショ

10

20

30

40

50

ン P 1 に割り当てることができる。プロセッサ 1 0 2 ~ 1 0 3、ローカル・メモリ 1 6 0 ~ 1 6 3 からのメモリの一部、I / O アダプタ 1 2 1 および 1 3 6 は、パーティション P 2 に割り当てることができる。プロセッサ 1 0 4、ローカル・メモリ 1 6 0 ~ 1 6 3 からのメモリの一部、グラフィック・アダプタ 1 4 8、およびハード・ディスク・アダプタ 1 4 9 は、論理パーティション P 3 に割り当てることができる。

【 0 0 1 8 】

データ処理システム 1 0 0 内で実行する各オペレーティング・システムは、異なる論理パーティションに割り当てられる。このため、データ処理システム 1 0 0 内で実行する各オペレーティング・システムは、その論理パーティション内にある I / O ユニットにのみアクセスすることができる。このため、例えば、A I X (R) オペレーティング・システム 10 の一つのインスタンスをパーティション P 1 内で実行し、A I X (R) オペレーティング・システムの第 2 のインスタンス (イメージ) をパーティション P 2 内で実行し、L i n u x (R) または O S / 4 0 0 (R) オペレーティング・システムを論理パーティション P 3 内で動作することができる。

【 0 0 1 9 】

I / O バス 1 1 2 に接続された周辺機器相互接続 (P C I) ホスト・ブリッジ 1 1 4 は、P C I ローカル・バス 1 1 5 に対するインタフェースを提供する。複数の P C I 入力 / 出力アダプタ 1 2 0 ~ 1 2 1 を、P C I - P C I ブリッジ 1 1 6、P C I バス 1 1 8、P C I バス 1 1 9、I / O スロット 1 7 0、および I / O スロット 1 7 1 を介して、P C I バス 1 1 5 に接続することができる。P C I - P C I ブリッジ 1 1 6 は、P C I バス 1 1 8 および P C I バス 1 1 9 に対するインタフェースを与える。P C I I / O アダプタ 1 2 0 および 1 2 1 は、I / O スロット 1 7 0 および 1 7 1 内にそれぞれ配置される。典型的な P C I バス実装は、4 個乃至 8 個の I / O アダプタ (すなわちアドイン・コネクタのための拡張スロット) をサポートする。各 P C I I / O アダプタ 1 2 0 ~ 1 2 1 は、データ処理システム 1 0 0 と、例えばデータ処理システム 1 0 0 のクライアントである他のネットワーク・コンピュータ等の入出力デバイスとの間のインタフェースを与える。 20

【 0 0 2 0 】

追加の P C I ホスト・ブリッジ 1 2 2 は、追加の P C I バス 1 2 3 のためのインタフェースを提供する。P C I バス 1 2 3 は、複数の P C I I / O アダプタ 1 2 8 ~ 1 2 9 に接続されている。P C I I / O アダプタ 1 2 8 ~ 1 2 9 は、P C I - P C I ブリッジ 1 2 4、P C I バス 1 2 6、P C I バス 1 2 7、I / O スロット 1 7 2、および I / O スロット 1 7 3 を介して、P C I バス 1 2 3 に接続することができる。P C I - P C I ブリッジ 1 2 4 は、P C I バス 1 2 6 および P C I バス 1 2 7 に対するインタフェースを提供する。P C I I / O アダプタ 1 2 8 および 1 2 9 は、I / O スロット 1 7 2 および 1 7 3 内にそれぞれ配置される。このように、例えばモデムまたはネットワーク・アダプタ等の追加の I / O デバイスを、P C I I / O アダプタ 1 2 8 ~ 1 2 9 の各々を介してサポートすることができる。このように、データ処理システム 1 0 0 は、複数のネットワーク・コンピュータに対する接続を可能とする。 30

【 0 0 2 1 】

I / O スロット 1 7 4 に挿入されたメモリ・マップ・グラフィック・アダプタ 1 4 8 は、P C I バス 1 4 4、P C I - P C I ブリッジ 1 4 2、P C I バス 1 4 1、および P C I ホスト・ブリッジ 1 4 0 を介して I / O バス 1 1 2 に接続することができる。ハード・ディスク・アダプタ 1 4 9 は、P C I バス 1 4 5 に接続された I / O スロット 1 7 5 内に配置される。そして、このバスは、P C I - P C I ブリッジ 1 4 2 に接続され、これは、P C I バス 1 4 1 によって P C I バス・ブリッジ 1 4 0 に接続されている。 40

【 0 0 2 2 】

P C I ホスト・ブリッジ 1 3 0 は、P C I バス 1 3 1 のためのインタフェースを提供して、I / O バス 1 1 2 に接続する。P C I I / O アダプタ 1 3 6 は、I / O スロット 1 7 6 に接続され、これは、P C I バス 1 3 3 によって P C I - P C I ブリッジ 1 3 2 に接続されている。P C I - P C I ブリッジ 1 3 2 は、P C I バス 1 3 1 に接続されている。 50

また、このPCIバスは、PCIホスト・ブリッジ130を、サービス・プロセッサ・メールボックス・インタフェースおよびISAバス・アクセス通過論理194およびPCI-PCIブリッジ132に接続する。サービス・プロセッサ・メールボックス・インタフェースおよびISAバス・アクセス通過論理194は、PCI/ISAブリッジ193に宛てられたPCIアクセスを転送する。NVRAMストレージ192は、ISAバス196に接続されている。サービス・プロセッサ135は、サービス・プロセッサ・メールボックス・インタフェースおよびISAバス・アクセス通過論理194に、そのローカルPCIバス195を介して結合されている。また、サービス・プロセッサ135は、複数のJTAG/I²Cバス134を介してプロセッサ101~104に接続されている。JTAG/I²Cバス134は、JTAG/スキャン・バス(IEEE1149.1を参照)およびPhillips I²Cバスを組み合わせたものである。しかしながら、代替として、JTAG/I²Cバス134は、Phillips I²CバスのみまたはJTAG/スキャン・バスのみによって置換することができる。ホスト・プロセッサ101、102、103、および104の全てのSP-ATTN信号は、共に、サービス・プロセッサの割り込み入力信号に接続されている。サービス・プロセッサ135は、それ自身のローカル・メモリ191を有し、ハードウェアOPパネル190に対するアクセスを有する。

【0023】

データ処理システム100に最初に電源投入した場合、サービス・プロセッサ135は、JTAG/I²Cバス134を用いて、システム(ホスト)・プロセッサ101~104、メモリ・コントローラ/キャッシュ108、およびI/Oブリッジ110に質問する。このステップが完了すると、サービス・プロセッサ135は、データ処理システム100のインベントリおよびトポロジの理解を得る。また、サービス・プロセッサ135は、ホスト・プロセッサ101~104、メモリ・コントローラ/キャッシュ108、およびI/Oブリッジ110に質問することによって見出された全ての要素に対して、内蔵自己テスト(BIST)、基本検証テスト(BAT)、およびメモリ・テストを実行する。BIST、BAT、およびメモリ・テストの間に検出された故障についてのエラー情報は全て、サービス・プロセッサ135によって収集されて報告される。

【0024】

BIST、BAT、およびメモリ・テストの間に故障と見なされた要素を除外した後に、システム・リソースの意味のある有効な構成が可能である場合、データ処理システム100は、先に進んで、実行可能コードをローカル(ホスト)・メモリ160~163にロードすることができる。次いで、サービス・プロセッサ135は、ローカル・メモリ160~163にロードされたコードを実行するため、ホスト・プロセッサ101~104を解放する。ホスト・プロセッサ101~104がデータ処理システム100内の各オペレーティング・システムからのコードを実行している間、サービス・プロセッサ135はエラーの監視および報告のモードに入る。サービス・プロセッサ135が監視するアイテムのタイプは、例えば、冷却ファンの速度および動作、熱センサ、電源レギュレータ、ならびに、プロセッサ101~104、ローカル・メモリ160~163、およびI/Oブリッジ110によって報告された回復可能および回復不可能エラーを含む。

【0025】

サービス・プロセッサ135は、データ処理システム100内の監視されている全アイテムに関するエラー情報を保存し報告する責任がある。また、サービス・プロセッサ135は、エラーのタイプおよび規定された閾値に基づいてアクションを取る。例えば、サービス・プロセッサ135は、プロセッサのキャッシュ・メモリ上の過度の回復可能エラーを記録し、これがハード故障を予測するものであることを決定することができる。この決定に基づいて、サービス・プロセッサ135は、現在の実行セッションおよび今後の初期プログラム・ロード(IPL)の間に構成解除するためにそのリソースをマークすることができる。また、IPLは、時に、「ブート」または「ブートストラップ」と呼ばれる。

【0026】

10

20

30

40

50

データ処理システム100は、様々な市販のコンピュータ・システムを用いて実施可能である。例えば、データ処理システム100は、IBM社から入手可能なIBM(R) eServer(R) iSeries(R) Model 840システムを用いて実装することができる。かかるシステムは、OS/400(R)オペレーティング・システムを用いた論理分割をサポートすることができる。OS/400(R)オペレーティング・システムも、IBM社から入手可能である。

【0027】

図1に示すハードウェアが様々に変化し得ることは、当業者には認められよう。例えば、図示したハードウェアに加えて、またはその代わりに、光ディスク・ドライブ等の他の周辺デバイスも使用可能である。図示した例は、本発明に対してアーキテクチャ上の限定を加えることを意味しない。

10

【0028】

ここで図2を参照すると、本発明の実施形態を実装することができる例示的な論理パーティション化プラットフォームのブロック図が示されている。論理パーティション化プラットフォーム200におけるハードウェアは、例えば図1のデータ処理システム100として実装可能である。論理パーティション化プラットフォーム200は、パーティション化ハードウェア230、オペレーティング・システム202、204、206、208、および、パーティション管理ファームウェア210を含む。オペレーティング・システム202、204、206、および208は、単一オペレーティング・システムの複数コピー、または、論理パーティション化プラットフォーム200上で同時に動作する複数の異種オペレーティング・システムとすることができる。これらのオペレーティング・システムは、OS/400(R)を用いて実装可能であり、これは、ハイパーバイザ等のパーティション管理ファームウェアとインタフェースするように設計されている。OS/400(R)は、これらの例示的な実施形態において、一例としてのみ用いられる。むしろ、具体的な実施に応じて、AIX(R)およびLinux(R)等の他のタイプのオペレーティング・システムも使用可能である。オペレーティング・システム202、204、206、および208は、パーティション203、205、207、および209に配置されている。ハイパーバイザ・ソフトウェアは、パーティション管理ファームウェア210を実装するために使用可能であるソフトウェアの一例であり、IBM社から入手可能である。ファームウェアは、例えばリード・オンリ・メモリ(ROM)、プログラマブルROM(PROM)、消去可能プログラマブルROM(EPROM)、電氣的消去可能プログラマブルROM(EEPROM)、および不揮発性ランダム・アクセス・メモリ(不揮発性RAM)等の、電力なしでその内容を保持するメモリ・チップにストアされる「ソフトウェア」である。

20

30

【0029】

更に、これらのパーティションは、パーティション・ファームウェア211、213、215、および217も含む。パーティション・ファームウェア211、213、215、および217は、初期ブート・ストラップ・コード、IEEE-1275標準オープン・ファームウェア、および、IBM社から入手可能なランタイム・アブストラクション・ソフトウェア(RTAS)を用いて実装可能である。パーティション203、205、207、および209をインスタンス化するとき、プラットフォーム・ファームウェア210によって、ブート・ストラップ・コードのコピーがパーティション203、205、207、および209にロードされる。その後、制御はブート・ストラップ・コードに移され、次いでブート・ストラップ・コードは、オープン・ファームウェアおよびRTASをロードする。次いで、パーティションに関連するまたはパーティションに割り当てられたプロセッサが、パーティションのメモリにディスパッチされて、パーティション・ファームウェアを実行する。

40

【0030】

パーティション化ハードウェア230は、複数のプロセッサ232~238、複数のシステム・メモリ・ユニット240~246、複数の入出力(I/O)アダプタ248~2

50

62、およびストレージ・ユニット270を含む。プロセッサ232~238、メモリ・ユニット240~246、NVRAMストレージ298、およびI/Oアダプタ248~262の各々は、論理パーティション化プラットフォーム200内の複数のパーティションのうち1つに割り当てることができる。パーティションの各々は、オペレーティング・システム202、204、206、および208に対応する。

【0031】

パーティション管理ファームウェア210は、パーティション203、205、207、および209のための多数の機能およびサービスを実行して、論理パーティション化プラットフォーム200のパーティショニングを実施する。パーティション管理ファームウェア210は、基礎にあるハードウェアと同一のファームウェア実装仮想マシンである。このため、パーティション管理ファームウェア210は、論理パーティション化プラットフォーム200の全てのハードウェア・リソースを仮想化することによって、独立したOSイメージ202、204、206、および208の同時実行を可能とする。

10

【0032】

サービス・プロセッサ290を用いて、パーティションにおけるプラットフォーム・エラーの処理等、様々なサービスを提供することができる。また、これらのサービスは、IBM社等のベンダーにエラーを報告するサービス・エージェントとして機能することができる。異なるパーティションの動作は、ハードウェア管理コンソール280等のハードウェア管理コンソールを介して制御することができる。ハードウェア管理コンソール280は、別個のデータ処理システムであり、これによって、システム・アドミニストレータは、異なるパーティションへのリソース再割り振りを含む様々な機能を実行することができる。

20

【0033】

かかる論理パーティション化環境において、カーネルと、現在のプロセス・アドレス空間以外のアドレス空間との間でデータを動かそうとする場合はいつでも、メモリ間カーネル・サービスを利用してこのデータ移動を実行する。アドレス空間の1つの領域内のデータ・エリアが、メモリ間動作のためのユーザ・バッファにアタッチするxmattnカーネル・サービスを呼び出すことによってアタッチされる。呼び出されると、xmattnカーネル・サービスはメモリ間記述子を発生する。その後、他のメモリ間カーネル・サービスを利用して、カーネルから現在のプロセス・アドレス以外のアドレス空間に、データを移動またはコピーすることができる。例えば、xmeminカーネル・サービスは、指定されたアドレス空間からカーネル・グローバル・メモリにデータをコピーすることによってメモリ間移動を実行する。xmemoutカーネル・サービスは、カーネル・グローバル・メモリから指定されたアドレス空間にデータをコピーすることによってメモリ間移動を実行する。xmemdmaカーネル・サービスは、DMA I/Oのためのメモリのページを準備するか、または、DMA I/Oが完了した後にページを処理する。

30

【0034】

既知のコンピューティング・システムにおいて、メモリ間記述子は、ローカル・メモリを記述するために用いられるデータ構造である。ローカル・メモリのメモリ間記述子は、オペレーティング・システムの仮想メモリ管理(VMN)コンポーネントによって維持される情報から発生される。

40

【0035】

本発明の好適な実施形態による拡張メモリ間記述子は、ブート時にオペレーティング・システムによって発生されパースされるデバイス・ツリーに基づいてカーネル・サービスによって発生される。デバイス・ツリーは、システム・ハードウェアおよびユーザ構成選択を記述する階層データ構造であるオープン・ファームウェア・デバイス・ツリーと同様である。また、オープン・ファームウェア・デバイス・ツリーは、ハードウェア・ドライバと、これらのドライバによって用いるサポート・ルーチンを含む。以下は、オープン・ファームウェア・デバイス・ツリーの一例である。

50

【 0 0 3 6 】	
root /	
ff8885f8 /rtas	
ff866bf4/rom@ff000000	
ff8627c0/flash@fff00000	
ff8513d0 /cpus	
ff88a1c8 /PowerPC,604ev@0	
ff88a788 /l2-cache	
ff84d8e0 /pci	
ff89552c /ethernet@4	10
ff8952e0 /display	
ff88cf44 /mac-io@2	
ff893e68 /misc@0	
ff894688 /iic	
ff893d7c /via@16000	
ff8939fc /escc@13000	
ff893c2c/ch-b@13000	
ff893af4/ch-a@13020	
ff88fd08 /scsi@10000	
ff892e1c /tape	20
ff8924d0 /disk	
ff88f944 /escc-legacy@12000	
ff88fba8/ch-b@12000	
ff88fa60/ch-a@12002	
ff88d78c /adb@11000	
ff88f364 /mouse@3	
ff88e6dc/keyboard@2	
ff88d638 /open-pic@40000	
ff88aad8 /ide@1,1	
ff88c504 /disk	30
ff85a534 /isa@1	
ff864208 /sound@i534	
ff8640e8 /midi@i330	
ff863ff8 /game@i200	
ff863368 /gpio@i800	
ff863008 /nvram@me0000	
ff862aa4 /rtc@i70	
ff85f644 /8042@i60	
ff8618b0 /mouse@aux	
ff860260 /keyboard@	40
ff85d804 /floppy@i3f0	
ff85d3b4 /parallel@i3bc	
ff85c704 /serial@i2f8	
ff85b9f4 /serial@i3f8	
ff85b490 /timer@i40	
ff85b01c /interrupt-controller@i20	
ff85ae08 /dma-controller@i00	
ff84a650 /mmu	
ff83f2e4/memory@0	
【 0 0 3 7 】	50

本発明の一実施形態によって、同様のデバイス・ツリー表現を用いて、拡張メモリ間記述子を発生させるための情報を得る。デバイス・ツリーは、オブジェクト指向であり、デバイス・ツリー内のエントリの各々が関連する特性を有し、これを調べることができるようになっている。デバイス・ツリーにおいて、遠隔メモリに関する情報をオペレーティング・システムに提示することができるノードの一例は、以下に示す通りである。

【 0 0 3 8 】

/devicenode

...

remote-memory-info10000000 00200000 00040000

20000000 003D0000 00020000

...

...

【 0 0 3 9 】

ここで、「remote-memory-info(遠隔メモリ情報)」は、ノード「/devicenode」の特性である。クライアント - サーバ仮想 I/O モデルの典型的な実装において、ノード「/devicenode」は、サーバ・デバイスのノードを表す。以下に述べるように、このノードはパースすることができ、このノードのパースから得た情報を用いて、本発明の好適な実施形態に従った拡張メモリ間記述子を発生することができる。このため、「10000000」エントリは、メモリの論理識別子であり、「00200000」エントリは、メモリの開始アドレスであり、「00400000」エントリは、最初のクライアントについてのバイト単位のメモリ・サイズである。同様に、「20000000」は、メモリの論理識別子であり、「03D00000」は、メモリ内の開始アドレスであり、「00200000」は、第2のクライアントについてのバイト単位のメモリ・サイズである。このトリプレット(論理 id、開始アドレス、サイズ)は、サーバが対応するクライアントごとに繰り返すことができる。

【 0 0 4 0 】

このデバイス・ツリー情報から、システム・ハードウェアおよびユーザ構成選択の詳細な記述を得ることができる。特に、先に示したように、デバイス・ツリーのメモリ・ノードは、メモリの論理識別子、メモリの開始アドレス、およびメモリの長さを示す特性を含む。この情報は、カーネル・サービスによって用いて、コンピューティング・システムのローカル・メモリ・リソースのためのメモリ間記述子を発生させることができる。

【 0 0 4 1 】

メモリ間記述子を、対象のメモリの所有者の知識を必要とせずに、2つの無関係なコンテキストからメモリにアクセスするための方法として用いる。例えば、割り込みハンドラがデータをプロセスに渡す必要がある場合、これは、メモリ間記述子を用いて、プロセスのメモリ空間に関する情報を得て、プロセスのメモリ空間にデータをコピーする。これが必要である理由は、割り込みハンドラが、必ずしもそれが所有するデータの処理に関連して動作しているわけではなく、従って、どのプロセスにデータが属するかに関して知らないからである。割り込みハンドラが知っているのは、データをコピーしなければならないメモリの記述、すなわちメモリ間記述子である。

【 0 0 4 2 】

ローカル・メモリのメモリ間記述子の使用の別の例は、ディスク I/O である。例えば、ユーザ・レベルのアプリケーションは、I/OCTL を呼び出して、ディスクに関する情報を読むことができる。この情報は、ユーザ空間内にある、アプリケーションが供給するバッファにストアされる。ディスク・ドライバが呼び出されると、ディスク・ドライバは xmatch を呼び出してアプリケーションのバッファをアタッチする。これは全て、ディスク・ドライバの上半分によって行われる。

【 0 0 4 3 】

将来のある時点で、I/O 要求が完了し、データが利用可能となった場合、ディスク・ドライバの下半分内にあるディスク・ドライバの「iodone」ルーチンが、xmount を呼び出して、メモリ間記述子を用いてユーザ空間にデータをコピーする。メモリ間記

10

20

30

40

50

述子は、ディスク・ドライバによってストアされ、上半分および下半分の双方に利用可能とする。

【 0 0 4 4 】

メモリ間記述子は、オペレーティング・システムの初期化の際にカーネル・サービスによって発生され、必要に応じて、コンピューティング・デバイスの I/O スタックにおいて、インタフェースからインタフェース（またはメモリ管理サービス）へと渡される。I/O スタックのインタフェースおよびメモリ管理サービスは、メモリ間記述子を調べて、メモリ間動作を実行するために必要な情報を得る。

【 0 0 4 5 】

オペレーティング・システムは、別のパーティションの遠隔メモリに関して同様の問題に直面する。なぜなら、遠隔メモリ・コピー/DMA動作は、別のパーティションにおける所有者デバイスの知識を持たないからである。このため、本明細書中に記載する手法は、メモリ間記述子の拡張バージョンを利用して、遠隔メモリ・コピー/DMA動作を実行する場合に、別のパーティションの遠隔メモリに関する情報を提供する。

10

【 0 0 4 6 】

ローカル・メモリ間記述子は、4つの主フィールドから成る。すなわち、アドレス空間識別子、セグメント数、開始セグメント識別子、およびセグメント内の開始アドレスである。アドレス空間識別子は、データをコピーするメモリの一意的識別子を提供する。セグメント数は、アドレス空間識別子が識別するアドレス空間のセグメントの合計数を提供する。開始セグメント識別子は、セグメントの合計数内で、コピー対象データの書き込みを始めるセグメントを識別する。セグメント内の開始アドレスは、開始セグメント識別子が識別したセグメント内で、コピー対象データの書き込みを始めるアドレスを識別する。

20

【 0 0 4 7 】

以下は、これらの主フィールドを例示するメモリ間データ構造の一例である。

【 0 0 4 8 】

```
enum asid={LOCAL=0, REMOTE=1};
struct cross_memory {
    enum asid address_space_id; /* LOCAL, REMOTE, etc. */
    int number_of_segments; /* number of segments if address_space_id=LOCAL */
    long segment_id; /* segment number if address_space_id=LOCAL */
    char *vaddr; /* starting address within the segment if address_space_id=LOCAL */
    uint flags; /* any special attributes of the segment or address range */
};
#define remote_size(number_of_segments) /* size of remote memory if address_space_id=REMOTE */
#define remote_logical_id(segment_id) /* remote memory's logical identifier */
#define remote_start_addr(vaddr) /* starting address within the remote memory */
```

30

【 0 0 4 9 】

上記のフィールドから明らかなのは、特定のアドレス空間の所有者を知る必要がなく、従って、アドレス空間にアクセスするために所有者にデータを渡す必要なしにメモリ間動作を実行するために使用可能であるという特異なメモリの記述を、メモリ間記述子と与えることである。本発明の実施形態は、この既知のデータ構造を利用して、完了のためにリソースの所有権に関する情報を必要としない異なる動作を容易にする。すなわち、メモリ間記述子の拡張バージョンを用いて、あるパーティションから別のパーティションに対する遠隔メモリ・コピー/DMA動作を容易にする（ここで、「遠隔」メモリは、現在実行しているプロセスのパーティションとは異なるパーティションにおけるメモリである）。かかる動作により、あるパーティションにおけるプロセスは、仮想入出力（VIO）プロセスの一部として別のパーティションのメモリにアクセスすることができ、遠隔メモリの所有権を有する特定のプロセスにデータを渡す必要はない。この結果、拡張メモリ間記述子を用いることによって、パーティション間の遠隔メモリ・コピーおよびDMA動作が可

40

50

能となる。

【0050】

メモリ間記述子フィールドを用いて、遠隔メモリに関する情報をストアする。メモリ間記述子の元の4つのフィールドをオーバーロードして、遠隔メモリを記述し、これによって、以下の4つのフィールド、すなわち、アドレス空間識別子、遠隔メモリのサイズ、遠隔メモリを識別する識別子、および遠隔メモリ内の開始アドレスを有する拡張メモリ間記述子を発生する。拡張メモリ間記述子のアドレス空間識別子は、元のメモリ間記述子のアドレス空間識別子と同様であるが、異なる点は、このフィールドに所定の値をストアして、拡張メモリ間記述子がクライアント・パーティションの遠隔メモリに対応することを示すことである。

10

【0051】

遠隔メモリのサイズは、バイトまたは他の任意のメモリ・サイズ単位で指定することができ、アクセスを実行する遠隔メモリの合計サイズを指定する。遠隔メモリ識別子は、特定のクライアント・パーティションにおける特定の遠隔メモリを識別する。これは、複数のクライアント・パーティションがあり、どの遠隔メモリにアクセスしているかを指定する必要がある場合に重要である。遠隔メモリ内の開始アドレスは、遠隔メモリ識別子が指定する遠隔メモリ内のアクセスを開始する位置を提供する。遠隔メモリ・コピー/DMA動作を実行するため、これらの4つのフィールドから、遠隔メモリの全体像を得る。

【0052】

拡張メモリ間記述子は、元のメモリ間記述子と同じ方法で、I/Oスタックを通してインタフェースからインタフェースへ渡すことができる。拡張メモリ間記述子を調べる必要があるインタフェースおよびメモリ管理サービスは、拡張メモリ間記述子进行处理するためのコードを含むように増強する。すなわち、それらのインタフェースおよびメモリ管理サービスは、アドレス空間識別子が拡張メモリ間記述子を識別する所定の値を含む場合、それを遠隔メモリ・アクセス要求のための遠隔メモリを記述するとして認識するコードを含むように増強する。次いで、インタフェースおよびメモリ管理サービスは、拡張メモリ間記述子が提供する遠隔メモリのサイズ、遠隔メモリの識別子、および遠隔メモリ内の開始アドレスに基づいて、適切なアクションを実行することができる。例えば、`d_map_page`および`d_map_list`インタフェースは、遠隔メモリ間記述子を用いて、サーバ・パーティションが所有する物理デバイスから遠隔メモリへのDMAを正しくセットアップする。このため、遠隔メモリ間記述子は、物理デバイスのドライバによって、`d_map_page/list`インタフェースに渡されて、DMAをセットアップする。サーバ・パーティションの定義により、物理デバイスはサーバ・パーティションに属する。この結果、遠隔メモリ・コピー/DMA動作をパーティション間で実行することができる。

20

30

【0053】

元のメモリ間記述子と同じフォーマットを本質的に有するが、メモリ間記述子のフィールドにストアされた異なるデータおよびこの異なるデータに関連する異なる含意 (connotation) を有する拡張メモリ間記述子に、既存のメモリ間記述子をオーバーロードすることによって、本発明の実施形態は、メモリ間記述子を必要とするがメモリ間記述子の内容を調べないインタフェースおよびメモリ管理サービスの多くを、変更なしで動作させることができる。すなわち、拡張メモリ間記述子によって、オペレーティング・システムは、メモリ間記述子を必要とするが実際にそれを調べないインタフェースおよびメモリ管理サービスとのバイナリ互換性を維持することができる。例えば、ほとんどのデバイス・ドライバは、メモリ間記述子を必要とするが、それを調べない。それらは、単に、メモリ間記述子をカーネルまたはPCIバス・ドライバに渡すだけである。かかるデバイス・ドライバは、メモリ間記述子に対する変更のため、再コンパイルする必要がなく、それらは、いつもと変わらず動作し続けることができる。

40

【0054】

図3は、本発明の1つの例示的な実施形態に従った拡張メモリ間記述子を発生させるための機構の一例を示す例示的な図である。図3に示すように、論理パーティション310

50

のブートアップの間等、論理パーティション化コンピューティング・システムにおいて論理パーティション310のオペレーティング・システム320を初期化すると、デバイス・ドライバのローディング等のブートアップ動作を実行するため、オペレーティング・システム・カーネル・サービス325が呼び出される。これらのオペレーティング・システム・カーネル・サービス325の一部として、カーネル・サービスを呼び出し、他の論理パーティション312～316の遠隔メモリを、本オペレーティング・システムにアタッチし、遠隔メモリ・コピーおよびDMA動作が他の論理パーティションの遠隔メモリによって実行可能であるようにする。本発明の1つの例示的な実施形態では、遠隔メモリをアタッチするカーネル・サービスは、xmattach_remioと称し、既知のxmattachメモリ間カーネル・サービスの拡張バージョンである。xmattach_remioは、既知のxmemoutおよびxmeminメモリ間カーネル・サービスならびにd_map_pageおよびd_map_listカーネル・サービスの拡張バージョンによって実行可能であるメモリ間動作のための遠隔メモリをアタッチする。

10

【0055】

拡張アタッチ・カーネル・サービスは、ハイパーバイザ330等に維持することができるコンピューティング・システムのデバイス・ツリー332～338を調べる。ハイパーバイザ330は、論理パーティション化コンピューティング・システムにおいて複数の仮想マシンを提供し管理する機構である。論理パーティションごとの個々のデバイス・ツリー332～338は、ハイパーバイザ330に維持される。更に、論理パーティションのためのデバイス・ツリー332～338のコピーは、論理パーティションのオペレーティング・システム・カーネルに利用可能であるように、その論理パーティションのローカル・メモリにストアすることができる。オペレーティング・システムは、デバイス・ツリー332～338をパースするためにカーネル拡張（例えばデバイス・ドライバ）およびユーザ・レベル・コードが使用可能である関数のライブラリを提供する。

20

【0056】

デバイス・ツリー332～338は、コンピューティング・システムおよびそれらのリソースのパーティションごとのノードを含む。これらのリソースは、パーティションごとのローカル・メモリ340～360を含む。パーティションのこれらのローカル・メモリを表すデバイス・ツリー332～338のノードは、ローカル・メモリの論理識別子、ローカル・メモリの開始アドレス、およびローカル・メモリの長さを含むローカル・メモリの特性を含む。この情報から、論理パーティションごとのローカル・メモリ（これは、現在の論理パーティションに対して「遠隔」メモリである）のサイズ、論理パーティションごとのローカル・メモリの識別子、および論理パーティションごとのローカル・メモリ内の開始アドレスを得ることができる。この情報は、各論理パーティションのローカル・メモリごとに、拡張メモリ間記述子370～390にパッケージすることができる。これらの拡張メモリ間記述子370～390は、次いで、他の論理パーティション312～316に対する遠隔メモリ・コピーおよびDMA動作を実行する際に用いるため、現行論理パーティション310のオペレーティング・システム320によって維持される。

30

【0057】

図4は、拡張メモリ間記述子を用いて遠隔メモリ・アクセスを必要とする動作を実行するための、本発明の1つの例示的な実施形態に従った機構の一例を示す例示的な図である。図4に示すように、クライアント・デバイス400のアプリケーション410等のプロセスが、例えば、read()システム・コールによって、遠隔メモリ・アクセス動作を実行することを求める入出力(I/O)要求を、オペレーティング・システム・カーネル空間420に発生すると、このI/O要求は、I/Oスタック422からクライアント仮想デバイス・ドライバ424に至るパスで処理される。クライアント仮想デバイス・ドライバ424は、この要求を、サーバ480の対応するもの、すなわちサーバ仮想ドライバ430に転送する。

40

【0058】

サーバ仮想ドライバ430は、ブート時に、例えば論理パーティション2のような遠隔

50

パーティションのメモリ 415 のための拡張メモリ間記述子 440 をすでに生成している。例えば、サーバ仮想ドライバ 430 は、ブート時に、xmattach_remioカーネル・サービスを呼び出して、拡張メモリ間記述子を発生することができる。サーバ仮想ドライバ 430 は、拡張メモリ間記述子 440 と共に、I/O要求を、サーバ側のI/Oスタック 450 に渡す。

【0059】

I/O要求は、最終的に、I/Oスタック 450 を介して、例えばSCSIアダプタ・ドライバのような物理デバイス・ドライバ 460 に渡される。これは、そのI/Oアダプタ 470 を使用するDMA動作をスケジューリングする。I/Oアダプタ 470 によってDMAが完了すると、データは遠隔メモリ 415 に書き込まれる。このため、先の例に示したように、データ自体に加えて、遠隔直接メモリ・アクセスを実行するために要求される他の必要なエレメントは、拡張メモリ間記述子によって与えられる遠隔メモリの記述である。

10

【0060】

図5および図6は、本発明の例示的な実施形態による、拡張メモリ間記述子を発生し用いるための例示的なプロセスを示すフローチャートである。フローチャート図の各ブロックおよびフローチャート図のブロックの組み合わせが、コンピュータ・プログラム命令によって実装可能であることは、理解されよう。これらのコンピュータ・プログラム命令を、プロセッサまたは他のプログラマブル・データ処理装置に提供して、マシンを生成し、それにより、プロセッサまたは他のプログラマブル・データ処理装置上で実行する命令で、フローチャートのブロックに明示される機能を実装するための手段を生成することができる。また、これらのコンピュータ・プログラム命令は、コンピュータ読み取り可能メモリまたはストレージ媒体にストアして、プロセッサまたは他のプログラマブル・データ処理装置を特定の方法で機能させることができ、それにより、コンピュータ読み取り可能メモリまたはストレージ媒体にストアされた命令が、フローチャートのブロックに明示される機能を実装する命令手段を含む製造物を生成するようになっている。

20

【0061】

従って、フローチャート図のブロックは、明示した機能を実行するための手段の組み合わせ、明示した機能を実行するためのステップの組み合わせ、および、明示した機能を実行するためのプログラムの組み合わせをサポートする。また、フローチャート図の各ブロックおよびフローチャート図のブロックの組み合わせは、明示した機能またはステップを実行する特殊目的のハードウェア・ベースのコンピュータ・システムによって、または、特殊目的の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って実装可能であることは理解されよう。

30

【0062】

図5は、本発明の1つの例示的な実施形態に従った拡張メモリ間記述子を発生するための例示的なプロセスの概要を示すフローチャートである。図5に示すように、動作は、オペレーティング・システムの初期化によって開始する(ステップ510)。次いで、オペレーティング・システム・カーネル・サービスは、他の論理パーティションの遠隔メモリを表すノードを探してコンピュータ・システムのためのデバイス・ツリーをパースする(ステップ520)。次いで、これらの遠隔メモリ・ノードのための特性を検索し(ステップ530)、他のパーティションの遠隔メモリごとに、識別子、サイズ、および開始アドレスを求めるために用いる(ステップ540)。次いで、識別子、サイズ、および開始アドレスを、遠隔識別子と共にパッケージし、他の論理パーティションの各々の遠隔メモリごとの拡張メモリ間記述子とする(ステップ550)。次いで、これらの拡張メモリ間記述子を、遠隔メモリ・アクセスを実行する際に用いるためにストアする(ステップ560)。

40

【0063】

図6は、本発明の1つの例示的な実施形態に従って拡張メモリ間記述子を用いるための例示的なプロセスの概要を示すフローチャートである。図6に示すように、動作は、遠隔

50

メモリ・アクセスを求める要求を受信することによって開始する(ステップ610)。次いで、この要求のターゲットである遠隔メモリのための拡張メモリ間記述子を検索する(ステップ620)。次いで、拡張メモリ間記述子を要求と共にI/Oスタックに渡す(ステップ630)。これは、拡張メモリ間記述子および要求をそれに応じて処理する。次いで、拡張メモリ間記述子を用いて、I/Oアダプタによって、遠隔メモリ・アクセスを実行する(ステップ640)。

【0064】

従って、既知のメモリ間記述子を拡張して、現在プロセスが実行しているパーティションとは異なるパーティションの遠隔メモリを記述することができるようにするための機構が提供される。本発明の実施形態が提供する拡張は、既知のメモリ間記述子と同一のフォーマットを用いるが、データおよびデータに関連する含意は異なるので、コンピューティング・デバイスに存在するインタフェースおよびメモリ管理サービスの多くは、拡張メモリ間記述子を扱うために変更する必要はない。メモリ間記述子を調べる他のインタフェースおよびメモリ管理サービスは、拡張メモリ間記述子を処理するためのコードを含むように増強させて、遠隔メモリ・アクセスを実行するための機能を提供することができる。

【0065】

本発明の例示的な実施形態について、典型的なDMA動作との関連で説明してきたが、本発明の実施形態はこれに限定されないことは認められよう。例えば、拡張メモリ間記述子は、遠隔コピー動作にも用いることができる。コピー動作は、ハイパーバイザ、または他のハードウェアもしくはソフトウェア機構またはその両方によって、容易にすることができる。この場合、サーバ・パーティションは、例えばハイパーバイザのような機構に対して、オペレーティング・システムとハイパーバイザとの間の詳細に記録されたインタフェースを介して情報を供給する。供給される情報は、拡張メモリ間記述子からのものである(例えば、クライアント・パーティションの遠隔メモリを識別する識別子)。次いで、ハイパーバイザは、サーバからクライアントに、またはクライアントからサーバに、実際のコピー動作を実行することができる。

【0066】

動作は、サーバ側によって開始することができ、従って、サーバは、データをプルする(例えばxmeninシステム・コールによって)か、またはクライアント・パーティションにデータをプッシュする(例えばxmemoutシステム・コールによって)ことができる。ハイパーバイザは、実際のコピー動作を実行することができる。なぜなら、その上で動作しているオペレーティング・システムよりも高い管理特権を有するからである。これによって、ハイパーバイザは、どのパーティションのどのメモリにもアクセスすることが可能となる。こうした遠隔コピー動作によって、サーバとクライアントとの間のいっそう洗練されたデータ交換プロトコルの設計および実装が容易になる。

【0067】

本発明の実施形態は、完全に機能するデータ処理システムとの関連で説明してきたが、このプロセスは、命令のコンピュータ読み取り可能媒体の形態で、および、様々な形態で分配することができ、この分配を実行するために実際に用いる信号担持媒体の特定のタイプには無関係に、ここに記載した手法が等しく適用されることは、当業者には認められよう。コンピュータ読み取り可能媒体の例は、フレキシブル・ディスク、ハード・ディスク・ドライブ、RAM、CD-ROM、DVD-ROM等の記録可能タイプの媒体、ならびに、デジタルおよびアナログ通信リンク、例えば無線周波数および光波伝送等の伝送形態を用いた有線または無線通信リンク等の伝送タイプの媒体を含む。コンピュータ読み取り可能媒体は、特定のデータ処理システムにおいて実際に用いるために復号される符号化フォーマットの形態を取ることができる。

【0068】

ここでの記載は、例示および説明の目的のために提示するものであり、網羅的であったり、開示した形態に本発明を限定したりすることは意図していない。当業者には、多く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、実際の適用を最良に説明

10

20

30

40

50

するため、および、考えられる特定の使用法に適した様々な変更を用いる様々な実施形態について当業者が本発明を理解することを可能とするために、選択して記載した。

【図面の簡単な説明】

【0069】

【図1】本発明の実施形態を実装可能なデータ処理システムのブロック図である。

【図2】本発明の実施形態を実装可能である例示的な論理パーティション化プラットフォームのブロック図である。

【図3】本発明の1つの例示的な実施形態に従って拡張メモリ間記述子を発生させるための機構の一例を示す例示的な図である。

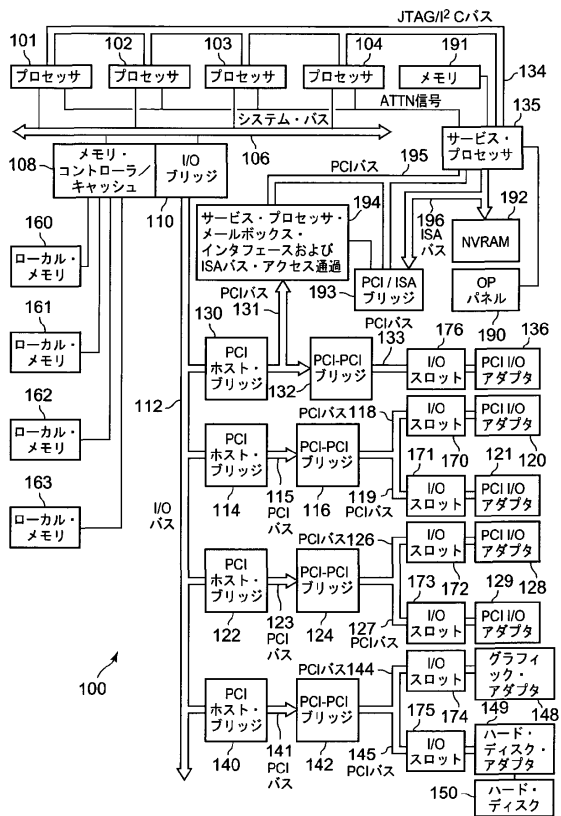
【図4】拡張メモリ間記述子を用いて遠隔メモリ・アクセスを必要とする動作を実行するための、本発明の1つの例示的な実施形態に従った機構の一例を示す例示的な図である。

【図5】本発明の1つの例示的な実施形態に従って拡張メモリ間記述子を発生させるための例示的なプロセスの概要を示すフローチャートである。

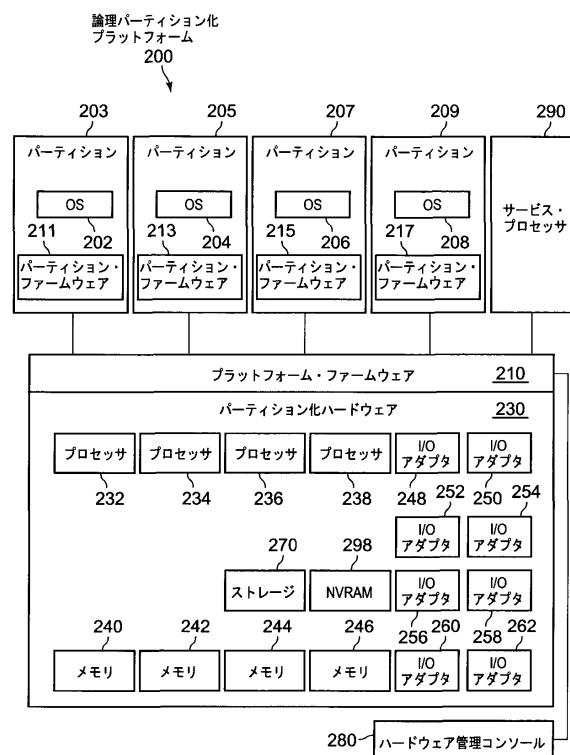
【図6】本発明の1つの例示的な実施形態に従って拡張メモリ間記述子を用いるための例示的なプロセスの概要を示すフローチャートである。

10

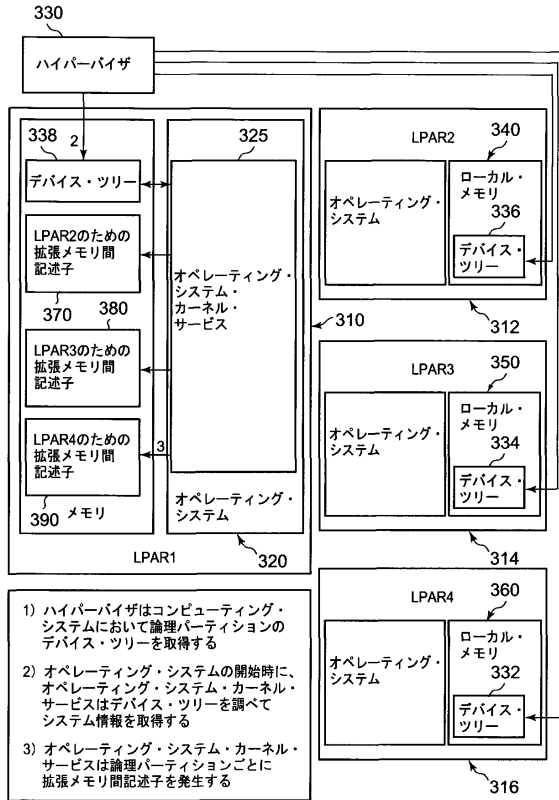
【図1】



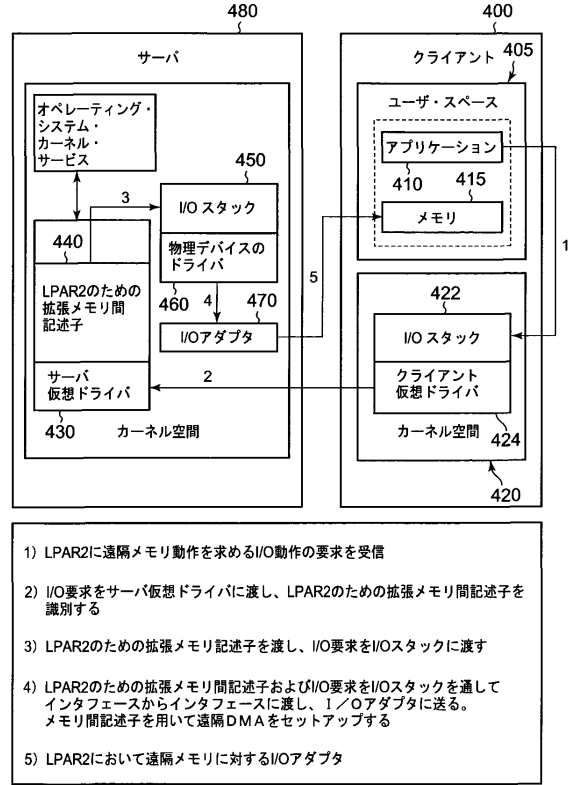
【図2】



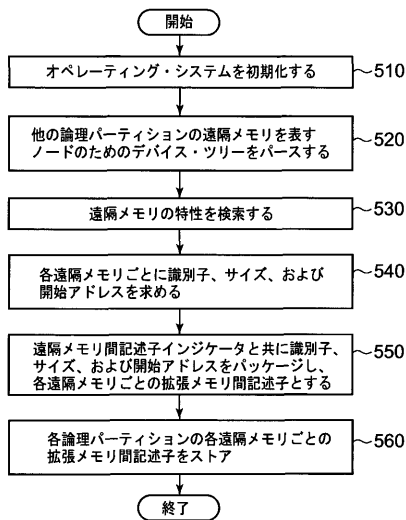
【図3】



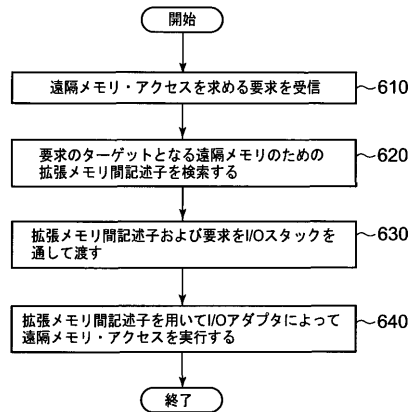
【図4】



【図5】



【図6】



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 アスロット、ヴィシヤル、チットランジャン

アメリカ合衆国78758 テキサス州オースティン グレーシー・ファームズ・レーン2600
アパートメント 1312

(72)発明者 メアリー、ブルース

アメリカ合衆国78759 テキサス州オースティン アミー・コート3914

審査官 鈴木 修治

(56)参考文献 特開2002-073358(JP,A)

特開2002-157133(JP,A)

特開2004-005443(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 9/46-9/54

G06F 12/10

G06F 13/28