

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4921018号  
(P4921018)

(45) 発行日 平成24年4月18日 (2012. 4. 18)

(24) 登録日 平成24年2月10日 (2012. 2. 10)

(51) Int. Cl.

F I

G O 6 F 12/10 (2006. 01)

G O 6 F 12/00 (2006. 01)

G O 6 F 9/46 (2006. 01)

G O 6 F 13/10 (2006. 01)

G O 6 F 12/10 5 5 9

G O 6 F 12/00 5 1 4 Z

G O 6 F 9/46 3 5 0

G O 6 F 12/00 5 9 9

G O 6 F 12/00 5 9 7 U

請求項の数 2 (全 21 頁) 最終頁に続く

(21) 出願番号 特願2006-102266 (P2006-102266)  
 (22) 出願日 平成18年4月3日 (2006. 4. 3)  
 (65) 公開番号 特開2006-294028 (P2006-294028A)  
 (43) 公開日 平成18年10月26日 (2006. 10. 26)  
 審査請求日 平成21年1月22日 (2009. 1. 22)  
 (31) 優先権主張番号 05102679.7  
 (32) 優先日 平成17年4月5日 (2005. 4. 5)  
 (33) 優先権主張国 欧州特許庁 (EP)

(73) 特許権者 390009531  
 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  
 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  
 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  
 (74) 代理人 100108501  
 弁理士 上野 剛史  
 (74) 代理人 100112690  
 弁理士 太佐 種一  
 (74) 代理人 100091568  
 弁理士 市位 嘉宏

最終頁に続く

(54) 【発明の名称】 直接実行機能を提供するためのシステム、コンピュータシステム、方法およびプログラム

(57) 【特許請求の範囲】

【請求項 1】

メモリアドレス指定デバイスに対する直接実行機能を自動的に提供する方法であって、アプリケーションプログラムによりファイルにアクセスするための要求を受理するステップと、

前記ファイルを保持するファイルシステムを識別するステップと、

前記ファイルシステムを管理するファイルシステムドライバが、前記メモリアドレス指定デバイスにあるファイルの内容のシステムメモリアドレスを検索するためのファイルシステムダイレクトアクセスインターフェースを提供するかどうかを決定するステップと、

前記ファイルシステムドライバが、前記ファイルシステムダイレクトアクセスインターフェースを提供しない場合には前記直接実行機能を使用せず、前記ファイルシステムドライバが、前記ファイルシステムダイレクトアクセスインターフェースを提供する場合には前記ファイルシステムが存在するデバイスを識別するステップと、

前記デバイスを管理するデバイスドライバが、前記メモリアドレス指定デバイスの指定されたブロックのシステムメモリアドレスを検索するためのデバイスダイレクトアクセスインターフェースを提供するかどうかを決定するステップと、

前記デバイスドライバが、前記デバイスダイレクトアクセスインターフェースを提供しない場合には前記直接実行機能を使用せず、前記デバイスドライバが、前記デバイスダイレクトアクセスインターフェースを提供する場合には前記ファイルシステムが、前記直接実行機能を可能にするように構成されるかどうかを決定するステップと、

10

20

前記ファイルシステムが、前記直接実行機能を可能にするように構成されると決定することに対応して、前記ファイルシステムダイレクトアクセスインターフェースを使用して、前記直接実行機能を提供するステップと、を含む方法。

【請求項 2】

ディジタルコンピュータの内部メモリ内に格納されたコンピュータプログラムであって、請求項 1 に記載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】

10

【0001】

本発明は、全体としてオペレーティングシステム、特に、直接実行機能をインプリメントするためのシステムおよび方法を提供するオペレーティングシステムに関する。

【背景技術】

【0002】

先行技術のコンピュータシステムは、プログラムおよびデータファイルを保持するための不揮発性大容量記憶装置（たとえば、ハードディスクドライブ）を含む。これらのファイルの内容は、RAM（「ランダムアクセスメモリ」）タイプのシステムメモリ内にロードされ、CPU（「中央処理装置」）によりアクセスまたは実行される。この動作は、一般に、アプリケーションプログラムに代わってオペレーティングシステムにより実行される。先行技術のコンピュータシステムおよびオペレーティングシステムは、仮想メモリおよびデマンドページングをサポートしている。アプリケーションは、システムメモリアドレスを直接使用して、アプリケーションが使用するコードおよびデータを指示するのではなく、代わりに「仮想アドレス」を使用して記憶場所を指示し、記憶場所は、CPU回路によりインプリメントされ、オペレーティングシステムにより制御されるページング機構により、システムメモリアドレスに変換される。その結果、オペレーティングシステムは、プログラムおよびデータファイル全体をRAM内にロードする必要性をなくすることができる。むしろ、システムメモリは特定サイズの塊（「ページ」と呼ばれる）に分割され、オペレーティングシステムは、この特定のページにアクセスした時に、ファイル内容の対応する塊を各々のメモリページ内にロードする。このプロセスは、通常、「デマンドページング」と呼ばれる。

20

30

【0003】

この方法の不利な点の 1 つに、RAM がプログラムおよびデータファイルの内容を保持する必要があり、他の目的に使用可能な RAM の量が減少することがある。さらに、一般に、場合に応じて内容を RAM 内にダウンロードする必要がある。したがって、先行技術のコンピュータシステムによっては、RAM と同様に、CPU が直接アクセスできる別のタイプの不揮発性記憶装置（「メモリアドレス指定デバイス」）が提供されている。メモリアドレス指定デバイスの先行技術の一実施態様は、フラッシュメモリカードである。メモリアドレス指定デバイスを用いれば、CPU は、最初に RAM 内に内容をダウンロードすることなく、メモリアドレス指定デバイス上に格納されたコードおよびアクセスデータを直接実行できるようになる。メモリアドレス指定デバイス上に存在するコードを直接実行する方法は、「直接実行」と呼ばれる。仮想メモリをサポートするオペレーティングシステム上で実行されるアプリケーションに直接実行機能を提供するには、オペレーティングシステムは、アプリケーションのアドレス空間の特定の仮想アドレスが、メモリアドレス指定デバイスによりサポートされるアドレスの範囲内のシステムメモリアドレスにマップされるように、ページング機構を制御する必要がある。

40

【0004】

先行技術のその他のコンピュータシステムは、仮想化機能を提供する。仮想化は、単一のコンピュータシステムを実行する「ハイパーバイザ」と呼ばれることの多いソフトウェアプログラムによってインプリメントされるが、複数の「ゲスト」オペレーティングシス

50

テムが各々別個の「仮想マシン」で同時に実行をすることができるようにする。オペレーティングシステムから各々の仮想マシンを見ると、仮想マシン自体が、CPU、RAMおよびI/Oデバイスを備えた実際のコンピュータシステムとして内部で動作しているように見える。これらの仮想コンポーネントに対するアクセスは、ハイパーバイザによって傍受され、実際のコンポーネントに対するアクセスに変換される。これによって、コンピュータシステムのリソースを複数のゲストオペレーティングシステム間で共有し、システムリソース全体の利用率を増加させることができる。

#### 【0005】

いくつかの先行技術の仮想化コンピュータシステムの不利な点の1つとして、同じプログラムまたはデータが同じハイパーバイザ下で動作する複数のゲストにより同時にアクセスされることと、各々のゲストオペレーティングシステムは、それらの内容を保持するための仮想RAMを別個に割り当て、ひいては、ハイパーバイザが、前記内容の複数の同じコピーを物理RAMに割り当てる必要性が出てくる点が挙げられる。これは、他の目的に使用することのできるメモリが減少することを意味し、これによって効率的に同時に動作可能なゲストの数が制限されてしまう。したがって、先行技術のハイパーバイザによっては、複数のゲストから同時にアクセスできる物理メモリのセグメント（「共有メモリセグメント」）が提供されている。プログラムまたはデータファイルを共有メモリセグメント内に格納することにより、複数のゲストが、前記ファイルに同時にアクセスすることができ、最初に内容を仮想RAM内にダウンロードする必要はなくなる。ゲストオペレーティングシステムから前記共有セグメントを見ると、物理メモリアドレス指定デバイスであるかのように見える。

#### 【0006】

データおよびプログラムファイルは、一般に、標準のファイルシステムレイアウトを使用するデバイス上に格納され、オペレーティングシステムによっては、異なる利用シナリオのために最適化された複数の異なるファイルシステムレイアウトを使用することができる。これをできるようにするため、先行技術のオペレーティングシステムは、一般に、複数のコンポーネント状に構造化される。オペレーティングシステムによっては、中央ファイルおよびメモリ管理コンポーネント、複数のファイルシステムドライバ、並びに複数のI/Oデバイスドライバが存在する。したがって、オペレーティングシステムは、ファイルおよびメモリ管理コンポーネントと組み合わせて、ファイルシステムドライバおよびI/Oデバイスドライバの適切な対を使用することにより、サポートされる何れかのI/Oデバイス上で、サポートされる何れかのファイルシステムレイアウトを使用できるようにする。しかし、先行技術のオペレーティングシステムは、既存のファイルシステムドライバを使用して、直接実行をできるようにする形態でメモリアドレス指定デバイスにアクセスすることはできない。メモリアドレス指定デバイスに対するアクセスをサポートする直接実行は、モノリシック形態でインプリメントされる。

#### 【0007】

実際、先行技術のオペレーティングシステムのインプリメンテーションによっては、標準ファイルシステムレイアウトを使用して、メモリアドレス指定デバイス上にデータを格納することができず、むしろ、このようなデバイスに特有の方法で、これらのデバイス上のデータを配列する必要がある。この配列には、特にI/Oデバイスおよびメモリアドレス指定デバイスの両方を使用するコンピュータシステムの場合、不利な点が複数ある。異なるファイルシステムレイアウトをサポートすることにより、システム管理はさらに難しくなる可能性がある。異なるレイアウトをフォーマット、管理、バックアップおよび復元するには、異なるツールが必要になる。既存のファイルの集合をI/Oデバイスからメモリアドレス指定デバイス、またはこの逆に移行することはさらに難しくなる。メモリアドレス指定デバイスにより要求される特定のレイアウトは、存在するすべての特徴（たとえば、アクセス制御および特権チェックをインプリメントする）に、標準のファイルシステムレイアウトを提供するわけではない。

#### 【0008】

10

20

30

40

50

もう1つの先行技術のインプリメンテーション(z Series (IBM Corporation)の商標)上のLinux(Linus Torvaldsの米国およびその他の国における商標、以下、同様なので略)用XIP2FSファイルシステム)は、Linuxオペレーティングシステムが提供する標準ファイルシステムフォーマットの1つである第2拡張ファイルシステム(「ext2」)を使用して、プログラムおよびデータを仮想メモリアドレス指定デバイス上(z/VM(IBM Corporationの商標)ハイパーバイザにより提供される共有メモリセグメント)に格納するためのサポートを提供する。しかし、この方法には、上記の段落に記載した不利な点の殆どがあり、他の標準Linuxファイルシステムフォーマットを使用することができないほか、XIP2FSは、ext2のすべての特徴を提供するわけではない(たとえば、XIP2FSは、書き込みアクセスをサポートしない)。

10

#### 【0009】

XIP2FSのもう1つの不利な点は、XIP2FSがオペレーティングシステムの上記のコンポーネント構造に組み込まれないことであり、XIP2FSが、ext2ファイルシステムレイアウトを使用してファイルにアクセスした場合も、XIP2FSは、Linux ext2ファイルシステムドライバを使用してアクセスするのではなく、代わりにext2ファイルシステムレイアウトにアクセスするのに必要なアクセスロジックを再インプリメントする。この場合もやはり、完全なext2ロジックの一部のみが再インプリメントされるため、XIP2FSはext2のすべての特徴をサポートするわけではない。さらに他の不利な点として、Linuxオペレーティングシステムの標準ext2ファイルシステムコンポーネントは、長期にわたって開発され、新しい特徴が追加され、たとえば、Linuxカーネルバージョン2.6が提供されたext2ファイルシステムドライバのバージョンは、非常に大きいディレクトリ構造、およびより精巧なアクセス制御機構に対してより迅速にアクセスするためのサポートを追加した。XIP2FSは、ext2ドライバのこのような強化から自動的に効果を得るのではなく、必要なすべての特徴をXIP2FSコード内に再インプリメントする必要がある。

20

#### 【発明の開示】

#### 【発明が解決しようとする課題】

#### 【0010】

本発明の目的は、オペレーティングシステムによる直接実行機能を提供して、上記の先行技術の不利な点を解消する方法を提供することである。

30

#### 【課題を解決するための手段】

#### 【0011】

本発明は、直接実行機能をインプリメントするための新しいシステムおよび方法を提供するオペレーティングシステムを開示する。

#### 【0012】

本発明の基礎となる先行技術のオペレーティングシステムは、アプリケーションプログラムに対するインターフェースを有するメモリ/ファイルマネージャと、メモリ/ファイルマネージャに対するファイルシステムI/Oインターフェースを有する少なくとも1個のファイルシステムドライバと、ファイルシステムドライバに対するデバイスI/Oインターフェースを有する少なくとも1個のデバイスドライバとを含み、少なくとも1個のデバイスドライバは、少なくとも1個のI/Oベースデバイス、およびファイルシステムドライバに対するデバイスI/Oインターフェースを有する少なくとも1個のデバイスドライバに対するアクセスを提供し、少なくとも1個のデバイスドライバは、少なくとも1個のメモリアドレス指定デバイスに対するアクセスを提供し、オペレーティングシステムは、少なくとも1個のメモリアドレス指定デバイスにアクセスするための直接実行機能を提供する。

40

#### 【0013】

先行技術のオペレーティングシステムは、直接実行機能をインプリメントするための以下の新しい独創的な機能コンポーネント：

50

メモリ/ファイルマネージャと少なくとも１個のファイルシステムドライバとの間のファイルシステムのダイレクトアクセスインターフェースであって、ファイルシステムのダイレクトアクセスインターフェースが、指定のオフセットにおけるファイルの内容のシステムメモリアドレスを検索する機能を提供し、ファイルがメモリアドレス指定デバイス上に存在するインターフェースと、

少なくとも１個のファイルシステムドライバと、少なくとも１個のメモリアドレス指定デバイスに対するアクセスを提供する少なくとも１個のデバイスドライバとの間のデバイスダイレクトアクセスインターフェースであって、デバイスダイレクトアクセスインターフェースが、少なくとも１個のメモリアドレス指定デバイスの指定ブロックのシステムメモリアドレスを検索する機能を提供するインターフェースとにより拡張され、

10

直接実行機能は、メモリ/ファイルマネージャ、少なくとも１個のファイルシステムドライバ、ファイルシステムのダイレクトアクセスインターフェースおよびデバイスダイレクトアクセスインターフェースを使用することにより、少なくとも１個のメモリアドレス指定デバイスに対するアクセスを提供する少なくとも１個のデバイスドライバにより提供される。

#### 【 0 0 1 4 】

本発明の上記および追加の目的、特徴および効果は、以下の詳細な説明で明らかになる。

#### 【 0 0 1 5 】

本発明の新たな特徴は、添付の請求の範囲に記載する。しかし、本発明自体、および好ましい使用形態、他の目的、および効果は、具体的な実施態様に関する以下の詳細な説明を参照し、添付の図面に関連して読むと最も良く理解されるであろう。

20

#### 【発明を実施するための最良の形態】

#### 【 0 0 1 6 】

図１は、コンピュータシステム１０のブロック図を示す。コンピュータシステム１０は、パーソナルコンピュータ、メインフレームコンピュータ、または他の何らかのタイプのコンピュータもしくはデータ処理システムで良い。コンピュータシステム１０は、以下に述べるように、別のコンピュータシステム上で動作するハイパーバイザにより提供される仮想マシンでも良い。コンピュータシステム１０は、中央演算処理装置（「ＣＰＵ」）１１と、ランダム・アクセスメモリ（「ＲＡＭ」）１２と、メモリアドレス指定デバイス１３と、Ｉ／Ｏベースデバイス１４とを備える。一実施態様では、メモリアドレス指定デバイス１３はフラッシュメモリカードで良い。他の実施態様では、メモリアドレス指定デバイス１３は、ＣＰＵ１１がメモリ動作のために直接アクセスできる何らかのデバイスで良い。一実施態様では、Ｉ／Ｏベースデバイス１４はハードディスクドライブで良い。他の実施態様では、Ｉ／Ｏベースデバイス１４は、Ｉ／Ｏ操作を使用して、データをＲＡＭ１２からコピーするか、またはＲＡＭ１２にコピーできる何らかのデバイスで良い。コンピュータシステム１０は、メモリアドレス指定デバイス１３またはＩ／Ｏベースデバイス１４あるいはその両方の複数のインスタンスも含む場合がある。ＣＰＵ１１、ＲＡＭ１２、並びにデバイス１３および１４は、システムバス１５に結合される。ＣＰＵ１１は、メモリ動作のために、ＲＡＭ１２およびメモリアドレス指定デバイス１３に直接アクセスすることができる。ＣＰＵ１１は、メモリ動作のためにＩ／Ｏベースデバイス１４に直接アクセスすることはできないが、データは、Ｉ／Ｏ操作を使用して、デバイス１４からＲＡＭ１２に、あるいは逆にコピーすることができる。コンピュータシステム１０は、１つまたは複数のアプリケーションプログラムを実行させることが可能なオペレーティングシステムを実行させる（以下で詳細に説明する）。オペレーティングシステムは、コンピュータシステム１０の様々なリソース（ＣＰＵ１１、ＲＡＭ１２、デバイス１３および１４）に対するアプリケーションプログラムによるアクセスを管理および調整する。

30

40

#### 【 0 0 1 7 】

いくつかの実施態様では、コンピュータシステム１０は、別のコンピュータシステム上で動作するハイパーバイザによりエミュレートされる仮想マシンで良い。図２は、コンピ

50

ユータシステム10が、仮想CPU11、仮想RAM12、並びに仮想デバイス13および14から成るような一実施態様を示す。すべての仮想コンポーネントはハイパーバイザ21により提供され、ハイパーバイザ21は、別のコンピュータシステム20上で動作するソフトウェアプログラムであり、それ自体はCPU、RAMおよびデバイスから成る。ハイパーバイザ21は、完全にソフトウェアによって、またはコンピュータシステム20が提供する視覚化ハードウェア支援機能を使用することによって、仮想コンポーネント10~14を提供する。仮想マシン10のほかに、他の仮想マシン22が、コンピュータシステム20上のハイパーバイザ21下で同時に動作する。一実施態様では、ハイパーバイザ21は、IBM Corp. が製造販売するz/VM (IBM Corporationの商標)ソフトウェアプログラムであって、IBM eServer (IBM Corporationの商標) zSeries (IBM Corporationの商標)メインフレームコンピュータ上で動作するプログラムで良い。この実施態様では、コンピュータシステム10は、z/VMにより提供される仮想マシンで良く、メモリアドレス指定デバイス13は、z/VMに基づいて画定される非連続保管セグメント(「DCSS」)で良い。DCSSは、複数の仮想マシンが同時に使用可能なz/VMにより管理されるメモリのセグメントである。z/VMは、複数の仮想マシンがDCSSを同時に使用可能である場合でも、コンピュータシステム20の実RAM内におけるDCSSの内容の1つのコピーのみを保持する。

#### 【0018】

CPU11がメモリ動作のために直接アクセス可能なすべてのコンポーネントは、図3に示すように、コンピュータシステム10のシステムメモリアドレス空間30を構成する。RAM12およびメモリアドレス指定デバイス13は、システムメモリアドレス空間30の一部である。さらに、その他のコンポーネント(図示しない)は、システムメモリアドレス空間30、たとえば読出し専用メモリ(「ROM」)またはビデオカードフレームバッファである。CPU11により実行されるすべてのプログラムインストラクション、およびCPU11により実行されるインストラクションによりアクセスされるすべてのメモリデータは、実行が生じた時点で、システムメモリアドレス空間30内に存在しなければならない。使用可能なメモリの見掛けのサイズを増加し、同じコンピュータシステム上で同時に動作する異なるアプリケーションプログラムが、偶発的に互いに他のメモリを修正するのを防止するため、オペレーティングシステムは、各々のアプリケーションプログラムに「仮想アドレス空間」を提供し、仮想アドレス空間内のシステムメモリアドレス空間の選択部分に対するアクセスを提供する。CPU11は、アプリケーションコードを実行するが、アプリケーションの仮想アドレス空間内に存在するメモリにのみアクセスすることができる。仮想アドレス空間とシステムメモリアドレス空間との間で変換することができるよう、仮想アドレス空間とシステムメモリアドレス空間の両方が、一般に「ページ」と呼ばれる等サイズの塊状に分割される。

#### 【0019】

図3は、コンピュータシステム10のシステムメモリアドレス空間30を示す。さらに、アプリケーションプログラムの仮想アドレス空間31が示されている。仮想アドレス空間31でアドレス指定できる各々のページについては、当該ページの状態を実際に記述するページ記述子が存在しなければならない。ページは、ある時点でシステムメモリアドレス空間30内に存在しても、あるいは存在しなくても良い。ページがシステムメモリアドレス空間30内に存在する場合、ページ記述子はその事実を指示し、システムメモリアドレス空間30内のページの位置を指示することになる。ページがシステムメモリアドレス空間30内に存在しない場合、記述子はその事実を指示し、アプリケーションプログラムがこのページを介してアクセスしようとしている内容の場所をオペレーティングシステムが確認できるようにする追加の情報を含むことになる。仮想アドレス空間に対するすべてのページ記述子の集合は、ページテーブルと呼ばれる。仮想アドレス空間のページテーブルは、オペレーティングシステムにより維持される。

#### 【0020】

CPU 11 が仮想アドレス空間 31 でアプリケーションコードを実行している間、CPU 11 によりアクセスされるすべてのメモリアドレスは、仮想アドレス空間 31 用のページテーブルを使用して、仮想アドレス空間 31 内の仮想アドレスからシステムメモリアドレス空間 30 内のシステムメモリアドレスに変換される。この変換はページング機構により実行され、ページング機構は、ハードウェアもしくはソフトウェアインプリメンテーション、またはこれらの組合せで良い。実施態様によっては、ページング機構は、CPU 11 内に存在するページングユニットによりインプリメントされる。CPU 11 が、システムメモリアドレス空間 30 に対応するページを現在持たない仮想アドレス空間 31 内のページにアクセスする場合、ページング機構は、CPU 11 にページ障害割込みを生成させる。この割込みにより、「ページ障害ハンドラ」と呼ばれるオペレーティングシステムコードが、必要な内容をシステムメモリアドレス空間 30 内のいくつかのページ内にもたらし、ページテーブルを更新して、仮想ページをシステムメモリアドレス空間 30 内のこのページにマップする。次に、アプリケーションは、ページの実行およびアクセスを継続する。このプロセスは、一般に、「デマンドページング」と呼ばれる。

#### 【0021】

図 3 に示す実施例の状況では、仮想アドレス空間 31 は、現在 4 つのページを保持している。3 つのページ (32a ~ 32c) は、現在実行されているアプリケーションコードを含み、1 つのページ (32d) は、アプリケーションによりアクセスされるデータを含む。ページ 32a ~ 32c に含まれるアプリケーションコードは、ブロック 34a ~ 34c 内の I/O ベースデバイス 14 上に存在するアプリケーションプログラムファイルからロードされる。アプリケーションコードページ 32a および 32b は、システムメモリアドレス空間 30 内に現在実際に存在しており、RAM 12 のページ 33a および 33b 内に存在する。同様に、データページ 32d はシステムメモリアドレス空間内に存在しており、RAM 12 のページ 33d 内に存在する。アプリケーションコードページ 32c は、現在、システムメモリアドレス空間 30 内に存在しない。アプリケーションがページ 32c にアクセスすると、オペレーティングシステムのページ障害ハンドラは、RAM 12 に新しいページ 33c (図示しない) を割り当て、I/O 操作を実行してブロック 34c の内容をページ 33c にコピーし、ページテーブルを更新して、仮想アドレス空間 31 のページ 32c をシステムメモリアドレス空間 30 のページ 33c にマップする。ページ 32a および 32b の場合、このプロセスは、図 3 に示す時点で既に完了している。データページ 32d / 33d は、アプリケーションにより実行時に生成された内容を保持する。この内容は、I/O ベースデバイス 14 からロードされたものではない。

#### 【0022】

図 3 に示すように、CPU 11 が I/O ベースデバイス 14 上に存在するアプリケーションプログラムを実行している時、RAM 12 のページは、実行しているアプリケーションプログラムのコードを保持する必要がある。しかし、アプリケーションプログラムがメモリアドレス指定デバイス上に存在する場合、その必要はなく、RAM 12 のより多くのページを他の目的に使用できる (あるいは別法によると、より少ない RAM の全体量で、意図されたタスクを履行することができる)。このプロセスは、一般に「直接実行」と呼ばれる。図 4 は、仮想アドレス空間 31 で実行されている図 3 と同じアプリケーションを示すが、この場合、アプリケーションは、I/O ベースデバイス 14 ではなくメモリアドレス指定デバイス上に存在し、直接実行されている。図 4 に示すように、アプリケーションプログラムファイルからのアプリケーションコードは、メモリアドレス指定デバイス 13 のブロック 35a ~ c に存在する。メモリアドレス指定デバイス 13 は、システムメモリアドレス空間 30 内に存在するため、ブロック 35a および 35b は、それぞれ仮想アドレス空間 31 のページ 32a および 32b に直接マップされる。同様に、アプリケーションがページ 32c にアクセスすると、オペレーティングシステムのページ障害ハンドラは、単に、仮想アドレス空間 31 のページ 32c とシステムメモリアドレス空間 30 のページ 35c との間の別のマッピングを確立し、オペレーティングシステムは、どのページも RAM 12 に割り当てて必要はない。一方、データページ 32d は、図 3 に示されたシ

10

20

30

40

50

ナリオと同様に、RAM 12のページ33dにマップされる点に注目する。

【0023】

図5は、本発明を具現しない先行技術のオペレーティングシステム40のコンポーネント構造を視覚化する。図3および図4に示すデマンドページングおよび直接実行機能をインプリメントするために必要なコンポーネントのみを図示する。オペレーティングシステム40は、インターフェース48を介したアプリケーションプログラム49から発せられるI/Oベースデバイス14またはメモリアドレス指定デバイス13上に存在するファイルへのアクセス要求を処理するメモリ/ファイルマネージャ41を含む。I/Oベースデバイス14上のファイルにアクセスするために、メモリ/ファイルマネージャ41は、ファイルシステムI/Oインターフェース45を介して、ファイルシステムドライバ43と対話する。オペレーティングシステム40は、複数のバージョンのファイルシステムドライバを含み、各々のファイルシステムドライバは、特定のファイルシステムタイプの処理を担い、すべてのファイルシステムドライバは、同じファイルシステムI/Oインターフェース45をインプリメントする。同様に、オペレーティングシステム40は複数のバージョンのデバイスドライバを含み、各々のデバイスドライバは特定タイプのデバイスの処理を担い、すべてのデバイスドライバは同じデバイスI/Oインターフェース46をインプリメントする。

10

【0024】

図5に示すように、オペレーティングシステム40のコンポーネント構造内では、デバイスドライバ44（およびすべてのデバイスドライバ）が担っているのは、データをI/Oベースデバイス上の指定の場所からRAM内に、あるいはこの逆にコピーすることである。デバイスドライバは、デバイスの内容またはこれらの内容が組織化される方法に関する知識は持たない。デバイス上に存在するデータは、一般に、各々が「ブロック数」により識別される「ブロック」と呼ばれる塊状に分割される。デバイスI/Oインターフェース46は、デバイスドライバの要求により、ブロック数Bで識別されるブロックから、アドレスAで開始するRAMのブロックに、あるいはこの逆にデータをコピーできるようにする。

20

【0025】

デバイス上にデータを構造化ストレージで格納できるように、オペレーティングシステム40は、ファイルシステムドライバを提供する。ファイルシステムドライバは、「ファイル」の集合としてデバイス上に格納されたデータへのアクセスを可能にし、各々のファイルは、順に並んだバイトのシーケンスの抽象化を提供する。すべてのファイルは、一般に「ファイル名」と呼ばれるファイルシステムにより提供されるいくつかの手段により識別される。ファイルシステムは、各々のファイルのどのバイトが、下のデバイスのどのブロック内に格納されるかというトラックを保持する。この記録保持を実行するのに必要な情報は、通常、「ファイルシステムメタデータ」と呼ばれ、下にあるデバイス上に格納される。ファイルシステムメタデータの特定のレイアウトは、ファイルシステムごとに異なる。図5に示すように、オペレーティングシステム40のコンポーネント構造内では、ファイルシステムメタデータのレイアウトを処理するために必要なすべてのロジックをインプリメントするのは、ファイルシステムドライバ43（およびすべてのファイルシステムドライバ）が担っている。ファイルシステムI/Oインターフェース45は、あるファイル名で識別されるファイルFからデータを、ファイルF内の指定のオフセット0で開始して、アドレスAで開始するRAMのブロック内に、あるいはこの逆にコピーするというファイルシステムドライバの要求を認める。このような要求を処理するため、ファイルシステムドライバは、ファイルシステムメタデータを調査して、下にあるデバイスのどのブロックBが、ファイルF内のオフセット0に対応するデータを保持するかを決定し、下にあるデバイスを処理するデバイスドライバのデバイスI/Oインターフェース46を使用して、前記デバイスとメモリの指定ブロックとの間でコピーを行うかを決定する。

30

40

【0026】

したがって、図5に示すように、オペレーティングシステム40のコンポーネント構造

50



内では、図 3 に示すデマンドページング機能は以下のようにインプリメントされる：アプリケーションが、システムメモリアドレス空間内で現在使用されていないページにアクセスすると、オペレーティングシステムのページ障害ハンドラ（通常、メモリ/ファイルマネージャ 4 1 の一部）は、ページ記述子から、アプリケーションがページ中のどの内容を探しているのか決定する。ページ記述子は、ファイル F、および前記内容が発見されるファイル F 内のオフセット 0 を指定することにより、前記内容を識別する。次に、メモリ/ファイルマネージャ 4 1 は、新しいページを RAM 1 2 に割り当て、ファイルシステムドライバ 4 3 からファイルシステム I/O インターフェース 4 5 を介して、データをファイル F 内のオフセット 0 から前記の新しいページにコピーするように要求する。上記のとおり、ファイルシステムドライバ 4 3 は、ファイルシステムメタデータを調査して、下にあるデバイス（この場合、I/O ベースデバイス 1 4）のどのブロック B がデータを保持しているかを決定し、デバイスドライバ 4 4 のデバイス I/O インターフェース 4 6 を使用して、そのデータを前記の新しいページにコピーする。I/O 操作が完了すると、メモリ/ファイルマネージャ 4 1 は、それに応じてページテーブルを更新する。

【 0 0 2 7 】

しかし、先行技術のファイルシステム I/O インターフェース 4 5 は、この目的に適していないため、ファイルシステムドライバ 4 3 を使用して、メモリアドレス指定デバイス 1 3 上に存在するファイルに対する直接実行アクセスを実行するのは不可能である。その代わりに、図 5 に示すように、オペレーティングシステム 4 0 のコンポーネント構造内では、直接実行アクセスは、メモリ/ファイルマネージャ 4 1 内にしっかりと組み込まれている XIP マネージャ 4 2 により処理され、メモリアドレス指定デバイス 1 3 に直接アクセスする。したがって、XIP マネージャ 4 2 は、メモリアドレス指定デバイス 1 3 にアクセスし、前記デバイス上に格納されたデータのファイルシステムレイアウトを処理することの両方を担っている。XIP マネージャ 4 2 によりサポートされるファイルシステムレイアウト以外のファイルシステムレイアウトを使用して、メモリアドレス指定デバイス 1 3 上のデータにアクセスすることは、オペレーティングシステム 4 0 が、さもなければ、このようなファイルシステムにファイルシステムドライバを提供すると思われる場合でも不可能である。

【 0 0 2 8 】

本発明は、このような制約を除去する。図 6 は、本発明を具現するメモリアドレス指定デバイス 1 3 に対する直接実行アクセスをインプリメントするオペレーティングシステム 5 0 のコンポーネント構造を示す。図 5 に示す先行技術のオペレーティングシステム 4 0 と比べて、オペレーティングシステム 5 0 は、コンピュータシステム 1 0 のすべてのハードウェアコンポーネント（RAM 1 2、I/O ベースデバイス 1 4、メモリアドレス指定デバイス 1 3）に対して同じインターフェースを保持する。また、オペレーティングシステム 5 0 は、アプリケーションプログラム 4 9 に対して同じインターフェース 4 8 を使用する。メモリ/ファイルマネージャ 5 1 は、先行技術（図 5 参照）により提供される修正バージョンのメモリ/ファイルマネージャ 4 1 であり、ファイルシステムドライバ 5 2 は、先行技術（図 5 参照）により提供される修正バージョンのファイルシステムドライバ 4 3 である。オペレーティングシステム 5 0 は、メモリアドレス指定デバイス 1 3 にアクセスするデバイスドライバ 5 3 も提供する。オペレーティングシステム 4 0 の先行技術のいくつかのバージョンは、同様にメモリアドレス指定デバイス 1 3 にアクセスするデバイスドライバも提供していたが、このようなドライバは、デバイス I/O インターフェース 4 6（図 5 参照）のみをインプリメントし、したがって、直接実行機能を提供することはできなかった点に注目する。この機能は、XIP マネージャ 4 2 によりインプリメントされたが、XIP マネージャ 4 2 は、デバイスドライバを使用せずに、メモリアドレス指定デバイス 1 3 に直接アクセスする。図 6 に示すように、本発明は、XIP マネージャのコンポーネントを必要としない。その代わりに、直接実行機能は、この場合、メモリ/ファイルマネージャ 5 1、ファイルシステムドライバ 5 2 およびデバイスドライバ 5 3 内に組み込まれる。これは、2 つの新しいインターフェース：ファイルシステムドライバ 5 2 により提供される

ファイルシステムのダイレクトアクセスインターフェース 54、およびデバイスドライバ 53 により提供されるデバイスダイレクトアクセスインターフェース 55 を使用することにより可能である。デバイスダイレクトアクセスインターフェース 55 の基本的な特徴は、システムメモリアドレス空間内に存在するメモリアドレス指定デバイスのブロック B のシステムメモリアドレス A を検索する手段を提供することである。同様に、ファイルシステムのダイレクトアクセスインターフェース 54 の基本的な特徴は、ファイル F が、システムメモリアドレス空間内に存在するメモリアドレス指定デバイス上に存在する場合、オフセット 0 におけるファイル F の内容のシステムメモリアドレス A を検索する手段を提供することである。これらのダイレクトアクセスインターフェースの使用については、以下に詳細に説明する。

10

#### 【0029】

図 6 に示すように、オペレーティングシステム 50 は、インターフェース 48 を通してアプリケーションプログラム 49 と対話する。インターフェース 48 の一部は、アプリケーションプログラム 49 によりトリガされるデマンドページング要求から成り、アプリケーションプログラム 49 は、システムメモリアドレス空間 30 のページに現在マップされない仮想アドレス空間 31 のページにアクセスする。(仮想アドレス空間 31 およびシステムメモリアドレス空間 30 は、図 3 および図 4 に示され、仮想アドレス空間 31 は、アプリケーション 49 の仮想アドレス空間に対応する点に注目する。) インターフェース 48 のその他の部分は、アプリケーションプログラム 49 によるオペレーティングシステム 50 に対する要求(「システムコール」)であって、I/O ベースデバイス 14 またはメモリアドレス指定デバイス 13 上に存在するファイルの内容を読み取るか、書き込むか、さもなければアクセスするための要求から成る。メモリ/ファイルマネージャ 51 は、インターフェース 48 を介したアプリケーションプログラム 49 による要求を処理するオペレーティングシステム 50 のコンポーネントである。I/O ベースデバイス 14 またはメモリアドレス指定デバイス 13 上に存在するファイルの内容にアクセスするには、メモリ/ファイルマネージャ 51 は、ファイルシステム I/O インターフェース 45 またはファイルシステムのダイレクトアクセスインターフェース 54 あるいはその両方を介して、ファイルシステムドライバ 52 と対話する。オペレーティングシステム 50 は、各々が特定のファイルシステムレイアウトの処理を担う複数のバージョンのファイルシステムドライバを含む。すべてのファイルシステムドライバは同じファイルシステム I/O インターフェース 45 をインプリメントし、いくつかのファイルシステムドライバは、さらに、メモリアドレス指定デバイス上に存在するファイルに対する直接実行アクセスを実行するのに適するファイルシステムのダイレクトアクセスインターフェース 54 をインプリメントする。ファイルシステムドライバ 52 は、デバイス I/O インターフェース 46 またはデバイスダイレクトアクセスインターフェース 55 あるいはその両方を介して、デバイスドライバ 44 および 53 と対話する。オペレーティングシステム 50 は、各々が特定タイプのデバイスの処理を担う複数のバージョンのデバイスドライバも含む。すべてのデバイスドライバは同じデバイス I/O インターフェース 46 をインプリメントし、メモリアドレス指定デバイス 13 のデバイスドライバは、さらに、メモリアドレス指定デバイス 13 上に存在するファイルに対する直接実行アクセスを実行するのに適するデバイスダイレクトアクセスインターフェース 55 をインプリメントする。I/O インターフェース 45 およびダイレクトアクセスインターフェース 54 の両方を提供するファイルシステムドライバの場合、メモリ/ファイルマネージャ 51 は、このようなファイルシステムドライバにより処理されるファイルに対する定期的なアクセスおよび直接実行アクセスの両方を実行すること可能である。図 7 および図 8 は、それぞれ定期的なアクセスおよび直接実行アクセスを実行するのに必要な制御フローを詳細に示す。図 9 は、2 つの方法のどちらを特定ファイルのアクセスに使用するかを選択するために必要な決定ロジックを詳細に示す。

20

30

40

#### 【0030】

図 7 は、オペレーティングシステム 50 が、I/O ベースデバイス 14 上に存在するファイルにアクセスするための要求を処理する時の制御フローを示す。この制御フローは、

50

先行技術のオペレーティングシステム 40 が対応するアクセスに使用する制御フローと同じである点に注目する。制御フローのステップは、順にアクション 60 a ~ f により表現される。ここに示す特定のアクションは、アプリケーション 49 が、仮想アドレス空間 31 のページ 32 c にアクセスしようと試みた後、図 3 に示す状況に対応する。このページは、システムメモリアドレス空間 30 のどのページにも現在マップされていないため、ページ障害割込みが生成され、オペレーティングシステム 50 のメモリ/ファイルマネージャ 51 コンポーネントにより処理される。これは、図 7 にアクション 60 a で表現される。メモリ/ファイルマネージャ 51 は、ページ 32 c のページ記述子から、アプリケーションが、ページ 32 c の内容が、オフセット 0 におけるファイル F の内容に対応するはずであると予想していることを読み取る。メモリ/ファイルマネージャ 51 は、ファイル F が、ファイルシステムドライバ 52 により処理されるファイルシステム上に存在すると決定する。メモリ/ファイルマネージャ 51 は、さらに、ファイル F が直接実行をサポートしないと決定する（この点について、以下で詳細に説明する）。メモリ/ファイルマネージャ 51 は、次に、RAM 13 内の自由ページ 33 c を割り当て、そのシステムメモリアドレス A を決定し、ファイルシステムドライバ 52 のファイルシステム I/O インターフェース 45 を介して、オフセット 0 におけるファイル F の内容をシステムメモリアドレス A のページにコピーするように要求する（アクション 60 b）。ファイルシステムドライバ 52 は、オフセット 0 におけるファイル F のデータが、I/O ベースデバイス 14 のブロック B 上に存在し、デバイスドライバ 44 が、I/O ベースデバイス 14 へのアクセスを担うことを決定する。ファイルシステムドライバ 52 は、次に、デバイスドライバ 44 のデバイス I/O インターフェース 46 を介して、I/O ベースデバイス 14 のブロック B をシステムメモリアドレス A のページにコピーするように要求する（アクション 60 c）。デバイスドライバ 44 は、I/O ベースデバイス 44 上の I/O 操作 61 に、そのコピーを実行させる。I/O 操作 61 が完了すると、デバイスドライバ 44 は、デバイス I/O インターフェース 46 を介して、ファイルシステムドライバ 52 に完了を報告し（アクション 60 d）、ファイルシステムドライバ 52 は、同様に、ファイルシステム I/O インターフェース 45 を介してメモリ/ファイルマネージャ 51 に完了を報告する（アクション 60 e）。メモリ/ファイルマネージャ 51 は、最後に、システムメモリアドレス空間 30 におけるページ 33 c に対する仮想アドレス空間 31 のページ 32 c のマッピングを確立する。

#### 【0031】

図 8 は、同様に、オペレーティングシステム 50 が、メモリアドレス指定デバイス 13 上に存在するファイルに対する直接実行アクセスを実行するという要求を処理する時の制御フローを示す。このフローは、先行技術のオペレーティングシステム 40 が対応するアクセスに使用するフローとは異なり、本発明により記述される新しいダイレクトアクセスインターフェースを使用する点に注目する。さらに、メモリアドレス指定デバイス 13 上のファイルに対するいくつかのアクセスは、直接実行アクセスに適さず、この場合、その代わりに、図 7 に示し、上記で説明した制御フローと等価な制御フローを使用して、ファイルに対する定期的な I/O アクセスが実行される。これは、デバイスドライバ 53 が、デバイスドライバ 44 と同様にデバイス I/O インターフェース 46 もインプリメントするために可能である。

#### 【0032】

図 8 に示す制御フローのステップは、順にアクション 70 a ~ f で表現される。図 8 に示す特定のアクションは、アプリケーション 49 が仮想アドレス空間 31 のページにアクセスすることを試みた後の図 4 に示す状況に対応する。図 7 のように、ページ障害割込みが生成されて、メモリ/ファイルマネージャ 51 により処理され（アクション 70 a）、メモリ/ファイルマネージャ 51 は、やはり、ページ 32 c のページ記述子から、アプリケーションが、ページ 32 c の内容がオフセット 0 におけるファイル F の内容に対応するはずであると予想していることを読み取る。メモリ/ファイルマネージャ 51 は、ファイル F が、ファイルシステムドライバ 52 により処理されるファイルシステム上に存在すると決定する。メモリ/ファイルマネージャ 51 は、さらに、ファイル F が直接実行をサポートしな

いと決定する（この点について、以下で詳細に説明する）。メモリ/ファイルマネージャ 51 は、次に、ファイルシステムドライバ 52 のファイルシステムのダイレクトアクセスインターフェースを介して、オフセット 0 におけるファイル F の内容のシステムメモリアドレスを検索するように要求する（アクション 70b）。ファイルシステムドライバ 52 は、オフセット 0 におけるファイル F の内容が、メモリアドレス指定デバイス 13 のブロック B 上に存在し、デバイスドライバ 53 がメモリアドレス指定デバイス 13 へのアクセスを担うと決定する。ファイルシステムドライバ 52 は、次に、デバイスドライバ 53 のデバイスダイレクトアクセスインターフェース 55 を介して、メモリアドレス指定デバイス 13 のブロック B のシステムメモリアドレスを検索するように要求する（アクション 70c）。デバイスドライバ 53 は、メモリアドレス指定デバイス 13 のブロック B が、システムメモリアドレス空間 30 のページ 35c に対応し、デバイスダイレクトアクセスインターフェース 55 を介して、そのアドレス A をファイルシステムドライバ 52 に返し（アクション 70d）、同様に、ファイルシステムのダイレクトアクセスインターフェース 54 を介して、アドレス A をメモリ/ファイルマネージャ 51 に返す（アクション 70e）と決定する。メモリ/ファイルマネージャ 51 は、最後に、システムメモリアドレス空間 30 のアドレス A におけるページ 35c に対する仮想アドレス空間 31 のページ 32c のマッピングを確立する。

#### 【0033】

図 9 は、ファイルシステム I/O インターフェースまたはファイルシステムのダイレクトアクセスインターフェースを使用してファイル F にアクセスするかどうかを決定する時のオペレーティングシステム 50（図 6）内の制御フローを示す。オペレーティングシステムは、最初に、どのファイルシステムドライバが、ファイル F を保持するファイルシステム FS に責任を負うかを決定する。ファイルシステムドライバがファイルシステムのダイレクトアクセスインターフェースを全然サポートしていない場合、ファイルシステム I/O インターフェースが使用される。さもなければ、オペレーティングシステムは、どのデバイスドライバが、ファイルシステム FS の下にあるデバイス D に責任を負うかを決定する。デバイスドライバが、デバイスダイレクトアクセスインターフェースを全然サポートしていない場合、同様にファイルシステム I/O インターフェースが使用される。さもなければ、オペレーティングシステムは、ファイルシステム FS が、ファイル F に直接実行アクセスするためにユーザにより構成されたかどうかを決定する。構成された場合は、ファイルシステムのダイレクトアクセスインターフェースが使用され、さもなければ、ファイルシステム I/O インターフェースが使用される。実施態様によっては、ユーザは、FS 上のすべてのファイルに対する直接実行アクセスをできるようにするか、または FS 上のすべてのファイルに対して直接実行アクセスできないようにするかのみ選択する。他の実施態様では、ユーザは、各々の単一ファイルに関して別個に、このような選択を行う。さらに、実施態様によっては、その他の構成可能なファイルシステムパラメータが、ユーザにより、直接実行と互換性のある値に設定された場合にのみ、直接実行アクセスを可能にすることができ、たとえば、ファイルシステムのブロックサイズは、複数のシステムメモリページのサイズに等しい値、または複数のシステムメモリページのサイズに構成する必要がある。

#### 【0034】

オペレーティングシステムは、必ずしも、ファイル F に 1 回アクセスするごとに、図 9 に示す完全な決定ロジックを実行する必要はない点に注目する。その代わりに、選択は、ファイル F が最初にアクセスされ、ファイル F に関連するオペレーティングシステムのデータ構造内に記憶されている時点で選択され、その後のアクセスは、前記データ構造内に格納されている決定結果を再利用する。

#### 【0035】

図 10～図 18 は、Linux オペレーティングシステム内における本発明の一実施態様をさらに詳細に示す。本発明は、インターフェースを拡張することにより、直接実行機能を標準オペレーティングシステムの I/O コンポーネント構造内に組み込むのであって

、コンポーネント構造全体を交換するのではない。Linuxは、I/O操作の実行に関連して、以下のコンポーネント層を提供する：

関連する特定ハードウェアを認識せずに、外部記憶装置にアクセスできるようにするデバイスドライバ層と、

外部記憶装置を認識せずに、ロジックファイルのビューを使用して外部記憶装置にアクセスすることが可能なファイルシステム層と、

アプリケーションが仮想メモリ、およびオペレーティングシステムが使用する動的アドレス変換技術を認識せずに、アプリケーションの実行をできるようにするメモリ管理層。

【0036】

図10は、多くの現代的なオペレーティングシステムにインプリメントされるデバイスの抽象化を示す。この実施例は、Linuxのデバイスドライバを示す。make\_request関数は、データをデバイスに読み書きするための要求を提出するために呼び出される。request\_queueおよびdo\_request関数は、Linux内部の最適化の一部である。デバイスドライバは、「要求をディスパッチする -> 要求を処理する -> 割り込みハンドラ -> 要求をディスパッチする」のループで動作し、デバイスドライバに提出されたすべての作業が実行される。

【0037】

デバイスドライバ層は、読み取り要求または書き込み要求を提出できるようにする。この層の頻繁なユーザは、データをファイルから転送するか、またはファイルに転送するファイルシステムのread\_page(s)/write\_page(s)動作である。データをアドレス指定するには、物理ブロック数が使用される。

【0038】

本発明は、図11に示されているデバイスドライバ層のインターフェースに対する拡張を提供する。既存のインターフェースを元の状態に維持しつつ、この拡張は、デバイス上のデータを直接参照するために使用できる機能を提供する。この参照は、要求を提出せずにデバイス上のデータにアクセスし、その完了を待つために使用することができる。この拡張は、任意に、メモリアドレス指定デバイスにアクセスするデバイスドライバによりインプリメントすることが可能である。新しい動作direct\_accessは、make\_requestに類似する物理ブロック番号を取得するが、どのデータも転送しない。その代わりに、データの参照が返される。この参照は、デバイスが、さらにデバイスドライバ層と対話せずに、再び閉鎖/アンマウントされるまで、何らかのデータをその後いつでも物理ブロックに読み書きするために使用することができる。新しいインターフェースの使用は任意であり、従来のmake\_requestインターフェースは、生デバイスドライバなどの新しいインターフェースをサポートしていないユーザのために、元の状態を維持する。従来のインターフェースは汎用ファイルシステムをサポートすることが重要であり、汎用ファイルシステムを使用して、ファイルシステムのメタデータを転送する(inode、directoryエントリなど)。

【0039】

図12は、Linux内の汎用ファイルシステムが、デバイスドライバのmake\_request関数を使用して、どのようにアドレス空間操作をインプリメントするかを示す。read\_page(s)/write\_page(s)関数は、get\_block関数を使用して[複数のページを処理する場合、繰り返し]、複数の対象ページに関連するデバイス上の物理ブロック数を識別する。図10に示すmake\_request関数は、データにアクセスするために使用される。Linuxでは、ファイルシステムは、一般に、ファイルシステムのライブラリ関数と共にアドレス空間操作を使用して、sys\_read()およびsys\_write()などのファイル操作を実行する。これらのアドレス空間操作およびライブラリ関数の使用は任意だが、殆どの汎用ファイルシステムはこれらを使用する。アドレス空間操作、read\_page()、read\_pages()、write\_page()、write\_pages()は、データの1つまたは複数のメモリページをデバイスドライバから読み取るか、またはデバイスドライバに書

き込むために使用される。メモリページをアドレス指定するには、ロジックファイルハンドルおよびオフセットが使用される。このアドレス指定は、ファイルシステムにより物理ブロック番号に変換される。

【0040】

本発明は、図13に示すように、あるメモリページの背後で記憶装置の参照を検索できるようにする `get_xip_page` という名称の関数により、アドレス空間操作インターフェースに対する拡張を提供する。この機能は、ファイルハンドル、およびアドレス指定用オフセットを使用し、`get_block` 関数を呼び出すことにより、ファイルハンドルおよびアドレス指定用オフセットを物理ブロック番号に変換し、その物理ブロックの背後で、デバイスドライバ層の `direct_access` 関数から記憶装置の参照を検索する(図11に示す)。この参照は、ファイルシステムまたはデバイスドライバ層とさらに対話せずに、ファイルが切り捨てられるか、またはファイルシステムがアンマウントされるまで、何らかのデータをその後いつでも物理ブロックに読み書きするために使用することができる。新しいインターフェースの使用は、サポートされている場合は必須であり、データの完全性の点で、従来の `read_page(s)/write_page(s)` インターフェースまたは新しい `get_xip_page` インターフェースがファイルに対してサポートされる。

10

【0041】

`read_page(s)/write_page(s)` 関数の主なユーザは、ファイルシステムのライブラリ関数である。図14は、ファイルシステムのライブラリ関数が、Linuxの汎用ファイルシステムのための `read_type` ファイル操作をどのように実行するかを示す。

20

【0042】

`generic_file_read` 関数は、`sys_read()` システムコールに関連するファイル操作を実行する。

【0043】

`generic_file_readv` 関数は、`sys_readv()` システムコールに関連するファイル操作を実行する。

【0044】

`generic_file_aio_read` 関数は、非同期IOシステムコールに関連する読み取り操作を実行する。

30

【0045】

`generic_file_sendfile` ファイル操作は、`sys_sendfile()` システムコールに関連するファイル操作を実行する。

【0046】

これらのすべての関数は `generic_mapping_read` を呼び出し、`generic_mapping_read` は、`read_page(s)` 関数を使用して(図12に示すように)、デバイスから1つまたは複数のページを読み取る。これらの関数の使用は、ファイルシステムの場合は任意だが、汎用ファイルシステムは、独自のファイル操作のインプリメンテーションを行う代わりに、すべての関数を使用する。

40

【0047】

図16は、ファイルシステムのライブラリ関数が、どのようにLinuxの汎用ファイルシステムのための `write_type` ファイル操作を実行するかを示す。

【0048】

`generic_file_write` 関数は、`sys_write()` システムコールに関連するファイル操作を実行する。

【0049】

`generic_file_writev` 関数は、`sys_writev()` システムコールに関連するファイル操作を実行する。

【0050】

50

`generic_file_aio_write`関数は、非同期I/Oシステムコールに関連する書き込み操作を実行する。

【0051】

これらのすべての関数は、`generic_file_direct_write`（オプションの`O_DIRECT`を使用して、対象ファイルを開く場合）、または`generic_file_buffered_write`を呼び出す。これらの関数は`write_page(s)`関数を使用して、1つまたは複数のページをデバイスに書き込む。

【0052】

本発明は、ライブラリ関数が、`get_xip_page`インターフェースがサポートされている場合、それを使用できるようにするために、ライブラリ関数の拡張を提供する。図15は、`read_type`操作のために、ファイルシステムのライブラリ関数に対して行われる拡張を示す。`get_xip_page`アドレス空間操作が存在するかどうかに応じて、操作を実行するために、`generic_mapping_read`関数、または新しい`do_xip_mapping_read`関数を使用して、操作が実行される。`do_xip_mapping_read`関数は、`get_xip_page`アドレス空間操作を使用して、デバイス上の対象データの参照を検索する。データ転送の場合、この参照が直接使用され、I/O操作を実行する必要はない。図17は、`write_type`操作のために、ファイルシステムライブラリ関数に対して行われる拡張を示す。`get_xip_page`アドレス空間操作が存在するかどうかに応じて、`generic_file_buffered_write/generic_file_buffered_write`関数を使用するか、または新しい`generic_file_xip_write`関数を使用して操作を行う。`generic_file_xip_write`関数は、`get_xip_page`アドレス空間操作を使用して、デバイス上の対象データに対する参照を検索する。データ転送の場合、この参照を直接使用して、I/O操作を実行する必要はない。

【0053】

`get_xip_page`アドレス空間操作をインプリメントするすべてのファイルシステムは、ライブラリ関数によりインプリメントされるすべてのファイル操作を使って、直接実行アクセスを実行するために、さらにコードを変更する必要はない。図15、図17および図18は、拡張されたライブラリ関数が、`get_xip_page`アドレス空間操作がインプリメントされるかどうかに応じて、それぞれの機能をどのようにインプリメントするかを示す。`get_xip_page`がインプリメントされる場合、すべてのライブラリ関数は、`get_xip_page`から検索された参照を使用して、記憶装置に対するすべてのデータ転送を直接実行する。図18は、ファイルメモリマッピングのために、ファイルシステムのライブラリ関数に対して行われる拡張を示す。このアプリケーションは、現在存在しない仮想メモリアドレス空間の一部にアクセスした。`Linux`の標準のアーキテクチャ依存およびコアメモリ管理機能は、結果として得られるページ障害を処理するために使用される。定期的な処理と違って、ファイルシステムは、対象ページに関連するファイルのために、`filemap_xip_nopage`ハンドラをインストールした。このハンドラは、`get_xip_page`を使用して、デバイス上の対象データの参照を検索する。この参照は`do_nopage`関数に返され、アプリケーションの仮想アドレス空間ページテーブル内にページテーブルエントリを生成し、アプリケーションが、オペレーティングシステムがさらに関与せずに、デバイス上のデータを直接使用できるようにする。ページテーブルのエントリは、ファイルマッピングがプライベートであることが選択された場合（標準機構が適用される）、後に、書き込み機構上のコピーの対象になる。

【0054】

直接実行の効果は、アプリケーションのバイナリファイルおよび共有ライブラリが、`Linux`のファイルマッピングの対象になり、その結果、上記のページ障害機構の対象になるため達成される。

10

20

30

40

50

## 【 0 0 5 5 】

上記の拡張は、LinuxオペレーティングシステムのI/O関連コンポーネント内で、全体の構造および層隔離を元の状態に維持する。デバイスドライバ層は、物理ブロック番号をデバイスにマップするが、ファイルまたはその他のロジックオブジェクトと連動しない。ファイルシステムは、ロジックファイルと物理ブロック番号のアドレス指定との間でマッピングを実行するが、デバイスの構造と連動しない。一方、データ転送自体は、上下逆転する。デバイスドライバは、新たな拡張を使用する場合、データを転送しない。データは、ファイル操作ライブラリ関数で直接背転送される。この解決方法の効果は、デバイスドライバ層および汎用ファイルシステムに対して、非常に小さくかつ非侵入的な変化のみを含むことである。汎用ファイルシステムは、メモリ消費量の減少および実行の増加など、直接機構の利点を容易に取得することができる。

10

## 【図面の簡単な説明】

## 【 0 0 5 6 】

【図1】本発明をインプリメントするために必要なコンピュータシステムのブロック図を示す。

【図2】本発明のいくつかの実施態様のゲストとして、図1に示すコンピュータシステムをホストする仮想マシン環境を示す。

【図3】先行技術のオペレーティングシステムにより提供される仮想メモリおよびデマンドページング機能を示す。

【図4】先行技術のオペレーティングシステムにより提供される直接実行機能を示す。

20

【図5】直接実行をインプリメントする先行技術のオペレーティングシステムのコンポーネント構造を示す。

【図6】本発明を使用する直接実行をインプリメントするオペレーティングシステムのコンポーネント構造を示す。

【図7】本発明によるI/Oベースデバイスにアクセスするために使用される図6のオペレーティングシステムのコンポーネントを通る制御フローを示す。

【図8】本発明によるメモリアドレス指定デバイスに対する直接実行アクセスを実行するために使用される図6のオペレーティングシステムのコンポーネントを通る制御フローを示す。

【図9】図7および図8に示される2つの制御フローのどちらを使用するかを選択するために使用される図6のオペレーティングシステムによる決定ロジックを示す。

30

【図10】先行技術のLinuxオペレーティングシステムにインプリメントされるデバイスの抽象化を示す。

【図11】本発明をインプリメントするLinuxオペレーティングシステム内のデバイスドライバ層に行われる拡張を示す。

【図12】汎用ファイルシステムが、1つまたは複数のページを先行技術のLinuxオペレーティングシステム内のデバイスに読み書きするためのアドレス空間操作をどのように提供するかを示す。

【図13】本発明をインプリメントするLinuxオペレーティングシステム内のファイルシステムのアドレス空間操作に行われる拡張を示す。

40

【図14】ファイルシステムのライブラリ関数が、先行技術のLinuxオペレーティングシステム内において、汎用ファイルシステムに対してread typeファイル操作をどのように実行するかを示す。

【図15】本発明をインプリメントするLinuxオペレーティングシステム内において、read type操作のためにファイルシステムのライブラリ関数に対して行われる拡張を示す。

【図16】ファイルシステムのライブラリ関数が、先行技術のLinuxオペレーティングシステム内において、汎用ファイルシステムのためのwrite typeファイル操作をどのように実行するかを示す。

【図17】本発明をインプリメントするLinuxオペレーティングシステム内において

50



、write type 操作のためにファイルシステムのライブラリ関数に行われる拡張を示す。

【図 18】本発明をインプリメントする Linux オペレーティングシステムにおいて、ファイルメモリマッピングのためにファイルシステムのライブラリ関数に行われる拡張を示す。

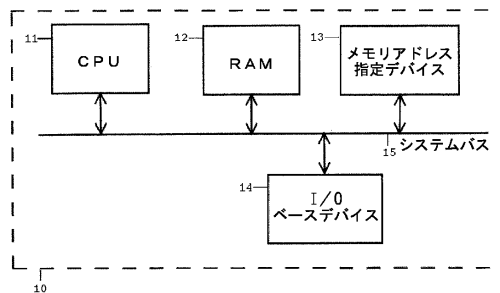
【符号の説明】

【 0 0 5 7 】

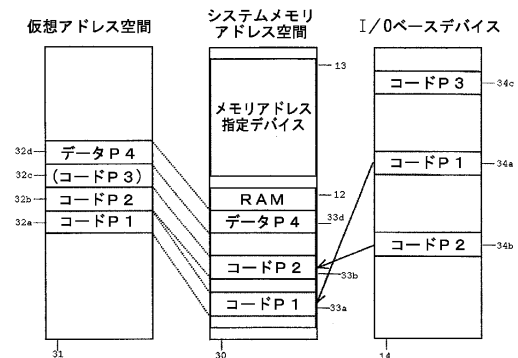
- 1 2    R A M
- 1 3    メモリアドレス指定デバイス
- 1 4    I / O ベースデバイス
- 4 4、5 3    デバイスドライバ
- 4 9    アプリケーションプログラム
- 5 0    オペレーティングシステム
- 5 1    メモリ / ファイルマネージャ
- 5 2    ファイルシステムドライバ

10

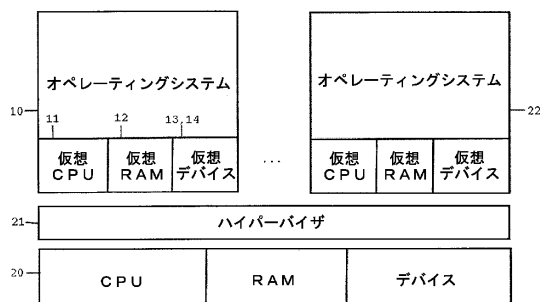
【図 1】



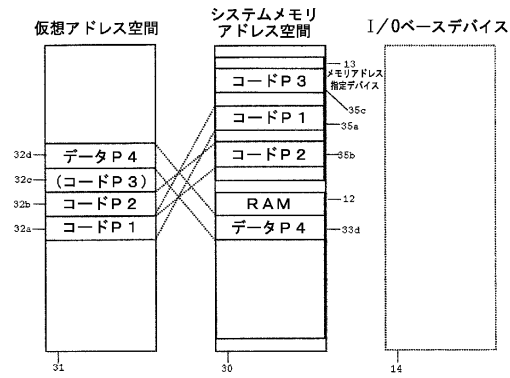
【図 3】



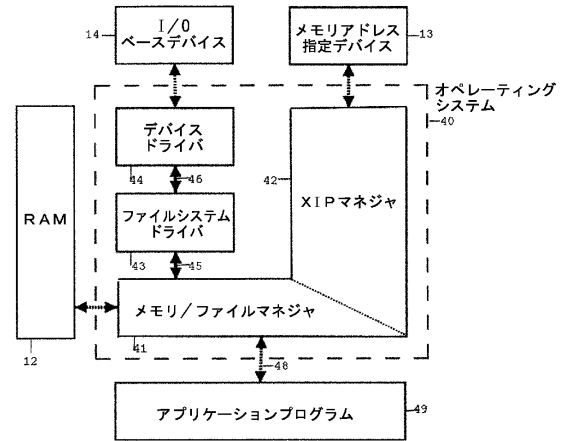
【図 2】



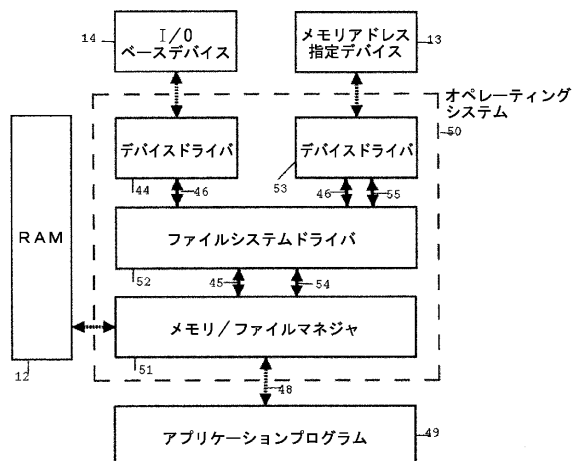
【図 4】



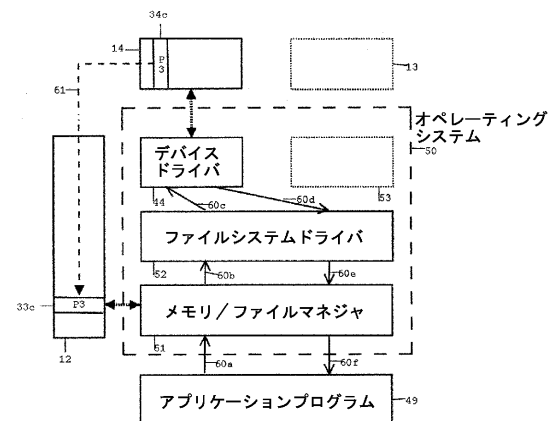
【図 5】



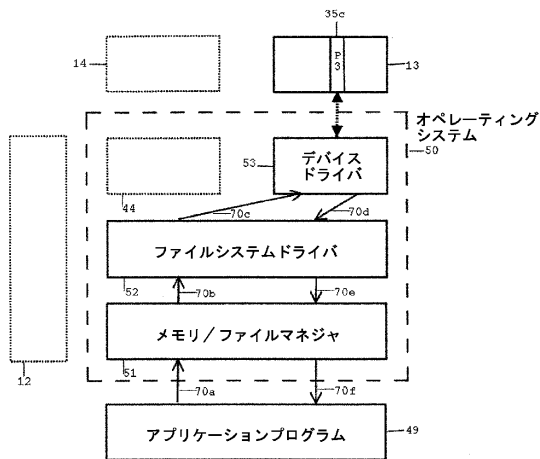
【図 6】



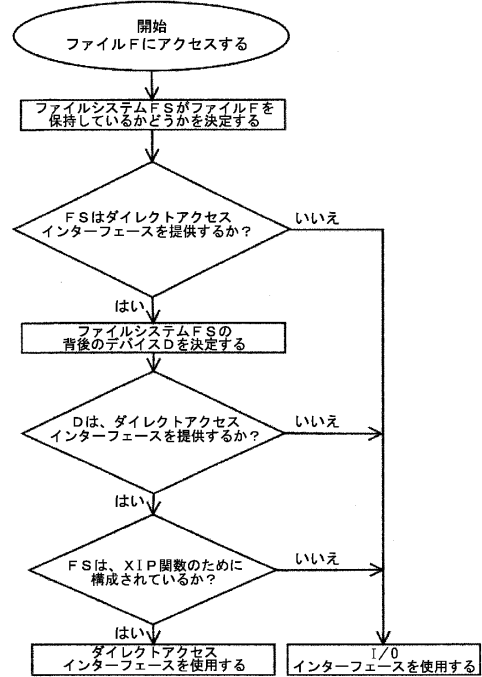
【図 7】



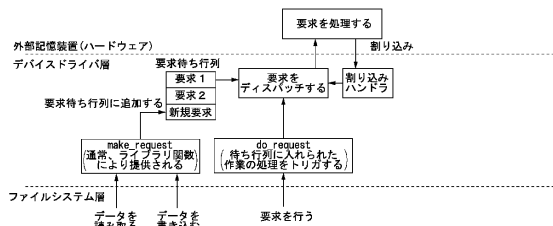
【図 8】



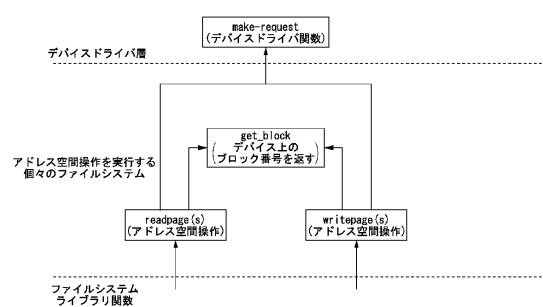
【図 9】



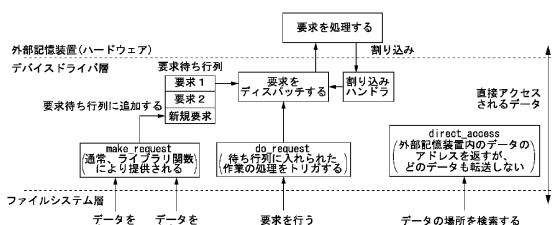
【図 10】



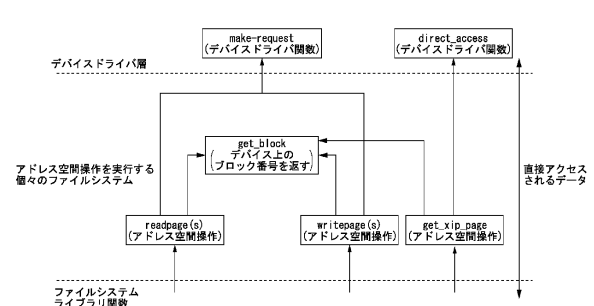
【図 12】



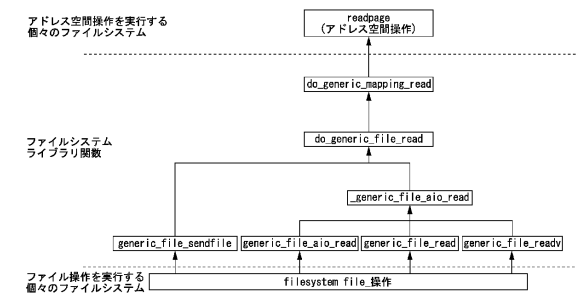
【図 11】



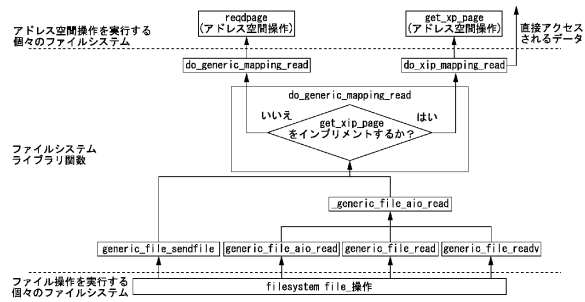
【図 13】



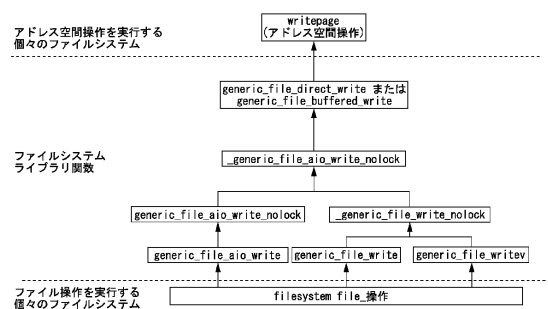
【図 14】



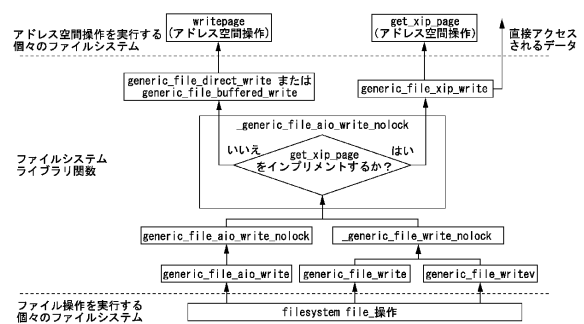
【図 15】



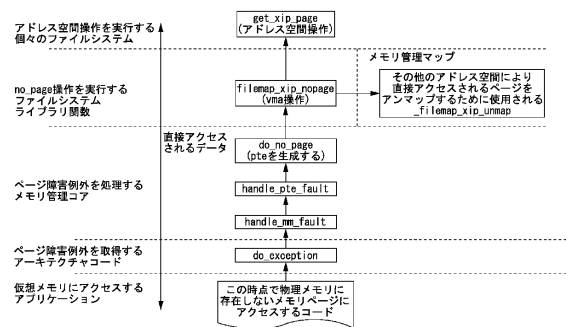
【図 16】



【図 17】



【図 18】



## フロントページの続き

(51)Int.Cl. F I  
G 0 6 F 13/10 3 4 0 Z

(74)代理人 100086243

弁理士 坂口 博

(72)発明者 カーステン・オット

ドイツ連邦共和国 D - 7 1 0 3 4、ボエブリンゲン、モーレネル・ウェグ 1 1 7

(72)発明者 ウルリヒ・ウェイガンド

ドイツ連邦共和国 D - 7 1 0 8 8、ホルツガーリンゲン、リヒテンステイン通り 8

審査官 原 秀人

(56)参考文献 特開平6 - 1 3 1 2 6 6 ( J P , A )

米国特許出願公開第2 0 0 2 / 0 0 6 9 3 4 2 ( U S , A 1 )

前川 守 外2名,分散オペレーティングシステム,日本,共立出版株式会社,1 9 9 3年1 0月 5日,第1版,p . 1 2 6 ~ 1 2 9

後藤 大地,ファイルシステムって、なに? より広範な可用性・堅牢性を実現するための方法,TransTECH,日本,株式会社翔泳社,1 9 9 9年1 2月 1日,第8巻 第1 4号,p . 3 0 ~ 3 9

Frank G. Soltis,Inside the AS/400,日本,(株)インフォ・クリエイツ,2 0 0 0年 2月1 0日,日本語版 第2版,p . 2 2 7 ~ 2 4 1

(58)調査した分野(Int.Cl., DB名)

G 0 6 F 1 2 / 1 0

G 0 6 F 1 2 / 0 0

G 0 6 F 9 / 4 6

G 0 6 F 1 3 / 1 0