

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5385977号  
(P5385977)

(45) 発行日 平成26年1月8日(2014.1.8)

(24) 登録日 平成25年10月11日(2013.10.11)

(51) Int. Cl. F I  
**G06F 9/48 (2006.01)** G O 6 F 9/46 4 5 2 H  
**G06F 13/10 (2006.01)** G O 6 F 13/10 3 4 O B

請求項の数 34 (全 31 頁)

(21) 出願番号	特願2011-511688 (P2011-511688)	(73) 特許権者	500551079
(86) (22) 出願日	平成21年5月8日(2009.5.8)		ソニー コンピュータ エンタテインメン ト アメリカ リミテッド ライアビリテ イ カンパニー
(65) 公表番号	特表2011-525010 (P2011-525010A)		アメリカ合衆国、カリフォルニア州 94 404-2175、フォスター・シティー 、セカンド・フロアー、イースト・ヒルス デイル・ブルバード 919
(43) 公表日	平成23年9月8日(2011.9.8)	(74) 代理人	100105924
(86) 国際出願番号	PCT/US2009/043387		弁理士 森下 賢樹
(87) 国際公開番号	W02009/148764	(72) 発明者	ターラー、アンドリュウ、アール、 アメリカ合衆国、カリフォルニア州 94 404、フォスター・シティー、セカンド ・フロアー、イースト・ヒルスデイル・ブ ルバード 919
(87) 国際公開日	平成21年12月10日(2009.12.10)		最終頁に続く
審査請求日	平成23年1月31日(2011.1.31)		
(31) 優先権主張番号	12/130,892		
(32) 優先日	平成20年5月30日(2008.5.30)		
(33) 優先権主張国	米国 (US)		

(54) 【発明の名称】 ファイル入出力スケジューラ

(57) 【特許請求の範囲】

【請求項1】

プロセッサユニット、メモリ、およびメディアデバイスをもつシステムにおいて、前記プロセッサユニットは、メインプロセッサと、関連づけられたローカルメモリをもつ補助プロセッサとを含み、メディアデバイスに対する入出力(I/O)をハンドリングする方法であって、

a) メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力I/Oリクエストを受け取るステップと、

b) 前記入力I/Oリクエストを完了するための予測時間 $T_p$ を計算するステップと、

c) 前記入力I/Oリクエストのスケジュール内の位置が予測時間 $T_p$ に少なくとも部分的に依存するように、前記入力I/Oリクエストをメモリ内に具体化されたスケジュールに挿入するステップと、

d) メディアデバイスからローカルメモリにデータを転送することによって、前記スケジュールにしたがって前記入力I/Oリクエストをサービスし、補助プロセッサを用いてローカルメモリのデータを変換するステップとを含み、

前記入力I/Oリクエストを前記スケジュールに挿入するステップは、

複数の異なるI/Oリクエストをサービスするための合計時間Tを決定するステップと、

前記複数のI/Oリクエストの2以上の異なる順序に対して前記合計時間Tの2以上の異なる値を計算するステップと、

10

20

前記合計時間の最適値にもとづいて前記入力 I / O リクエストを前記スケジュールに挿入するステップとを含み、

ステップ d ) は、データオーバーレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられることを特徴とする方法。

【請求項 2】

ステップ b ) は、メディアデバイスのパラメータと前記 I / O リクエストから決定される係数を用いた性能特性式 ( P C E ) にしたがって予測時刻  $T_p$  を計算するステップを含む請求項 1 の方法。

【請求項 3】

前記パラメータは、リクエストオーバーヘッド  $x$ 、ヘッド運動レート  $y$ 、およびスループット  $z$  を含む請求項 2 の方法。

【請求項 4】

前記係数は、オーバーヘッド係数  $A$ 、ヘッド移動距離  $B$ 、および転送データ量  $C$  を含む請求項 3 の方法。

【請求項 5】

前記性能特性式は、 $T_p = A x + B y + C z$  の形である請求項 4 の方法。

【請求項 6】

ステップ c ) は、前記スケジュール内の前記複数の異なる I / O リクエストに対する異なる係数をもつ前記性能特性式を繰り返すことで前記複数の I / O リクエストをサービスするための前記合計時間  $T$  を決定するステップを含む請求項 2 の方法。

【請求項 7】

前記合計時間  $T$  の最適値の値は、前記複数の I / O リクエストの 2 以上の異なる順序に対する合計時間  $T$  の最小値である請求項 1 の方法。

【請求項 8】

ステップ b ) は、I / O リクエストに対してかかった現実の時間  $T_a$  を測定するステップと、前記現実の時間  $T_a$  を前記係数とともに用いて、未知パラメータの値に対して P C E を反復して解くステップとを含む請求項 2 の方法。

【請求項 9】

ステップ d ) は、圧縮されたネストしたアーカイブから特定のデータを抽出するステップをさらに含む請求項 1 の方法。

【請求項 10】

前記システムはさらにグラフィックプロセッサと関連づけられたグラフィックメモリとを含み、ステップ d ) は、メディアデバイスからグラフィックメモリにデータを読み込むステップと、グラフィックメモリからローカルメモリにデータを転送するステップと、補助プロセッサを用いてデータ上でデータ変換を実行して、変換されたデータを生成するステップと、変換されたデータをグラフィックメモリに返送するステップとを含む請求項 1 の方法。

【請求項 11】

前記プライマリデータユニットのソースと前記セカンダリデータソースは同じメディア上にある請求項 10 の方法。

【請求項 12】

前記プライマリデータユニットのソースと前記セカンダリデータソースは異なるメディア上にある請求項 10 の方法。

【請求項 13】

前記セカンダリデータソースは仮想的なデータソースである請求項 10 の方法。

【請求項 14】

セカンダリデータユニットはコールバックによって指定される請求項 13 の方法。

【請求項 15】

前記システムは高速ストレージデバイスをさらに含み、ステップ d ) は、メディアデバ

10

20

30

40

50

イスから高速ストレージデバイスへ1以上のデータユニットをプリフェッチするステップと、それに続いて高速ストレージデバイスからそのデータユニットを転送するステップとを含む請求項1の方法。

【請求項16】

高速ストレージデバイスはハードディスクドライブである請求項15の方法。

【請求項17】

メディアデバイスは大容量の光ディスクドライブである請求項16の方法。

【請求項18】

1以上のデータユニットをプリフェッチするステップは、I/Oがアイドルであるとき、ある時間間隔の間、1以上のデータユニットをプリフェッチする請求項15の方法。

10

【請求項19】

1以上のデータユニットをプリフェッチするステップは、メディアデバイスから1以上のデータユニットの第1グループをプリフェッチするステップと、メディアデバイスがアイドルかどうかを調べるステップと、もしI/Oがアイドルなら、メディアデバイスから1以上のデータユニットの第2グループをプリフェッチし、もしI/Oがアイドルでないなら、前記第2グループのプリフェッチを遅延させるステップとを含む請求項18の方法。

【請求項20】

ステップb)は、前記システムに関連づけられたメディアデバイスよりも低速のメディアのオペレーションをエミュレートするために前記メディアデバイスのパラメータを調整するステップを含む請求項1の方法。

20

【請求項21】

プロセッサユニットと、  
前記プロセッサユニットと結合したメモリと、  
前記プロセッサユニットと結合したメディアデバイスと、  
前記メモリに具体化されたプロセッサ実行可能な命令セットとを含むシステムであって、

前記命令は実行されたときに、メディアデバイスに対する入出力(I/O)をハンドリングする方法を実行するように構成され、前記方法は、

a)メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力I/Oリクエストを受け取るステップと、

30

b)前記入力I/Oリクエストを完了するための予測時間 $T_p$ を計算するステップと、

c)前記入力I/Oリクエストのスケジュール内の位置が予測時間 $T_p$ に少なくとも部分的に依存するように、前記入力I/Oリクエストをメモリ内に具体化されたスケジュールに挿入するステップと、

d)前記スケジュールにしたがって前記入力I/Oリクエストをサービスするステップとを含み、

前記プロセッサ実行可能な命令セットは、デアーカイバレイヤ、RAMキャッシュレイヤ、スケジューラキャッシュレイヤ、オーバレイレイヤ、カタログキャッシュレイヤ、データ変換レイヤ、またはシミュレートされたデータレイヤを含む1以上のメディアフィルタレイヤをさらに含み、

40

前記入力I/Oリクエストを前記スケジュールに挿入するステップは、複数の異なるI/Oリクエストをサービスするための合計時間 $T$ を決定するステップと

、前記複数のI/Oリクエストの2以上の異なる順序に対して前記合計時間 $T$ の2以上の異なる値を計算するステップと、

前記合計時間の最適値にもとづいて前記入力I/Oリクエストを前記スケジュールに挿入するステップとを含み、

ステップd)は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカン

50

ダリデータユニットで置き換えられることを特徴とするシステム。

【請求項 2 2】

メインプロセッサと、関連づけられたローカルメモリをもつ補助プロセッサと、メインメモリと、メディアデバイスとを含むシステムにおいて、メディアデバイスに対する入出力（I/O）をハンドリングする方法であって、

a) メディアデバイスに対してデータを転送するために、メインプロセッサ上で動作するアプリケーションから入力 I/O リクエストを受け取るステップと、

b) 前記入力 I/O リクエストをメインメモリ内に具体化されたスケジュールに挿入するステップと、

c) スケジュールと 1 以上のフィルタにしたがって前記入力 I/O リクエストをサービスするステップとを含み、1 以上のフィルタの少なくとも 1 つは補助プロセッサによって実行され、

前記入力 I/O リクエストを前記スケジュールに挿入するステップは、

複数の異なる I/O リクエストをサービスするための合計時間 T を決定するステップと

、  
前記複数の I/O リクエストの 2 以上の異なる順序に対して前記合計時間 T の 2 以上の異なる値を計算するステップと、

前記合計時間の最適値にもとづいて前記入力 I/O リクエストを前記スケジュールに挿入するステップとを含み、

ステップ d) は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられることを特徴とする方法。

【請求項 2 3】

プロセッサユニット、メモリ、およびメディアデバイスをもつシステムにおいて、メディアデバイスに対する入出力（I/O）をハンドリングする方法であって、

a) メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力 I/O リクエストを受け取るステップと、

b) 前記入力 I/O リクエストを完了するための予測時間  $T_p$  を計算するステップと、

c) 前記入力 I/O リクエストのスケジュール内の位置が予測時間  $T_p$  に少なくとも部分的に依存するように、前記入力 I/O リクエストをメモリ内に具体化されたスケジュールに挿入するステップと、

d) 前記スケジュールにしたがって前記入力 I/O リクエストを実行するステップとを含み、ステップ d) は、圧縮されたネストしたアーカイブから特定のデータを抽出するステップを含み、

前記入力 I/O リクエストを前記スケジュールに挿入するステップは、

複数の異なる I/O リクエストをサービスするための合計時間 T を決定するステップと

、  
前記複数の I/O リクエストの 2 以上の異なる順序に対して前記合計時間 T の 2 以上の異なる値を計算するステップと、

前記合計時間の最適値にもとづいて前記入力 I/O リクエストを前記スケジュールに挿入するステップとを含み、

ステップ d) は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられることを特徴とする方法。

【請求項 2 4】

プロセッサユニット、メモリ、およびメディアデバイスをもつシステムにおいて、メディアデバイスに対する入出力（I/O）をハンドリングする方法であって、

a) メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力 I/O リクエストを受け取るステップと、

b) 前記入力 I/O リクエストを完了するための予測時間  $T_p$  を計算するステップと、

10

20

30

40

50

c) 前記入力 I / O リクエストのスケジュール内の位置が予測時間  $T_p$  に少なくとも部分的に依存するように、前記入力 I / O リクエストをメモリ内に具体化されたスケジュールに挿入するステップと、

d) 前記スケジュールにしたがって前記入力 I / O リクエストを実行するステップとを含み、ステップ d) は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられ、

前記入力 I / O リクエストを前記スケジュールに挿入するステップは、

複数の異なる I / O リクエストをサービスするための合計時間 T を決定するステップと

、  
前記複数の I / O リクエストの 2 以上の異なる順序に対して前記合計時間 T の 2 以上の異なる値を計算するステップと、

前記合計時間の最適値にもとづいて前記入力 I / O リクエストを前記スケジュールに挿入するステップとを含み、

ステップ d) は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられることを特徴とする方法。

【請求項 25】

前記プライマリデータユニットのソースと前記セカンダリデータソースは同じメディア上にある請求項 24 の方法。

【請求項 26】

前記プライマリデータユニットのソースと前記セカンダリデータソースは異なるメディア上にある請求項 24 の方法。

【請求項 27】

前記セカンダリデータソースは仮想的なデータソースである請求項 24 の方法。

【請求項 28】

セカンダリデータユニットはコールバックによって指定される請求項 27 の方法。

【請求項 29】

プロセッサユニット、メモリ、高速ストレージデバイス、およびメディアデバイスをもつシステムにおいて、メディアデバイスに対する入出力 (I / O) をハンドリングする方法であって、

a) メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力 I / O リクエストを受け取るステップと、

b) 前記入力 I / O リクエストを完了するための予測時間  $T_p$  を計算するステップと、

c) 前記入力 I / O リクエストのスケジュール内の位置が予測時間  $T_p$  に少なくとも部分的に依存するように、前記入力 I / O リクエストをメモリ内に具体化されたスケジュールに挿入するステップと、

d) 前記スケジュールにしたがって前記入力 I / O リクエストを実行するステップとを含み、ステップ d) は、メディアデバイスから高速ストレージデバイスへ 1 以上のデータユニットをプリフェッチするステップと、それに続いて高速ストレージデバイスからそのデータユニットを転送するステップとを含み、

前記入力 I / O リクエストを前記スケジュールに挿入するステップは、

複数の異なる I / O リクエストをサービスするための合計時間 T を決定するステップと

、  
前記複数の I / O リクエストの 2 以上の異なる順序に対して前記合計時間 T の 2 以上の異なる値を計算するステップと、

前記合計時間の最適値にもとづいて前記入力 I / O リクエストを前記スケジュールに挿入するステップとを含み、

ステップ d) は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカン

10

20

30

40

50

ダリデータユニットで置き換えられることを特徴とする方法。

【請求項 30】

高速ストレージデバイスはハードディスクドライブである請求項 29 の方法。

【請求項 31】

メディアデバイスは大容量の光ディスクドライブである請求項 30 の方法。

【請求項 32】

1 以上のデータユニットをプリフェッチするステップは、I/O がアイドルであるとき、ある時間間隔の間、1 以上のデータユニットをプリフェッチする請求項 31 の方法。

【請求項 33】

1 以上のデータユニットをプリフェッチするステップは、メディアデバイスから 1 以上のデータユニットの第 1 グループをプリフェッチするステップと、メディアデバイスがアイドルかどうかを調べるステップと、もし I/O がアイドルなら、メディアデバイスから 1 以上のデータユニットの第 2 グループをプリフェッチし、もし I/O がアイドルでないなら、前記第 2 グループのプリフェッチを遅延させるステップとを含む請求項 31 の方法。

10

【請求項 34】

プロセッサユニット、メモリ、高速ストレージデバイス、およびメディアデバイスをもつシステムにおいて、メディアデバイスに対する入出力 (I/O) をハンドリングする方法であって、

a) メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力 I/O リクエストを受け取るステップと、

b) 前記入力 I/O リクエストを完了するための予測時間  $T_p$  を計算するステップと、

c) 前記入力 I/O リクエストのスケジュール内の位置が予測時間  $T_p$  に少なくとも部分的に依存するように、前記入力 I/O リクエストをメモリ内に具体化されたスケジュールに挿入するステップと、

d) 前記スケジュールにしたがって前記入力 I/O リクエストを実行するステップとを含み、ステップ b) は、前記システムに関連づけられたメディアデバイスよりも低速のメディアのオペレーションをエミュレートするために前記 メディアデバイスのパラメータを調整するステップを含み、

20

前記入力 I/O リクエストを前記スケジュールに挿入するステップは、

複数の異なる I/O リクエストをサービスするための合計時間  $T$  を決定するステップと

30

、前記複数の I/O リクエストの 2 以上の異なる順序に対して前記合計時間  $T$  の 2 以上の異なる値を計算するステップと、

前記合計時間の最適値にもとづいて前記入力 I/O リクエストを前記スケジュールに挿入するステップとを含み、

ステップ d) は、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられることを特徴とする方法。

【発明の詳細な説明】

40

【技術分野】

【0001】

[優先権主張]

本出願は、2008年5月30日に出願された、同一譲受人による米国特許出願第 12 / 130,892 号の優先権の利益を主張し、開示内容全体を参照により組み入れる。

【0002】

[技術分野]

本発明の実施の形態はコンピュータゲームおよび関連するアプリケーションに関し、特に、コンピュータゲームおよび関連するアプリケーションにおけるファイル入出力 (I/O) 管理に関する。

50

## 【背景技術】

## 【0003】

ビデオゲームのような多くのソフトウェアアプリケーションは、アプリケーション内のメディアアクセスをより効率よくかつ信頼あるものとするためのファイル入力/出力（入出力；I/O）スケジューラを含む。ファイル入出力システム（file I/O system: FIOS）は、いくつかのパーツからなるファイルにアクセスするための、スケジューラおよび任意の入出力フィルタレイヤを含むミドルウェアレイヤである。スケジューラは、典型的には任意の期限および優先度制約による条件の下、可能な限り最短時間で完了するように入出力リクエストを最適に並べ替えるように設計されている。

## 【0004】

多くの異なるゲームコンポーネントは、ストレージメディア内のファイルへの入出力アクセスを要求する。オーディオコンポーネントはオーディオファイルを読み込み、ゲームプレイのエンジンはレベルの定義を読み込み、グラフィックコンポーネントはテクスチャマップおよびモデルを読み込み、動画コンポーネントは音声映像ファイルを読み込み、サブシステムは巨大なWADファイルを読み込む。メディアは、UMD（universal media disc）、CD（compact disc）、DVD（digital video disc）、Blu-ray（登録商標）ディスク等のような光学ディスク、ハードディスクのような中間ストレージメディア、あるいは発達するプラットフォームの等の他のメディアタイプのような、ゲーム配信メディアであってもよい。

## 【0005】

ビデオゲームのような単一のアプリケーションは、典型的には多数のコンポーネントを持っており、各コンポーネントは独自の入出力要求を持っている。あるものはメディアへのストリーミングアクセスを要求する。ここで入出力システムは、コンポーネントがゲームプレイにストリーム化されたデータを連続的に提示できるようにファイル内のデータをストリーム化する。例えば、音声コンポーネントは典型的にはサウンドトラックの再生のためにオーディオファイルをストリーミングする。動画コンポーネントはプレイヤに動画を再生するために音声映像コンテンツをストリーミングする。他のコンポーネントは、そのコンポーネントが処理するためのデータをファイルからチャンクという形で取ってくる非ストリーミングアクセスのみを必要とする。これらのコンポーネントは安定したデータフローを要求しないが、チャンクの配信のタイミングはしばしば極めて重要となる。もしアプリケーションがデータが必要になるときにそれを受け取れないと、パフォーマンスを下げられるかも知れない。

## 【0006】

ビデオゲームはしばしばかなりの量の入出力を行う。ゲームのイベントはしばしば必ず間に合わなければならない期限を持つタイムドリブン型のものである。例えば、もしプレイヤがA地点からB地点まで移動すると、ゲームはプレイヤがB地点に到達する前にB地点に関連するデータ用のファイルを持っておく必要がある。ゲームコンポーネントは低水準の入出力基本命令（primitives）を使ってメディアからデータを読み出すことができるが、これらの基本命令は多数のコンポーネントが同じデバイスからのデータを同時に要求するときにデバイス接続を扱うことはしない。混雑した入出力リクエストはストリーミングデータフローを中断したり、必要なときにコンポーネントがデータのチャンクを取得することを妨げたりする可能性がある。989 SoundのStreamsafeのような高レベルシステムは少なくともファイルに対して信頼できるストリームアクセスを供給するため、ストリーム化されたデータは中断されることはない。残念ながら、そのようなシステムは一つだけデータストリームを割り当てるので、非ストリームデータに対する競合をうまく扱えない。SCEAは普通のデータアクセスに対する競合を解消する標準的なレイヤを提供しない。

## 【発明の概要】

## 【発明が解決しようとする課題】

## 【0007】

アプリケーションのコンポーネントは、典型的にはお互いの入出力活動を関知していないので、それらの入出力リクエストはかなりの量の過度な探索（データを探すために行きつ戻りつすること）を含む、とても非効率な入出力パターンにしばしば陥る。ゲームが大きくなると、過度な探索および非効率なデータ取得が増加する。CDのような配布メディアに対するディスクのマスター化では、頻繁にアクセスされるデータブロックを重複させたり、ディスク全体にデータブロックを分散させることで、ストレージデバイスの読み取りヘッドがどこにあってもコピーがいつも近くにあるようにする技術を通してこの非効率性のある程度避けることができる。ファイルシステムは、データが物理的にディスクにある場所を明らかにするとは限らないので、I/Oシステムはそのデータを最大限活用することができるとは限らない。

10

【0008】

本発明の実施の形態はこのような文脈の中で生じた。

【課題を解決するための手段】

【0009】

本発明の実施の形態によれば、プロセッサユニット、メモリ、およびメディアデバイスをもつシステムにおいて、メディアデバイスに対する入出力（I/O）をハンドリングする方法が実装される。メディアデバイスに対してデータを転送するために、プロセッサ上で動作するアプリケーションから入力I/Oリクエストを受け取る。前記入力I/Oリクエストを完了するための予測時間 $T_p$ を計算する。前記入力I/Oリクエストをメモリ内に具体化されたスケジュールに挿入する。前記入力I/Oリクエストのスケジュール内の位置は予測時間 $T_p$ に少なくとも部分的に依存する。前記スケジュールにしたがって前記入力I/Oリクエストをサービスする。

20

【0010】

ある実施の形態では、プロセッサ実行可能な命令セットがメモリに具体化される。前記命令は実行されたときに、上述の方法を実行するように構成されてもよい。プロセッサ実行可能な命令セットは1以上のメディアフィルタレイヤを含む。メディアフィルタレイヤは、デアーカイバレイヤ、RAMキャッシュレイヤ、スケジューラキャッシュ（たとえば、ハードディスクドライブ（HDD）キャッシュ）レイヤ、オーバレイヤ、カタログキャッシュレイヤ、データ変換レイヤ、またはシミュレートされたデータレイヤを含んでもよい。

30

【0011】

ある実施の形態では、メディアデバイスのパラメータと入力I/Oリクエストから決定される係数を用いた性能特性式（performance characteristic equation；PCE）にしたがって予測時間 $T_p$ を決定してもよい。一例として、一般性を失うことなく、前記パラメータは、リクエストオーバーヘッド $x$ 、ヘッド運動レート $y$ 、およびスループット $z$ を含んでもよい。そのような例において、前記係数は、オーバーヘッド係数 $A$ 、ヘッド移動距離 $B$ 、および転送データ量 $C$ を含んでもよい。前記性能特性式は、 $T_p = Ax + By + Cz$ の形であってもよい。

【0012】

ある実施の形態では、I/Oリクエストに対するPCE係数はI/Oリクエストと初期メディア状態にもとづいて決定されてもよい。初期メディア状態は、性能に影響するデバイスの状態に関する情報を含む。一例として、一般性を失うことなく、その情報には、デバイスのレディ状態とアクセスされる最後の論理ブロックアドレス（LBA）が含まれる。ある実施の形態では、I/Oリクエストと初期メディア状態からのPCE係数の計算は、更新されたメディア状態を返し、これは別のPCE計算に対する初期メディア状態として用いることができる。

40

【0013】

ある実施の形態では、PCE係数は、所与のI/Oリクエストに対する予測時間 $T_p$ を計算するために、PCEにおける未知の値と結合されてもよい。その後、I/Oリクエストの順序付けされたリストに対する合計時間 $T$ は、リスト内の各I/Oリクエストに対し

50

てPCEとメディア状態を反復し、その結果得られる時間 $T_p$ を加算することにより計算されてもよい。この処理は、I/Oリクエストの順序付けされたリストのある順列またはすべての順列に対して繰り返されてもよい。リストの好ましい順序付けは、各順序づけに対して計算された時間 $T$ と場合によっては他の基準とにもとづいて、複数の順列の中から選択されてもよい。ある実施の形態では、好ましい順序付けは、最小時間 $T$ を与える順列であってもよい。順序付けに影響する付加的な基準は、これに限定しないが、I/Oリクエスト優先度、ストリームバッファ状態、およびレーテンシ要件であってもよい。

【0014】

ある実施の形態では、I/Oリクエストに対してかかった現実の時間 $T_a$ を測定し、前記現実の時間 $T_a$ を前記係数とともに用いて、未知パラメータの値に対してPCEを反復して解いてもよい。これにより、モデルが現実のドライブ性能で絶えず更新される。

10

【0015】

ある実施の形態では、I/Oリクエストを実行するステップは、圧縮されたネストしたアーカイブから特定のデータを抽出するステップを含んでもよい。

【0016】

ある実施の形態では、前記プロセッサユニットは、メインプロセッサと、関連づけられたローカルメモリをもつ補助プロセッサとを含んでもよい。そのような実施の形態では、I/Oリクエストを実行するステップは、メディアデバイスからローカルメモリにデータを転送するステップと、補助プロセッサを用いてローカルメモリのデータを変換するステップとを含んでもよい。

20

【0017】

ある実施の形態では、前記システムはさらにグラフィックスプロセッサと関連づけられたグラフィックスメモリとを含んでもよい。そのような実施の形態では、I/Oリクエストを実行するステップは、メディアデバイスからグラフィックスメモリにデータを読み込むステップと、グラフィックスメモリからローカルメモリにデータを転送するステップと、補助プロセッサを用いてデータ上でデータ変換を実行して、変換されたデータを生成するステップと、変換されたデータをグラフィックスメモリに返送するステップとを含んでもよい。

【0018】

ある実施の形態では、I/Oリクエストを実行するステップは、データオーバレイを実行するステップを含み、メディアデバイスとメモリ間で転送されるプライマリデータユニットは、セカンダリデータソースからのセカンダリデータユニットで置き換えられてもよい。前記プライマリデータユニットのソースと前記セカンダリデータソースは同じメディア上であっても、異なるメディア上であってもよい。前記セカンダリデータソースは仮想的なデータソースであってもよい。セカンダリデータユニットはコールバックによって指定されてもよい。

30

【0019】

ある実施の形態では、前記システムは高速ストレージデバイスをさらに含んでもよい。そのような実施の形態では、I/Oリクエストを実行するステップは、メディアデバイスから高速ストレージデバイスへ1以上のデータユニットをプリフェッチするステップと、それに続いて高速ストレージデバイスからそのデータユニットを転送するステップとを含んでもよい。一例として、高速ストレージデバイスはハードディスクドライブであり、メディアデバイスはブルーレイ（登録商標）ディスクドライブのような大容量光ディスクドライブであってもよい。I/Oがアイドルであるとき、ある時間間隔の間、1以上のデータユニットをプリフェッチしてもよい。ある実施の形態では、プリフェッチするステップは、メディアデバイスから1以上のデータユニットの第1グループをプリフェッチするステップと、メディアデバイスがアイドルかどうかを調べるステップと、もしI/Oがアイドルなら、メディアデバイスから1以上のデータユニットの第2グループをプリフェッチし、もしI/Oがアイドルでないなら、前記第2グループのプリフェッチを遅延させるステップとを含んでもよい。

40

50

## 【 0 0 2 0 】

ある実施の形態では、予測時間  $T_p$  を計算するステップは、前記システムに関連づけられたメディアデバイスよりも低速のメディアのオペレーションをエミュレートするために前記パラメータを調整するステップを含んでもよい。

## 【 0 0 2 1 】

本発明の実施の形態によれば、メインプロセッサと、関連づけられたローカルメモリをもつ補助プロセッサと、メインメモリと、メディアデバイスとを含むシステムにおいて、メディアデバイスに対する入出力 (I/O) をハンドリングする方法が実装される。この方法は、a) メディアデバイスに対してデータを転送するために、メインプロセッサ上で動作するアプリケーションから入力 I/O リクエストを受け取るステップと、b) 前記入力 I/O リクエストをメインメモリ内に具体化されたスケジュールに挿入するステップと、c) スケジュールと 1 以上のフィルタにしたがって前記入力 I/O リクエストをサービスするステップとを含み、1 以上のフィルタの少なくとも 1 つは補助プロセッサによって実行される。

## 【 図面の簡単な説明 】

## 【 0 0 2 2 】

本発明の実施の形態は、添付の図面と併せて、以下の詳細な説明を考慮することによって容易に理解される。

【 図 1 】 本発明の実施の形態に係るファイル入出力システム (FIOS) を実装するシステムを示すブロック図である。

【 図 2 】 本発明の別の実施の形態に係る FIOS を実装する別のシステムを示すブロック図である。

【 図 3 】 本発明の実施の形態に係るファイル入出力システム (FIOS) ソフトウェアのブロック図である。

【 図 4 】 本発明の実施の形態に係る FIOS のオペレーションの流れを例示するデータフロー図である。

【 図 5 】 本発明の実施の形態に係る FIOS で使用されるソフトウェアスケジューラの例を示すフロー図である。

【 図 6 】 本発明の実施の形態に係る FIOS におけるプリフェッチの例を説明する模式図である。

【 図 7 A 】 本発明の実施の形態に係る関連づけられたローカルメモリをもつ補助プロセッサを用いた I/O バッファリングの例を説明するフロー図である。

【 図 7 B 】 本発明の実施の形態に係る関連づけられたローカルメモリをもつ補助プロセッサを用いた I/O バッファリングとデータ変換の例を説明するフロー図である。

【 図 8 】 本発明の実施の形態とともに用いられるネストしたアーカイブに対するデータ構造を説明するブロック図である。

【 図 9 】 本発明の実施の形態に係るデアーカイバを説明するフロー図である。

## 【 発明を実施するための形態 】

## 【 0 0 2 3 】

以下の詳細な説明には、説明のための多くの具体的な詳細事項が含まれるが、それらに対する多くの変更や修正も本発明の技術範囲に含まれることは、当業者であれば誰でも理解するであろう。したがって、以下に説明する本発明の例示的な実施の形態は、クレームされた発明の一般性をなんら損なうことなく、また限定も行うことなく記載されている。

## 【 0 0 2 4 】

本発明の実施の形態は、システム用の全ての入出力が通過する集中型のレイヤを提供するファイル入出力システム (file I/O system: FIOS) によって実装されうる。ファイル入出力システムは、ここの入出力リクエストの最適な処理のための木構造を作成する命令と、入出力リクエストをスケジュールし、入出力リクエストが最も効率的にサービスを提供する順序を決定する命令とを含む。

## 【 0 0 2 5 】

例として、図1に示すように、コンピュータで実装されたシステム100が本発明の実施の形態に係るファイル入出力システム(file I/O system: FIOS)を実装するように構成されてもよい。例として、一般性を失うことなく、システム100は本発明の実施の形態の実行に適したパーソナルコンピュータ、ビデオゲームのコンソール、PDA(personal digital assistant)、または他のデジタルデバイスとして実装されてもよい。システム100はCPU105と、処理ユニット105に接続するメインメモリ106とを含んでもよい。CPU105はソフトウェアアプリケーションおよび、オプションとしてオペレーティングシステムを実行するように構成されている。セルプロセッサアーキテクチャの例は、「セルブロードバンドエンジンアーキテクチャ」に詳細が記述されている。これは、IBM、ソニー・コンピュータエンタテインメント、東芝の2005年8月8日付けの著作権で保護されており、コピーは<http://cell.scei.co.jp/>からダウンロード可能であり、その内容全体を参照によってここに組み入れる。

10

**【0026】**

メインメモリ106は、CPU105に使用されるアプリケーションおよびデータを格納する。メインメモリ106は、例えばRAM、DRAM、ROMやそれと同様の集積回路の形態であってもよい。コンピュータプログラム101は、CPU105上で実行可能な命令の形態でメインメモリ106内に格納される。プログラム101の命令は他のもの、すなわち以下に記載する特徴を持ったファイル入出力システム(file input/output system: FIOS)で実装するよう構成されてもよい。メモリ106は、例えばスタックまたはキューの形態で、ファイル入出力プログラム101に使用される受信、スケジューリング、発行、完了、および解放の入出力リクエストである入出力キュー101Qを含む。それらのキューの例もまた以下に記載される。

20

**【0027】**

例として、FIOSプログラム101は、a)メディアデバイス118に関わる入力される(incoming)入出力(I/O)リクエストをアプリケーション103から受信し、b)その入出力リクエストを完了する予測時間 $T_p$ を計算し、c)少なくとも部分的に予測時間 $T_p$ にもとづいてメモリ106内のスケジュールのある位置にその入力された入出力リクエストを挿入し、d)そのスケジュールにしたがって入出力リクエストを実行する命令を含む。

**【0028】**

FIOSプログラム101は、インタラクティブな環境を実装するように構成された命令とともに機能する。一例として、そのような命令は、ビデオゲームプログラムのような、メインプログラム103のサブルーチンまたは呼び出し可能な関数である。あるいは、メインプログラム103は、仮想世界とインタフェースをもつためのプログラムである。メインプログラム103は、ビデオディスプレイ上でカメラPOV(視点)からシミュレートされた環境の一部のシーンを表示し、シミュレートされた環境とのユーザの相互作用の間、カメラ経路に沿ってカメラPOVの動きに応じてカメラPOVが変化するようにシーンを変化させるように構成される。メインプログラムは、物理シミュレーション104とカメラ管理107に対する命令を含む。メインプログラム103は、FIOSプログラム101、物理シミュレーション命令104、カメラ管理命令107、および広告インプレッション報告命令109をたとえば関数またはサブルーチンとして呼び出してもよい。

30

40

**【0029】**

クライアントシステム100はさらに、入出力(I/O)装置111、電源(P/S)112、クロック(CLK)113およびキャッシュ114などの公知のサポート機能510を備えてもよい。クライアントシステム100は、アプリケーションとデータ用の不揮発性のストレージを提供するハードディスクのような高速ストレージ115をさらに備えてもよい。高速ストレージ115は、他のものの中で、より遅いメディアデバイス118から読み出されたファイル116の一時的または長期間のストレージに使用されてもよい。さらに、高速ストレージ115上のファイル116は、より遅いメディアデバイス1

50

16以外の他のソースから来てもよい。限定されるものではないが、例えば、ファイル116はオペレーティングシステムのファイル、アプリケーションによって作成されたテンポラリファイル、写真、音声、映像のようなユーザのデータ、ダウンロードされたコンテンツ、その他を含んでもよい。例として、高速ストレージ115は固定ディスクドライブ、リムーバブルディスクドライブ、フラッシュメモリデバイス、あるいはテープドライブであってもよい。より遅いメディアデバイス118は、例えばCD-ROMドライブ、DVD-ROMドライブ、HD-DVD (high-definition digital versatile disc) ドライブ、Blu-Ray (登録商標) ドライブ、UMDドライブ、または他の光学ストレージデバイス等の大容量光学ディスクドライブであってもよい。メディアデバイス118からのプリフェッチファイル116は、メインメモリ106に素早く読み出されるために、高速ストレージ115のハードウェアキャッシュに一時的に格納されてもよい。

10

**【0030】**

1またはそれ以上の入力デバイス120は、1またはそれ以上のユーザからのユーザ入力をシステム100に伝えるために使われてもよい。例として、1またはそれ以上のユーザ入力デバイス120は入出力エレメント111を経由してクライアントデバイス100と接続されていてもよい。適切な入力デバイス120の例は、キーボード、マウス、ジョイスティック、タッチパッド、タッチスクリーン、光学ペン、静止または動画カメラ、および/またはマイクロフォンを含む。クライアントデバイス100は電子通信ネットワーク127を介してコミュニケーションを促進するためのネットワークインタフェース125を含んでもよい。ネットワークインタフェース125はローカルエリアネットワークおよびインターネットのようなワイドエリアネットワーク上の有線または無線通信を実装するように構成されていてもよい。システム100は、データを送信および受信し、および/またはネットワーク127を介して1またはそれ以上のメッセージパケット126を通じたファイルをリクエストしてもよい。

20

**【0031】**

システム100は、グラフィックプロセッシングユニット(GPU)135とグラフィックスメモリ140とを有するグラフィックス・サブシステム130を更に含む。グラフィックスメモリ140は、出力イメージの各画素に対する画素データを格納するために使われるディスプレイメモリ(例えばフレーム・バッファ)を含む。グラフィックスメモリ140は、GPU135と同じデバイス内に統合されてもよく、別のデバイスとしてGPU135に結合されてもよく、メモリ106内に実装されてもよい。画素データは、CPU105から直接グラフィックスメモリ140に提供される。あるいは、CPU105は、所望の出力イメージを定めるデータやインストラクションをGPU135に提供する。GPU135はそれを用いて1またはそれ以上の出力イメージの画素データを生成する。所望出力イメージを定めるデータやインストラクションは、メモリ106やグラフィックスメモリ140に格納される。ある実施の形態では、GPU135は、例えば、適当なプログラミングまたはハードウェア・コンフィギュレーションによって構成され、シーンのためにジオメトリ、ライティング(照明)、シェーディング、テクスチャリング、モーション、および/またはカメラパラメータを定めるインストラクションとデータから出力イメージに対する画素データを生成するための3Dレンダリング機能をもつ。GPU135は、シェーダプログラムを実行することができるプログラム可能な実行ユニットを更に含む。

30

40

**【0032】**

グラフィックス・サブシステム130は、表示装置150で表示されるべきグラフィックスメモリ140からのイメージに対する画素データを周期的に出力する。表示装置150は、システム100からの信号に応じてビジュアル情報を表示することができる任意のデバイスであり、CRT、LCD、プラズマ、およびOLEDディスプレイを含む。コンピュータシステム100は、表示装置150にアナログまたはデジタル信号を提供する。一例として、表示装置150は、テキスト、数字、グラフィカルシンボル、または画像を表示するCRTまたはフラットパネルスクリーンを含む。さらに、ディスプレイ150は

50

、可聴もしくは検知可能な音を出すオーディオスピーカーを含んでもよい。そのような音の生成を容易にするために、クライアントデバイス100は、CPU105、メモリ106、および/またはストレージ115によって提供されたインストラクションおよび/またはデータからアナログまたはデジタル音声出力を生成するために適したオーディオプロセッサ155を更に含む。

#### 【0033】

CPU105、メモリ106、支援機能110、データストレージデバイス115、ユーザ入力デバイス120、ネットワークインタフェース125、およびオーディオプロセッサ155を含むシステム100の構成要素は、データバス160を介して互いに機能的に接続されている。これらの構成要素は、ハードウェア、ソフトウェアまたはファームウェア、あるいはこれらのうちの2つ以上の組合せで実装される。

10

#### 【0034】

本発明のある実施の形態は、セルプロセッサアーキテクチャやそれと類似するプロセッサアーキテクチャをうまく利用する。図2は、本発明の別の実施の形態に係るファイル入出力システムを実装するセルプロセッサ200を示す。セルプロセッサ200は、メインメモリ202、ひとつのパワープロセッサ要素(power processor element: PPE)204、および8つのシナジスティックプロセッサ要素(synergistic processor element: SPE)206を備える。図4の例では、セルプロセッサ400は、1個のPPE404と8個のSPE406とを備える。このような構成では、7個のSPE406が並列処理に用いられ、1個は他の7個の一つが故障したときのバックアップとして確保されている。例示として、PPE204は、関連するキャッシュを持った64ビットのパワーPCプロセッサユニット(PPU)であってもよい。PPE204は、システム管理リソース(例えば、メモリ保護テーブルなど)にアクセス可能な汎用処理装置である。ハードウェアリソースは、PPEに見られるように実アドレス空間に明示的にマップされてもよい。したがって、PPEは適当な実効アドレス値を用いて、直接これらのリソースのいずれにでもアドレス指定できる。PPE204の主な機能は、セルプロセッサ200のSPE206に対するタスク管理とタスク割り当てである。ハードウェアリソースは、PPE204に見られるように実アドレス空間に明示的にマップされてもよい。したがって、PPE204は適当な実効アドレス値を用いて、直接これらのリソースのいずれにでもアドレス指定できる。PPE204の主な機能は、異なるSPE206に対するタスク管理とタスク割り当てである。PPUはFIOSプログラム205のコード化された命令を実行する。

20

30

#### 【0035】

なんらシステム管理機能を果たさないSPU206は、PPE204よりは簡単な構造のコンピュータ装置である。SPE206のそれぞれは、ときにシナジスティックプロセッサユニット(synergistic processor unit: SPU)として参照されるプロセッサユニットと、関連するローカル記憶領域LS(local store)とを含む。SPU206は、一般にSIMD(single instruction, multiple data)機能を有し、通常はデータ処理を行い、割り当てられたタスクを行うために(PPE204により設定されたアクセス特性にしたがって)要求されたデータ転送を開始する。SPE206はファイル入出力プログラム101の一部を実装するローカルストア命令207を自身のローカルストア内に格納する。SPU206の目的は、より高いコンピュータ装置密度を要求するアプリケーションを可能として、提供された命令セットを有効に利用できるようにすることである。本例では8個のSPEが示されているが、セルプロセッサ200は任意の数のSPEを用いて設定される。図2を参照すると、メモリ202、PPE204、およびSPE206は、リングタイプの要素相互接続バス210を介して、お互いにおよび入出力デバイス208と通信可能である。メモリ202は、メモリインタフェースコントローラ(memory interface controller: MIC)を介してPPE204およびSPE206からアクセスされる。

40

50

## 【0036】

メモリ202は、以下に記載するような受信(incoming)、スケジュール、発行、完了、解放、およびプリフェッチキュー等の入出力キュー203を含む。メモリ202は、本明細書に記載されているFIOSプログラム101と同様の特徴を持っているFIOSプログラム209の部分も含む。SPE206の少なくとも1つは、以下に述べるデータ圧縮解凍を実装するように構成されたコード207をローカルストア内に有する。PPE204は内部キャッシュ(L1)および外部キャッシュ(L2)を含んでもよい。PPEはこの内部キャッシュ(L1)にFIOSプログラム205の部分を格納する。FIOSプログラム205とSPEに実装されたファイル転送命令207、またはその一部とはまた、必要なときにSPE206とPPE204によってアクセスされるためにメモリ202

10

## 【0037】

例として、FIOSプログラム101/205は、ストレージデバイス115のようなハードウェアとのインタフェースを容易にするためにメディアスタックを含んでもよい。メディアスタックは図3に示すように実装される。具体的には、メディアスタック300は、スケジューラ302、1またはそれ以上のメディアフィルタレイヤ304、デバイスメディアレイヤ306、プロセッサファイルシステム(FS)読み込みレイヤ308、およびハードウェアレイヤ310を含んでもよい。デバイスメディアレイヤ306(ときどきメディアアクセスレイヤと呼ばれる)は、メディアフィルタレイヤ304から受け取ったメディアリクエストに回答し、要求されたデータを受け取り、そのリプライをスタック

20

## 【0038】

メディアフィルタレイヤ304は、アーカイバレイヤ301、RAMキャッシュレイヤ303、スケジューラキャッシュレイヤ305、オーバレイレイヤ307、カタログキャッシュレイヤ309、1以上のデータ変換レイヤ311、シミュレートされたデータレイヤ313、および速度エミュレーションレイヤ315を含む。

## 【0039】

アーカイバ301は、圧縮されたアーカイバから特定のASETファイルを抽出することを手助けするのに用いられてもよい。RAMキャッシュレイヤ303は、例えばメインメモリ106、202内にデータをキャッシュすることを実装するために用いられてもよい。スケジューラキャッシュ305は、光学ディスクのようなより遅いソースからのデータを一時的に格納するための、単純なディスクからディスクへの(disc-to-disc)キャッシュであってもよい。本明細書で用いられているように、「スケジューラキャッシュ」という用語は、メディアデバイス118へのアクセスよりも速いアクセスであるストレージメディア内への一時的なデータの格納をいう。スケジューラキャッシュ305内の全てのデータがプリフェッチされる必要はなく、いくらかは要求に応じて読み出されてキャッシュ内にコピーされてもよい。

30

## 【0040】

限定されるわけではないが、例として、スケジューラキャッシュレイヤ305は、そのような一時的な格納場所を生み出すために高速ストレージメディア115を利用する。高速ストレージメディアがハードディスクドライブ(hard disk drive: HDD)である特別なケースでは、スケジューラキャッシュ305はときにハードディスクドライブキャッシュと呼ばれる。

40

## 【0041】

スケジューラキャッシュ305は、単一のファイルまたは多数のファイルとして保持されてもよい。加えて、スケジューラキャッシュ305のコンテンツは、欠ける部分のない全体のファイルである必要はない。多数のファイルのキャッシュは、多くのデータを犠牲にすることなく必要なときにインテリジェントにディスクスペースを空けるために、いくつかの個々のファイルを消去することによって部分的に消去されてもよい。対称的に、単一ファイルのキャッシュは、典型的には切り捨てられるかあるいは完全に消去される。複

50

数ファイルのキャッシュは典型的には追加の帳簿作業 (bookkeeping work) を要求するため、単一ファイルのキャッシュは複数ファイルのキャッシュよりも高いパフォーマンスを提供する。帳簿作業は追加の入出力を要求し、ホストファイルシステム自体の内部で行われる。

#### 【 0 0 4 2 】

オーバーレイレイヤ 3 0 7 は、以下で述べられるように、ファイルシステムレベルにおいてファイルおよびディレクトリの任意のオーバーレイを可能にするために用いられてもよい。カタログキャッシュレイヤ 3 0 9 は、プロセッサユニットまたはセルプロセッサのためのオペレーティングシステムが適切にキャッシュしないであろうデータ (例えばファイルの存在やサイズ、格納場所等) をキャッシュするために用いられてもよい。データ変換レイヤ 3 1 1 は、メディアデバイス 1 1 8 から読み出したり、メディアデバイス 1 1 8 に書き込まれるデータに暗号、暗号解読、圧縮または圧縮伸張のようなデータ変換を実行する。シミュレートされたデータレイヤ 3 1 3 は、たとえば、そのようなシミュレートされたデータを要求する I / O リクエストに対してゼロまたは乱数を生成するために使われる。速度エミュレーションレイヤ 3 1 5 は、たとえばソフトウェアの開発過程で、以下で説明するように、メディアデバイスの異なる速度をシミュレートするために用いられる。

#### 【 0 0 4 3 】

スケジューラ 3 0 2 のオペレーションは、以下に述べられる図 4 を参照することで理解されよう。クライアントアプリケーションの観点から、入出力リクエスト 4 0 3 は比較的直接的な方法で実行される。例えば、ビデオゲームのようなアプリケーション 4 0 1 内のスレッド 4 0 2 は、関数 (例えば、readFile()) を呼び出すことによって F I O S 1 0 1 / 2 0 5 を通して入出力をリクエストしてもよい。この関数は、F I O S リクエスト 4 0 3 およびリクエストが完了されるべき期限を特定してもよい。この関数は、バッファ 4 0 5 へのポインタを提供してもよい。そのような実装は、I / O リクエスト 4 0 3 に対する I / O オペレーションの構造体を自動的に割り当て、新しく割り当てられたオペレーションを受信 (incoming) キュー 4 0 4 に移動させる。例として、受信キューは F I F O (first-in-first-out) キューであり、キュー 4 0 4 内に最初に置かれたオペレーションが、スケジューラ 3 0 2 によって最初にスケジューリングされるオペレーションである。受信キュー 4 0 4 内へのリクエスト 4 0 3 の挿入は、非同期的入出力を要求するときにスレッドがブロックされることを防ぐために、アトミックな操作 (atomic operation) によって実装されてもよい。ある実施の形態では、受信キュー 4 0 4 は、アトミックな交換 (atomic exchange) とともに補助プロセッサ 1 0 5 B にしたがうアトミックなスタック (atomic stack) の形態であってもよい。

#### 【 0 0 4 4 】

限定されないが、例として、入出力リクエスト 4 0 3 が読み込み命令の場合、クライアントアプリケーション 4 0 1 はメディアデバイスから読み込んだデータを受信するためのバッファ 4 0 5 を提供してもよい。入出力リクエストが完了したとき、クライアントアプリケーションは提供したバッファから読み出されたデータを読み込んでよい。クライアントが入出力リクエスト 4 0 3 の処理を終了したとき、クライアントは、将来の入出力の使用においてリソースが利用可能となるように、オペレーションの構造体の割り当てを解除 (de-allocate) してもよい。

#### 【 0 0 4 5 】

入出力リクエスト 4 0 3 は、スケジューラ 3 0 3 がまだ起動していない場合、スケジューラ 3 0 3 を起動してもよい。その後、アプリケーション 4 0 1 は入出力リクエスト 4 0 3 が完了するまで待機してもよい。例えば、クライアントアプリケーション 4 0 1 は、アトミックなメソッド (atomic method、例えば isDone()) を用いて周期的にポーリングしてもよい。もう一つの方法としては、クライアントは、入出力リクエスト 4 0 3 が完了したときにファイル入出力プログラム 1 0 1 が呼び出すコールバック関数 4 1 1 を待ってもよい。

#### 【 0 0 4 6 】

10

20

30

40

50

受信キュー 404 内に入出力リクエスト 403 が挿入された後に、F I O S プログラム 101 / 205 はスケジュール挿入 407 を実行するためにスケジューラ 302 を呼び出してよい。スケジューラによって実装された演算の順序は図 5 を参照することで理解されるであろう。502 で示すように、スケジューラ 302 は、扱うべき入出力がない場合、スリープ（例えば、依然として起動していない）状態であってもよい。新しい入出力リクエスト 403 が受信キュー 404 に入ると、スケジューラはスリープを解除してそれを扱う。ひとたびスケジューラがスリープを解除（あるいは既に解除）すると、受信キュー内の新しいオペレーションに気づいて、そのリクエストをスケジューリングキュー 406 内の適切な場所に移動する。

#### 【0047】

503 に示すように、入出力リクエスト 403 のキューの位置を決定するために、スケジューラ 302 はオプションとして、例えばデバイスメディアレイヤ 306 や F I O S スタック 300 内の他のレイヤに問い合わせることにより、メディアデバイスの現在の状態を決定してもよい。状態データは、F I O S スタック 300 によって扱われるメディアデバイスのタイプおよびスタック内にある様々なレイヤに依存して変わりうる。例としては、もっとも最近にアクセスされた論理ブロックアドレス（logical block address: LBA）、R A M キャッシュの最後にアクセスされたパスおよびオフセット、現在のディスクレイヤ、先頭の位置、ストリーミングのモード、およびストリーミングではないモード（non-streaming mode）が含まれる。

#### 【0048】

本発明の実施の形態によれば、スケジューラ 302 はメディアデバイスのパフォーマンスモデル 302 A に基づいてもよい。既存の入出力スケジューラとは異なり、ドライブモデル（drive model）は、リクエスト 403 の最適なスケジュールを決定する際にメディアデバイス 118 のパフォーマンスに関するデータを考慮に入れる。ドライブモデルを用いて入出力リクエストをスケジューリングすることは、C D を生成する際のマスタ化処理の逆と比較されてもよい。パフォーマンスモデル 302 A は、オーバーヘッド、ディスクの移動時間、読み込み時間、およびそれと同様のもののような要素を考慮するように構成されてもよい。パフォーマンスモデル 302 A は、スループット、レーザの細かな揺れ（wiggles）、読み取りヘッドの動き、レーザの変更、およびリクエストのオーバーヘッドのようなメディアデバイス 118 の任意の複雑な特性をモデル化してもよい。パフォーマンスモデル 302 A はまた、入出力リクエストによって引き起こされたメディアデバイス 118 によって読み込まれたり書き込まれたりする特定のメディアの他のパラメータを考慮に入れてもよい。例えば、パフォーマンスモデル 302 A は、デバイスが単一のレイヤのディスクを読み込んでいるのか、あるいは多数のレイヤのディスク（例えば B l u - r a y（登録商標）D M D のような二重のレイヤのディスク）を読み込んでいるのかを考慮してもよい。

#### 【0049】

スケジューラ 302 はまた任意であるが、符号 504 で示すように、新しい I / O オペレーションの時間要件、優先度、および見込まれる効率性を調べる。符号 506 で示すように、スケジューラは、受信 I / O リクエストを完了する予測時間  $T_p$  を計算するように構成される。好ましい実施の形態によれば、スケジューラ 302 は、メディアデバイス 118 に関するパラメータと I / O リクエスト 403 に関する係数を変数としてもつ性能特性式（performance characteristic equation; P C E）にしたがって、I / O リクエストを実行するための予測時間  $T_p$  を決定してもよい。

#### 【0050】

P C E は次の一般的な形式である。

$$T_p = A f_1(x) + B f_2(y) + C f_3(z) + D f_4(w) + E f_5(v) + \dots$$

ここで、A、B、C、D および E は係数であり、x、y、z、w および v はメディアデバイスの性能特性に関するパラメータを表し、 $f_1(x)$ 、 $f_2(y)$ 、 $f_3(z)$ 、 $f_4$

10

20

30

40

50

( $w$ ) および  $f_5(v)$  は  $x$ 、 $y$ 、 $z$ 、 $w$  および  $v$  の関数である。一例として、限定しないが、 $PCE$  は 1 以上の項の和の形式の線形式であり、各項はパラメータ値 ( $x$ 、 $y$ 、 $z$  など) を表す変数を乗じた係数 ( $A$ 、 $B$ 、 $C$  など) である。そのような式は次の形式である。

$$T_p = Ax + By + Cz + Dw + Ev + \dots$$

【0051】

上記式は任意の数の項目を有してもよい。変数は、要求されたタスクを実行するための時間に影響するメディアデバイスの任意のパラメータに対応する。そのようなパラメータの例には、これに限定しないが、オーバーヘッド、シーク時間、データ転送レートが含まれる。実際上、 $PCE$  は 3 より大きい数の項目、たとえば 10 以上の項目を用いる。たとえば、一つの読み取り (リード) オペレーションを記述するためにはしばしば 6 つの重要な項目が必要である。1. ファイルオープン、2. オーバーヘッド、3. シーク、4. 転送、5. 圧縮解凍、6. ファイルクローズ。他のオペレーション (オープン、クローズ、アンリンク (unlink)、リードディレクトリ (read directory)) には、典型的には異なる項目の集合が必要である。

10

【0052】

線形式の代わりとして、1 以上の項の和であり、各項は何乗かした対応する変数によって乗じた係数 ( $A$ ) である一般的な多項式に  $PCE$  を拡張してもよい。たとえば、 $PCE$  は次の形式をもつ。

$$T_p = Ax + By + Cz^2 + Dw^{10} + Ev^{-2} \dots$$

20

【0053】

別の代案として、1 以上の項の和として表され、各項は対応する変数の任意の関数によって乗じた係数 ( $A$ ) である一般的な式に  $PCE$  を拡張してもよい。たとえば、 $PCE$  は次の形式をもつ。

$$T_p = Ax + B \log y + Ce^z + D(w^2 + 2w + 1) + \dots$$

【0054】

単純な例として、線形  $PCE$  は  $T_p = Ax + By + Cz$  の形式であり、ここで、 $x$  はリクエスト 403 に対するオーバーヘッドに関するパラメータであり、 $y$  はセクタ単位のヘッドの動きに関するパラメータであり、 $z$  はデータスループットレート (たとえば、所与のデータ量を読み出すまたは書き込む時間) に関するパラメータである。係数  $A$ 、 $B$  および  $C$  の値は、スケジュールキュー 404 内のリクエスト 403 の位置に依存する。

30

【0055】

一例として、係数  $A$ 、 $B$ 、 $C$  は、状態 (ステート) だけでなくリクエスト 403 内の情報からも決定される。たとえば、一般性を失うことなく、その情報には、デバイスのレディ状態とアクセスされる最後の論理ブロックアドレス (LBA) が含まれる。最後の論理ブロックアドレスは、ドライブの読み取り (リード) ヘッドの位置を示し、これにより、シーク距離 (ヘッド移動距離) を計算することが可能になる。

【0056】

ある例では、係数  $A$  はオーバーヘッド係数と呼ばれる。この係数は、オーバーヘッドが部分的には、I/O リクエスト 403 の性質だけでなく、メディアデバイス 118 の特性にも依存するという事実を考慮に入れている。係数  $B$  は、データが読み取られるべきメディア上の位置に到達するために、メディアデバイス 118 上の読み取りヘッドが取らなければならないヘッドの運動量に関する。たとえば、リクエスト 403 は、メディアデバイス 118 がセクタ 1000 からスタートする 10 セクタを読むように要求する。そのような場合において、データ係数  $C$  は読み取るべきセクタ数である。運動係数  $B$  は、デバイスメディアレイヤ 306 から取得される状態 (状態) 情報から決定される。状態情報には、I/O リクエストの開始時における読み取りヘッドの初期位置が含まれる。たとえば読み取りヘッドが初期的にセクタ 500 にあるなら、運動係数  $B$  は、読み取りが開始するべきセクタ (セクタ 1000) と現在のセクタ (セクタ 500) の差、たとえば  $B = 1000 - 500 = 500$  から決定される。スケジューラ 302 は、I/O リクエスト 4

40

50

03を完了した後で最終ステートを決定する。たとえば、10セクタが1000番から読み取られるなら、最終ステータは、読み取りヘッドがセクタ1010で終わることを示す。

【0057】

PCEはまたI/Oリクエストを書き込むためのありうる異なるシナリオを考慮に入れてもよい。たとえば、メディアデバイス118が複数のレイヤ媒体を用いるなら、PCEは、あるレイヤから別のレイヤに切り替えて読み取りヘッドを動かす予測時間 $T_p$ に与える影響を考慮に入れる。

【0058】

ある実施の形態において、任意のI/Oリクエストに対する現実にかかった時間 $T_a$ を測定し、PCEを反復して解いて未知変数の値を求めるためにPCE係数とともに用いてもよい。これにより、モデルを現実のドライブ性能で絶えず更新することができる。

10

【0059】

いったん予測時間 $T_p$ が計算されると、スケジューラ302は同様の方法でスケジュールキュー406における他のリクエストに対する予測時間を決定する。この処理は、符号507に示すように、リクエスト403に対する現実の時間と最適なスケジュールキューの順序が決定されるまで、複数の異なるキュー順序で反復して繰り返される。リクエスト403に対する初期状態は、スケジュール順位における以前のリクエストに対する最終状態から決定され、スケジュール順位における次のリクエストに対する初期状態は、リクエスト403の最終状態から決定される。

20

【0060】

スケジューラ302は、スケジューリングキュー内のリクエストを検証して、既にキュー内にあるリクエストに対してリクエスト403の特性を比較する。スケジューラ302は、キュー内におけるリクエスト403の考え得る各位置を試し、優先度を無視していること、期限を過ぎること、およびタイミングの判断を探す。スケジューラ302は、リクエストされた入出力オペレーションがいずれも期限を過ぎない順序を探すことによって、考え得る最適な新しいキューの順序を決定してもよい。もし1またはそれ以上のオペレーションが絶対に期限を途過するという場合、スケジューラは異なる基準を用いてキューの順序を決定してもよい。スケジュールキュー406内のリクエストの順序に影響を与える追加の基準は、限定はされないが、入出力リクエストの優先度、ストリームバッファの状態、および要求されるレイテンシを含んでもよい。ある実施の形態では、 $T_p$ 計算に対するPCE係数は、別のPCE係数計算に対する初期メディア状態として使うことができる更新されたメディア状態を返すために使われる。

30

【0061】

一例によれば、スケジューラ302は、異なる挿入順列(permutation)を調べる。たとえば、オペレーション1、2、3を含む既存スケジュールにオペレーション4を挿入する。各順列は異なる合計時間 $T$ をもたらす。ありうる順列と時間には以下がある。

- 4, 1, 2, 3  $T_1$
- 1, 4, 2, 3  $T_2$
- 1, 2, 4, 3  $T_3$
- 1, 2, 3, 4  $T_4$

40

【0062】

スケジューラ302は、どの順列が最小予測合計時間 $T$ を生成するかを決め、この順列を好ましいものとして使う。

【0063】

さらに、スケジューラ302は、たとえば、スケジュールされた各アイテムに対する時間を加算することにより、徒過が予測される期限を探し、これにより、いずれかのスケジュールされたアイテムが自分の期限を徒過することになるかどうかを調べる。これには、スケジュールされた各オペレーションに対する完了時間を決定し、その完了時間をそのオペレーションに対する期限と比較する。いずれかの期限が徒過されるならば、オペレーシ

50

ョンはキュー内のより早い位置に再スケジュールされる。

【 0 0 6 4 】

ある実施の形態では、例えば期限が全て解決できない場合、優先度が使用される。例えば、ある期限が否応なく徒過する場合、最良のキューの順序は、最も優先度の低いリクエストがその期限を徒過する順序であろう。前述の考察を満たす多数の可能なキューの順序がある場合、同じ優先度である期限を徒過するリクエストの数が最小となるキューの順序が最良の順序であろう。前述した考察を全て満たす多数の可能なキューの順序がある場合、可能な限り少ない時間で全てのキューを実行するキューの順序が最良の順序であろう。いくつかの場合、高い優先度のリクエストが期限に間に合う限りにおいて、低い優先度のリクエストは高い優先度のリクエストの前にスケジューリングされてもよい。

10

【 0 0 6 5 】

前述した考察を全て等しく満たす多数の可能なキューの順序がある場合、スケジュールキュー 4 0 6 内の最新のリクエストはキューの末尾に来る順序が最良の順序であろう。スケジューラ 3 0 2 がキューに新しいオペレーションを挿入するために合計キュー実行時間を見積もらなければならない場合は、多数の異なるファクタを考慮する。たとえば、スケジューラは、キュー順序の評価の直前にスケジューラ 3 0 2 が取り出したデバイスドライバの現在の状態を考慮に入れる。それに加えて、スケジューラ 3 0 2 は、メディアフィルタレイヤと F I O S スタック 3 0 0 における他のレイヤのよって報告された性能係数を考慮する。これらの値には、たとえば、ヘッドの動き、ヘッドアライメント、DMA セットアップ、データ転送レート、ストレージレイヤに対するレイヤ変更などが含まれる。これらの値は、I / O 性能のランタイム測定によって供給されるので、見積もりを可能な限り真実なものにすることができる。

20

【 0 0 6 6 】

ひとたびスケジューラ 3 0 2 がリクエスト 4 0 3 におけるオペレーションに対してキュー 4 0 6 内の最適な位置を決定すると、図 4 の 4 0 7 および図 5 の 5 0 8 に示されるように、オペレーションはそこに挿入される。図 4 を再び参照すると、もし入出力リクエストを実行可能なリソースが存在すれば、4 0 9 に示されるように、スケジューラ 3 0 2 はスケジュールキュー 4 0 6 内の最初のオペレーションを発行キュー 4 0 8 に移動する。オペレーションの関連づけられたメディアリクエストは発行キュー 4 0 8 から実行される。メディアリクエストを実行するために、スケジューラ 3 0 2 は、F I O S スタック 3 0 0 内の最初のレイヤにリクエストを渡す。

30

【 0 0 6 7 】

スケジューラ 3 0 2 以下の各レイヤは、それよりも上のレイヤから渡されたメディアリクエスト 4 0 3 を見る。適宜、レイヤはそのデータを処理する。例えば、リクエストが読み出し (read) オペレーションの場合、デアーカイバレイヤ 3 0 1 は開かれているアーカイブのコンテンツに対して与えられたパスネームを確認し、もしファイルを見つけた場合、圧縮されたデータの読み出し (read) にそのリクエストを再マッピングする。F I O S スタック 3 0 0 の各レイヤは、リクエストが最終的にハードウェアレイヤ 3 1 0 に到達するまで、処理されたあるいは処理されていないメディアリクエストを次の下位のレイヤに渡す。ハードウェアレイヤ 3 1 0 がリクエストに応答するとき、その応答はスタック 3 0 0 内の各レイヤを下から上に向かって進み、各レイヤは適宜読み出したデータを処理する。例えば、デアーカイバレイヤは返されたデータは解凍されなければならないことを知っているため、応答をスタックに返す前にデータを解凍する。応答は最終的にスケジューラ 3 0 2 に戻されて完了する。

40

【 0 0 6 8 】

読み出されたデータがスタックに戻るとき、スケジューラ 3 0 2 がそれを受信して入出力オペレーションを完了キューに移動する。完了キューはアプリケーション 4 0 1 へのコールバック関数の契機となる (コールバックがアプリケーションによって設定されていた場合)。別の方法として、アプリケーション 4 0 1 は入出力オペレーションが完了したか否かを決定するために F I O S スタック 3 0 0 をポーリングしてもよい。ひとたび入出力

50

オペレーションが完了すると、フリーオペレーションプール412に移動される。フリーオペレーションプール412は使用されていない入出力オペレーション構造体の集合を含む。そのようなオペレーションには、クライアントアプリケーションに割り当てられていなかったりものや、クライアントアプリケーションによって使用された後、再び使用するために解放されたものが含まれる。クライアントアプリケーションが入出力リクエストを作ると、スケジューラ302はこのプールからクライアントに入出力オペレーションを割り当てる。フリー入出力プール412はスタックとして実装されてもよい。フリー入出力リクエストはフリー入出力プール412からポップされ、受信キュー404にプッシュされる。このような方法で、フリー入出力リクエストは再利用される。

#### 【0069】

10

本発明の実施の形態によれば、スケジューラ302は以下のようにスケジュールループを動作する。

1. 入出力の完了を調べる。
2. 新しい入出力リクエストを発行する。
3. 入出力コールバックを発行する（もしあれば）。
4. スケジュールに挿入する（受信からスケジュールに挿入する）。
5. 再度新規に発行する。
6. 予想される期限の途過を調べる。
7. 1に戻る。

#### 【0070】

20

各反復におけるスケジュールの挿入の数は、ある最大数（例えば16）で制限されていてもよい。

#### 【0071】

本発明のある実施の形態では、補助プロセッサユニット105B（例えばSPE206）が、受信する入出力用のリクエストを、メインプロセッサ105AまたはPPE204に関してアトミックな交換を通じて受信キュー404に追加することによって、自分自身に入出力をリクエストすることができる。例えば、従来型のセルプロセッサの実装において、SPE206はいずれの入出力能力も持っていない。しかしながら、受信キュー404が共通のデータ要素である場合、補助プロセッサ105Bは、メインプロセッサ105との標準的なアトミックな交換およびメインプロセッサ105に通知することを通じて、

30

入出力リクエストをキューに追加することができる。

#### 【0072】

多くの従来技術における実装では、補助プロセッサ105Bが即時に処理するデータを必要とする場合、そのデータはメインメモリになければならない。これとは対称的に、本発明の実施の形態では、補助プロセッサ105Bはファイル入出力システムメディアスタック300に、ハードメディアデバイス118やネットワーク127を介して得た所望のデータを取りに行かせる。

#### 【0073】

限定するわけではないが、例として、受信キュー404、完了キュー410、およびフリー（解放）プール412はアトミックなスタックであってもよい。補助プロセッサ105Bはフリープール412からオペレーションをポップし、オペレーションのパスとオフセットを書き込み、受信キュー406にオペレーションをプッシュし、さらに同期を実行し、スケジューラ302をウェイクアップする。補助プロセッサ105Bは入出力オペレーションの完了のための間欠的なポーリングの間に別の仕事を行ってもよい。

40

#### 【0074】

例として、セルプロセッサシステム200内において、入出力リクエストはSPE206によって、次のタイプの命令シーケンスを用いてサービスされる。

```
Op = pop(FREE)
Op path = .....
Offset = .....
```

50

Push (op, incoming)

【 0 0 7 5 】

S P E 2 0 6 によってリクエストされたデータは、P P E 2 0 4 によってアドレス可能な場所であればどこでも送信される。S P E 2 0 6 がロックされている場合、F I O S 3 0 0 は S P E 2 0 6 のローカルストア ( local store: LS ) にデータを直接書き込むことができる。ある実施の形態では、ひとつの S P E 2 0 6 はデアーカイバに他の S P E を用いて圧縮解凍するようリクエストしてもよい。これは例えばプロセッサ間の遠隔手続呼び出し ( remote procedure call: RPC ) を用いて実装されてもよい。この例では、S P E 2 0 6 が P P E 2 0 4 に何かをするように要求する。より従来型の遠隔手続呼び出しではこれと逆である。

10

【 0 0 7 6 】

[ プリフェッチ ]

ある実施の形態では、スケジューラキャッシュ 3 0 5 ( 例えばハードディスクキャッシュ ) は、ファイルが後に素早く読み出せるように、メディアデバイスからのファイルをストレージデバイス 1 1 5 内のキャッシュ 1 1 7 にプリフェッチするために用いられる。そのような場合、プリフェッチオペレーション 4 1 3 は発行キュー 4 0 2 に直接挿入されてもよい。プリフェッチは、P o w e r P C の C P U アーキテクチャにおける「dcbt」命令のような多くの種類のキャッシュによって実装される標準的な手動プリフェッチであってもよい。本発明の実施の形態によれば、プリフェッチは上述のスケジューラループの一部としてキューに入れられ、かつ実行されてもよい。プリフェッチは、比較的高いレイテンシでかつ低いスループットを伴う、他の入出力リクエストが完了しているときにアクセスされなければならない比較的遅いソースメディア ( 例えば B l u - r a y ( 登録商標 ) や U M D ) の仕事の手助けをする。

20

【 0 0 7 7 】

スケジューラ 3 0 2 は、メディアデバイス 1 1 8 がアイドル状態で、かつ他の入出力リクエストを邪魔しない、使用されていない時期にのみ実行するような、比較的低い優先度のプリフェッチを実装するように構成されてもよい。例えば、システム 3 0 0 は遅いメディアデバイス ( 例えば B l u - r a y ( 登録商用 ) ディスク ( B D ) のような光学ディスク ) と、ハードディスクのようなより高速なストレージデバイス 1 1 5 とを持っているとする。スケジューラキャッシュレイヤ 3 0 5 は、光学ディスクドライブからハードディスクへの非同期なコピーのために用いられてもよい。ファイルが後にアクセスされる時、入り低速な光学ディスクの速度 ( 例えば 8 M i B / s ) の代わりに、より高速なハードディスクドライブの速度 ( 例えば 2 0 M i B / S ) で読み込まれる。

30

【 0 0 7 8 】

キャッシュプリフェッチは通常ハードディスクキャッシュで行われるが、プリフェッチはメインメモリ 1 0 6、2 0 2 で行われてもよい。プリフェッチオペレーションは、スケジューリングが完了した後に受信された順序でスケジューリングキュー 4 0 6 の末尾に含まれてもよい。ある実施の形態では、プリフェッチオペレーションは、必要なときに遅延されてもよい。ある実施の形態では、スケジューラ 3 0 2 によって使用される性能特性式には、状態情報の一部として、キャッシュ 1 1 7 に格納されるものについての情報が含まれる。

40

【 0 0 7 9 】

プリフェッチオペレーションは他の入出力オペレーションと同様の方法ではスケジューリングされない。プリフェッチオペレーションは、スケジューリングキュー 4 0 6 が空のときのみ実行される別個のプリフェッチキュー 4 1 3 に保管される。プリフェッチオペレーションは必要があれば遅延されてもよい。そうでなければ、プリフェッチは受信された順序で実行してもよい。スケジューラ 3 0 2 はキューされたプリフェッチリクエストを保持し、入出力サブシステムが特定の長さの時間だけアイドル状態であるときにのみそれらのプリフェッチリクエストを実行する。これにより、プリフェッチが通常の入出力リクエストを邪魔するのを防ぐ。加えて、プリフェッチは単一のキャッシュブロックには限定され

50

ず、任意のサイズでよく、キャッシュにファイルのはじめから最後まで全体をロードすることを伝える特別な「全ファイル (whole-file)」値を渡してもよい。さらに、プリフェッチは任意のサイズでよいが、プリフェッチは、F I O S スタック 3 0 0 がアイドル状態を継続しているかを調べるためにスケジューラ 3 0 2 に戻る前には、せいぜいひとつのキャッシュブロックしか満たされないように実装されてもよい。

#### 【 0 0 8 0 】

F I O S スタック 3 0 0 は、不必要なプリフェッチを回避する方法でプリフェッチを実装するように構成される。たとえば、図 6 に示すように、どんなブロックもキャッシュ 6 0 1 にないならば、F I O S スケジューラキャッシュレイヤ 3 0 5 は、最初のブロックを埋めて、アイドル状態をチェックするために戻る。もしストレージデバイス 1 1 5 上のキャッシュ 1 1 7 にブロック 1 から N があるならば、( N + 1 ) 番目のブロックがメディアデバイス 1 1 8 からのデータで埋められる。スケジューラキャッシュレイヤ 3 0 5 はアイドル状態をチェックするために戻る。このプロセスは受信キューに新しい I / O が到着するまで繰り返される。これらの機能は任意の数のプリフェッチがキューされることを可能にするため、通常のスケジュールされた I / O リクエストを妨げることなく、アプリケーションがメディアデバイス 1 1 8 からストレージデバイス 1 1 5 (またはメインメモリ 1 0 6) へ大きな量のデータを高速かつ効率的にプリフェッチすることができるようになる。

#### 【 0 0 8 1 】

たとえば、ゲームアプリケーションにおいて、ゲームデータはしばしば、メディアデバイス 1 1 8 によって読まれる C D - R O M または H D - D V D のようなメディアに格納されている。ゲームは、プレイヤーが特定の行き先に動いていることを判定し、60 秒もすればそこに到着するであろうと計算する。その時間の間、F I O S スタック 3 0 0 は、メディアデバイスからその行き先に関してプリフェッチする。プレイヤーが向きを変えて異なる行き先に向かうことを決定したなら、アプリケーションはそのプリフェッチをキャンセルする。他の実施の形態では、アプリケーションは、キャッシュ 1 1 7 にある必要とされるデータが上書きされないように防ぐために事前通知を利用する。

#### 【 0 0 8 2 】

上述のキャッシュプリフェッチは、ビデオゲームのような入出力駆動型 ( I / O driven ) のアプリケーションや、現代的なビデオゲームコンソールのような特定の入出力を要するプラットフォームにおける入出力のパフォーマンスを向上しうる。特に、ゲームデータはしばしば光学ディスクやネットワークファイルサーバのような低速なメディアに格納されるが、ゲームはハードディスクドライブのような高速なローカルストレージへのアクセスをもってもよい。

#### 【 0 0 8 3 】

[ 補助プロセッサを用いたデータ変換 ]

別の実施の形態によれば、F I O S は、メインプロセッサと、関連づけられたローカルメモリをもつ補助プロセッサとを利用するあるプロセッサアーキテクチャを活用する。たとえば、圧縮伸張または暗号解読のようなデータ変換のためにそのような実施の形態が用いられる。比較のために、ある以前のアーカイブ解凍システムは、標準的なバックグラウンドスレッドの中で動作する ( z l i b のような ) 単純な圧縮伸張アルゴリズムを用いてきた。本発明の実施の形態は、それとは対照的に、標準的なデコンプレッサ ( decompressor ) 実装と、セル ( C e l l ) プロセッサアーキテクチャのような関連づけられたローカルメモリをもつ補助プロセッサを含むプロセッサアーキテクチャを活用するように構成された非標準的な実装の両方を可能にする抽象的なコーデックインタフェースを含む。そのような実施の形態によれば、デコンプレッサオブジェクトが多数のフラグと制約条件を設定することが可能になり、デアーカイバ ( de-archiver ) 3 0 1 は、それを活用して、メモリ使用と速度を最適化する。特別な取り扱いをすることになるこれらの制約条件と環境の多くは、ゲームプログラミング、ファイルシステム I / O 圧縮伸張、またはシステム 1 0 0 のアーキテクチャにとって固有のものである。

## 【 0 0 8 4 】

ビデオゲームアプリケーションのようなあるアプリケーションは、限られた量の利用可能なメモリをもち、可能な限り少ないメモリを用いることが望ましい。一つの共通の定式化は、I/O圧縮伸張のために3つのI/Oバッファをもつことである。一つのバッファは読み取りによって使われ、2つまでの一時バッファ（一つは入力用、もう一つは出力用）はデータ変換オペレーション（たとえば圧縮伸張や暗号解読）によって使われる。一例としてこれに限定しないが、ある実施の形態は、データ変換のために関連付けられたローカルメモリ（たとえばセルプロセッサ200におけるSPE206）をもつ補助プロセッサを用いることにより、一時バッファ全体を使用することを回避してもよい。

## 【 0 0 8 5 】

たとえば、圧縮伸張がメインプロセッサ（たとえば、PPE204）によって扱われるなら、典型的には、圧縮されたデータを一時的に格納するために一つのバッファが必要であり、伸張されたデータを一時的に格納するために別のバッファが必要である。十分なローカルメモリ106Bをもつ補助プロセッサ105Bを利用できるなら、一時的なバッファの代わりにローカルメモリ106Bを用いて同じ圧縮伸張を行うことができる。圧縮されたデータと伸張されたデータの両方をローカルメモリ106Bに格納してもよい。一例として、あるセルプロセッサアーキテクチャにおいて、SPE206は256キロバイトの容量のローカルストアを有する。これは、圧縮された入力、伸張された出力、および圧縮伸張アルゴリズムに対するコードのために十分なメモリスペースである。比較すれば、圧縮伸張がメインプロセッサ105A（たとえば、PPE204）によって扱われるなら、圧縮されたデータは第1バッファに読み込まれ、伸張され、伸張されたデータは第2バッファに書き込まれる。その後、伸張されたデータは、第2バッファから宛先へコピーされる。

## 【 0 0 8 6 】

図7Aおよび7Bは、補助プロセッサを利用してデータ変換を実装するものの中で2つの例を示す。具体的には、図7Aにおいて説明されるように、I/Oリクエストを遂行する過程で、メインプロセッサ105A上で動作するFIOS300が、たとえば、データ変換レイヤ311を介してデータ変換を起動する。符号702で示すように、メインプロセッサ105Aは補助プロセッサ105Bをセットアップする。具体的には、メインプロセッサ105Aは補助プロセッサ105Bに命じて、補助プロセッサのローカルメモリ106Bにデータ変換を埋め込むためのコード化された命令701をロードさせる。符号704で示すように、その後、メインプロセッサまたは補助プロセッサはFIOSスタック300を検査して、メディアデバイス118からメインメモリ106にデータ703を転送する。補助プロセッサ105Bはその後、符号706に示すようにデータ703をローカルメモリ106Bにロードする

## 【 0 0 8 7 】

次に補助プロセッサは、符号708で示すように、データ703に変換を実行し、変換されたデータ705を生成し、これをローカルメモリ106Bに格納する。一例として、コード化された命令701の実行により、データ703は暗号解読または圧縮伸張される。変換されたデータ705は、符号710で示すように、メインメモリ106に転送される。あるいは、変換されたデータ705は、他の宛先、たとえばストレージデバイス115に配信される。I/Oリクエストによって特定されるすべてのデータが処理されるまで、オペレーションの前述のシーケンスをデータの次のブロックに対して繰り返してもよい。いったん入力データ703がローカルメモリ106Bに転送されると、そのデータは、新しい入力データまたは変換されたデータ705によって、メインメモリにおいて上書きされる。

## 【 0 0 8 8 】

一例として、これに限定しないが、図2で図示したようなセルプロセッサアーキテクチャを用いて、図7Aのデータ圧縮伸張を実装してもよい。SPUプログラミングの特徴（データがSPE206のローカルストアLSにストリーム伝送され、その後、いったん処

10

20

30

40

50

理が完了するとストリーム返送される)は、ファイルシステム I/O 圧縮伸張の要求と組み合わせられ、I/O 性能を改良する機会を与える。たとえば、圧縮伸張のために S P E 2 0 6 を用いることにより、付加的なバッファを解放することの結果、性能が 3 7 M b p s から 4 2 M b p s に改善する。ソニープレイステーション 2 において使用されるプロセッサアーキテクチャを用いて類似した圧縮伸張を実装してもよい。そのような場合、メインプロセッサはエモーションエンジン ( E E ) と呼ばれ、補助プロセッサの一つは、関連づけられたローカルメモリをもつ I/O プロセッサ ( I O P ) として知られる。そのような場合、デアークイバ 3 0 1 は E E 上で動作し、圧縮伸張は、図 7 A に図示された方法にしたがって I O P 上で実行される。あるいは、図 7 A に図示されて技術は、独立してアドレス可能なメモリをもつ 2 つのプロセッサが関わる分散コンピューティングアーキテクチャに対して使用することもできる。

10

## 【 0 0 8 9 】

前述の例は、変換レイヤによって実装されたデータ変換を用いているが、図 7 A に関して説明した方法は、メディアフィルタレイヤ 3 0 4 のいずれかを実装する際に用いることもできることに留意する。それとは対照的に、先行技術のファイル I/O システムは典型的には、メインプロセッサ (たとえば、セルプロセッサの例では P P U ) で実装される。

## 【 0 0 9 0 】

代替的な実装では、メディアサービス 1 1 8 から読まれたデータの宛先が低速メモリ (たとえば、ビデオランダムアクセスメモリ ( V R A M ) のようなグラフィックスメモリ 1 4 0 ) であり、コーデックがその出力にランダムアクセスを実行する必要がある場合、F I O S デアークイバは、一時的な出力バッファを割り当てる必要がある。しかし、コーデックがランダムアクセスする必要がないなら (この例で言えば、たとえば S P U 実装されたデコンプレッサの場合)、F I O S デアークイバ 3 0 1 は、一時的な出力バッファを割り当てることを避けて、データ 7 0 3 を圧縮伸張して直接、宛先に送ることができる。さらに、ファイルシステム圧縮伸張には「スロップ (slop)」すなわち未使用バイトが必要である。すなわち、6 4 キロバイト全体を圧縮するとき、ほんのわずかなバイトが必要とされる。コーデックが、圧縮されたブロックの一部を出力することができる (この例では、S P U 実装されたデコンプレッサを用いて) ことを示しているなら、F I O S デアークイバは一時的な出力バッファを割り当てることを回避できる。それに加えて、関連づけられたローカルメモリ 1 0 6 B をもつ補助プロセッサ 1 0 5 B を上述の環境の任意のすべての組み合わせで用いることができる。

20

30

## 【 0 0 9 1 】

図 7 B は、補助プロセッサ 1 0 5 B とローカルメモリ 1 0 6 B を低速グラフィックスメモリ 3 4 0 とともに使用する例を示す。F I O S 3 0 0 はメディアデバイス 1 1 8 からグラフィックスメモリ 3 4 0 (たとえば、V R A M) ヘデータを読み出す。バッファを用いた V R A M へのデータ転送は低速になる傾向がある。しかし、V R A M データをあるタイプの補助プロセッサのローカルメモリ (たとえば、S P E 2 0 6 のローカルストア L S ) に読み込むことは問題ではない。図 7 B に示す例では、F I O S 3 0 0 は主としてメインプロセッサ 1 0 5 A 上で実装される。未変換データ 7 0 3 はグラフィックスメモリ 3 4 0 から、補助プロセッサ 1 0 5 B に関連づけられたローカルメモリ 1 0 6 B に転送される。F I O S 3 0 1 は、補助プロセッサ 1 0 5 B によって実行するために変換命令 7 0 1 をローカルメモリ 1 0 6 B にロードする。補助プロセッサ 1 0 5 B は、命令 7 0 1 を実行することによってデータ 7 0 3 を変換されたデータに変換し、変換されたデータ 7 0 5 を生成する。変換されたデータはその後、グラフィックスメモリ 3 4 0 に転送される。この実施の形態のバリエーションの中で、未変換 (たとえば暗号化または圧縮された) データ 7 0 3 は、メディアデバイス 1 1 8 から読み出され、ローカルメモリ 1 0 6 B に転送される前にメインメモリ 1 0 6 のバッファに格納される。

40

## 【 0 0 9 2 】

## [ ネストしたアーカイブ ]

以前の F I O S 実装には、ランダムアクセスデアークイブが含まれており、それによ

50

て、圧縮されたアーカイブのコンテンツがファイルシステム階層に仮想的に現れ、読み出しと圧縮伸張をオーバーラップさせることによって最適に効率化した方法でこれらのアーカイブからデータを読み出すことができた。本発明の実施の形態は、単一のデアーカイバレイヤ301をもち、性能ペナルティのないネストした(入れ子になった)アーカイブを取り扱う能力をもつことにより、このシステムを改善するものである。ここで使われるネストしたアーカイブという用語は、他のアーカイブ内に格納されるアーカイブのことを言い、他のアーカイブはさらにネストの任意の深いレベルにまで別のアーカイブに格納されていてもよい。

#### 【0093】

ネストしたアーカイブの支援は、ゲームI/Oおよびデータレイアウトの固有の要求を満たすために特に有益である。ネストしたアーカイブの使用は、スケジューラ302によって使用されるドライブモデル302Aを用いることと相乗作用がある。具体的には、大きなファイルは、ドライブヘッドの動きのより良いモデルを考慮に入れる。ネストしたアーカイブの支援は、圧縮されたアーカイブフォーマットがランダムアクセスであるならば、大いに役立つ。それとは対照的に、典型的には圧縮されたアーカイブはランダムアクセスではない。たとえば、zip、gzip、7z等の圧縮フォーマットは典型的には、アーカイブの始まりから線形アクセスすることが必要である。

#### 【0094】

ここで「アーカイブ」という用語は、単一の大きなファイルに結合された複数のデータファイルの集合のことである。ソースデータは変更されることなく格納されるか、圧縮または暗号化のような1以上のデータ変換がデータの任意の部分に適用される。インデックス(ときには目次と呼ばれる)は、アーカイブ内のパスおよびオフセットのような、ソースファイルの少なくともある一面を記述するものであり、アーカイブとともに維持されることで、ファイルを個別に調べてアクセスすることができるようになる。「ネストしたアーカイブ」という用語は、1以上のソースデータファイルが別のアーカイブ(またはネストしたアーカイブ)であるようなアーカイブのことである。このネスト構造のおかげで、ソースデータの所与のブロックは、それに適用された1以上のデータ変換を有する。たとえば、2つのレベル(レベル1とレベル2)をもつビデオゲームアプリケーションは、「level1.psarc」と「level2.psarc」と名付けられた2つのアーカイブを含む「game.psarc」と名付けられたネストしたアーカイブを用いる。これは、ゲーム開発者がデバッグやテストの目的で単一のアーカイブ(「level1.psarc」や「level2.psarc」など)を用いることを可能にし、いつも大きな「game.psarc」というアーカイブを生成するために時間を費やさなくてもよくなることで開発を容易にする。

#### 【0095】

開発者はネストしたアーカイブを生成することができるが、従来のI/Oリーダーは、複数のレベルをもつアーカイブを読むことができないか、効率的には読むことができない。これを克服するために、アーカイブは図8に示すように構成される。特に、外側のアーカイブ800は、ゲームの個別のレベルに対するデータをもったファイルが含まれる内側のアーカイブ801、802を含む。一例として、これに限られないが、内側のアーカイブ801、802は、そのレベルで利用可能な地図、武器やそのレベルで出会う敵のような特定のゲームレベルに関連するデータを含む。内側のアーカイブ801、802に対するデータは、個別に圧縮され、共通するデータ803とともにバンドルされる。外側のアーカイブ800内のファイルは、ブロック単位で圧縮される。各アーカイブは、デアーカイバが、内側のアーカイブ801、802における任意の場所を含めて、外側のアーカイブ800内の任意の場所の任意のデータブロックを見つけることを可能にするインデックスを含む。インデックスは、アーカイブに対する目次と等価なものとみなすことができる。

#### 【0096】

FIOS300がアーカイブにある特定のファイルに対するリクエストを発行するとき、リクエストはRequest:/game/level2/map(1000,10)のような形になる。「/game/level2/map」の部分はパスであり、特定のファイルを見つけるために使われる。「(1000,10)」

10

20

30

40

50

の部分はリクエストのオフセットと長さを特定する。外側のアーカイブ 800 またはその一部が読み出されるとき、デアーカイバ 301 は利用可能なファイルのリストを考慮する。このリストは、どのアーカイブを開いているかに依存して異なる。たとえば、外側のアーカイブ 800 も内側のアーカイブ 801、802 も開いていないなら、リストは次のようである。

```
/game.psarc
```

外側のアーカイブ 800 を開いているなら、リストは次のようである。

```
/game (folder)
/game/level1.psarc
/game/level2.psarc
/game/common.dat
```

10

内側のアーカイブを開いているなら、リストは次のようである。

```
/game/level1/map
/game/level1/weapons
/game/level1/enemies
/game/level2/map
/game/level2/weapons
/game/level2/enemies
/game/common.dat
```

#### 【0097】

20

各アーカイブ 800、801、802 は、コンテンツのパス（または部分的なパス）のハッシュのリストをもつ。ハッシュリストは、アーカイブに対する目次（table of contents; TOC）を生成するために使われる。たとえば、外側のアーカイブ 800 に対する目次は次のように構成される。

```
TOC:      hash (/level1.psarc)
           hash (/level2.psarc)
           hash (/common.dat)
```

ここで、hash() は括弧内の量のハッシュのことである。任意の適切なハッシュ、たとえば MD5 ハッシュを用いることができる。

#### 【0098】

30

同様に、レベル 1 の内側のアーカイブ 801 に対する目次は次のように構成される。

```
TOC:      hash (/map)
           hash (/weapons)
           hash (/enemies)
```

#### 【0099】

図 8 に示されたタイプのネストしたアーカイブが使われる場合、デアーカイバ 311 は、それらを取り扱うように構成される。一例として、デアーカイバ 311 は図 9 に示したアーカイブ解除方法 900 を実装する。この方法 900 は、符号 902 に示すように、アーカイブにおける初期ルックアップで始まる。一例として、これに限らないが、初期ルックアップは、探しているデータに対する初期パス、オフセット、長さ、および要求されるデータ変換の順序リストを返す。具体的には、I/O リクエストは次のようである。

```
(/game/level2/weapons, offset 0, length 50000, none)
```

このファイルに対する初期ルックアップは次を返す。

```
(/game/level2.psarc, offset 20000, length 30000, zlib decompress)
```

#### 【0100】

初期ルックアップの後、デアーカイバ 311 は、符号 904 で示すようにネストを解消する。ネストを解消する際、デアーカイバ 311 は、アーカイブ 802 が別のアーカイブ内に格納されているかどうか、それは圧縮されて格納されているか、圧縮されずに格納されているかを判定しようとする。

#### 【0101】

50

ネストを解消する際、デアーカイバ311は、探しているデータに対するパス、オフセット、長さ、および要求されるデータ変換の順序リストを返す。たとえば、パス/game.parc、オフセット120000、長さ30016、暗号解読+zlib圧縮伸張のような形式をもつ。スタックの制限を避けるために再帰よりもループでネストを解消する。ネストがいったん解消されると、メディアデバイス118からデータが転送される。一例として、これに限られないが、デアーカイバ311は、符号905で示すようにブロック単位でデータを繰り返し読み出して変換する。具体的には、デアーカイバは、符号906で示すようにデータのブロックを読み出し、その後、符号908で示すように読み出されたブロックにデータ変換を適用する。ネストしたアーカイブの複数のレベルが単一のインタリーブした読み取り-変換ループに解消され、もっとも外側のアーカイブからブロックを読み出すことで、ソースメディアを最大限効率的に利用する。各ブロックが読み取られ、変換された後、符号910に示すように、望ましい部分が適当な宛先、たとえば、メインメモリ106、ストレージデバイス115または他の場所にコピーされる。ブロックの望ましい部分は任意のサイズであり、たとえば、1バイトのように小さくても、ブロック全体であってもよい。

10

## 【0102】

## [オーバーレイレイヤ]

上述のように、FIOSスタック300におけるメディアフィルタレイヤ304は、オーバーレイレイヤ307を含む。本発明の実施の形態によれば、オーバーレイレイヤ307は、ファイルシステムレベルにおいてファイルとディレクトリの任意のオーバーレイができるように構成される。一般的なコンセプトは、より一層柔軟性があるという点を除けば、Unix(登録商標)のユニオン(union)マウントと比較することができる。典型的なユニオンマウントは、あるファイルシステムにおけるディレクトリを別のファイルシステムの全内容とマージすることにより実現される。オーバーレイレイヤ307は同様のことをするが、もっと柔軟であり、たくさんの特徴がある。たとえば、オーバーレイレイヤ307は、ディレクトリ、ファイルまたはファイル内のバイトのレンジの粒度で動作する。ディレクトリは隠され、そのコンテンツはオーバーレイレイヤ307を用いて置き換えられる。さらに、オーバーレイレイヤ307は、プライマリデータソースから要求されるファイルをセカンダリデータソースからの他のファイルで置き換えるために使われる。セカンダリデータソースは、ファイルまたはディレクトリであり、それは、プライマリデータソースと同じメディアまたは別のメディア上にある。さらに、第2データソースは仮想的なデータソースであり、データがコールバックによって指定されてもよい。この特定の文脈において、「コールバック」という用語は、他のコードへの引数として渡される実行可能なコードのことである。

20

30

## 【0103】

さらに、オーバーレイを生成するアプリケーションは、オーバーレイがオーバーレイレイヤ307によってどのようにして、いつ適用されるべきかの基準を特定する。たとえば、アプリケーション103が、ディレクトリBをディレクトリA上にオーバーレイしているならば、基準は、「オーバーレイ」=最初にB内のファイルを探す、「アンダーレイ」=最初にA内のファイルを探し、次にB内のファイルを探す、「より新しい」=両方の場所にファイルが存在するならば、より最新の更新日をもつ方を使う、「より古い」=両方の場所にファイルが存在するならば、より古い更新日をもつ方を使うなどである。

40

## 【0104】

オーバーレイレイヤ307は、特定のソースから別のソースへのI/Oリクエストに対するリクエストの任意のマッピングを可能にするように構成される。たとえば、メディアデバイス118から特定のファイルを読み出すリクエストは、ストレージデバイス115などにある対応するファイルにマップされる。ファイル内の特定のバイトから全体のディレクトリまでの範囲のそのような特徴データを用いていることは、隠されているか、入れ替えられている。どのディレクトリ、ファイル、またはバイトを置き換える必要があるかを特定することにより、オーバーレイレイヤ307は、特定のI/Oリクエストを処理するために

50

F I O Sスタック300によってなされなければならない仕事量を最小化することができる。

【0105】

オーバーレイレイヤ307は、たとえば、ソフトウェアおよび/またはデータに対するアップデートを実装するために使われる。多くの先行技術のアップデートシステムでは、実行可能なリンク可能ファイル(executable linkable file; ELF)がメディアデバイス118からストレージデバイス115にコピーされ、更新されたコードまたはデータでパッチが当てられる。本発明の実施の形態を用いれば、これとは対照的に、アプリケーション103が、メディアデバイス118からのデータや、ストレージデバイス115のような他の場所からの他のバイトレンジを用いて、あるバイトレンジを埋めるようにF I O S 101に指示する。したがって、パッチを当てるために全体のファイルをコピーする必要がない。その代わりに、パッチを当てることは、I/Oシステムレベルで実装される。

10

【0106】

たいていの先行するパッチの実装は、ファイル全体を置き換えることにもとづいている。これは実装するにはもっとも単純な方法であり、したがって、もっとも普通の方法によるものである。オーバーレイレイヤ307によって実装される強化したパッチは、それとは対照的に、圧縮され、かつ、暗号化されたアーカイブ内でパッチを当てるといった、より複雑で有益なパターンを可能にする。

【0107】

[速度エミュレーションレイヤ]

本発明のある実施の形態では、速度エミュレーションレイヤ315は、低速I/Oデバイスの性能をエミュレートするために開発過程でより速いストレージの速度を低下させるために用いられる。速度エミュレーションレイヤ315は、低速I/Oデバイスのシーク時間、スループット、および他の性能特性のモデルを利用する。たとえば、速度エミュレーションレイヤ315は、ストレージデバイス115(たとえばハードディスクドライブ)を低速化して、HD DVDやブルーレイ(登録商標)ドライブのようなメディアデバイス118の低速のシークタイムおよび/またはスループットをエミュレートする。

20

【0108】

本発明の実施の形態にしたがって、速度エミュレーションレイヤ315は、以前のセッションの間に集めたドライブ性能データを将来のセッションに与えることによって実装される。上述のように、スケジューラ302によって使用される性能モデルは、スループット、レーザの小刻みな動き、リードヘッドの動き、レイヤの変更、およびリクエストのオーバーヘッドのような複雑なドライブの特徴を適宜モデル化するように構成される。速度エミュレーションレイヤ315はこの情報を用いて、特定のメディアデバイスの高度に正確なエミュレーションを提供する。このエミュレーションは、ランタイム性能にもとづくため、シーク時間のようなより複雑な詳細を無視してスループットを単純に合わせようとする単純かつ素朴なエミュレーションレイヤよりもずっと正確である。

30

【0109】

本発明の実施の形態は、かなりの量のI/Oを利用するアプリケーションとシステムにおける改良されたI/O性能を提供する。上述のように、本発明の実施の形態は、ビデオゲームアプリケーションおよびビデオゲームシステムにおいて特に有益である。しかしながら、本発明の実施の形態はそのようなアプリケーションやシステムに限定されるものではない。

40

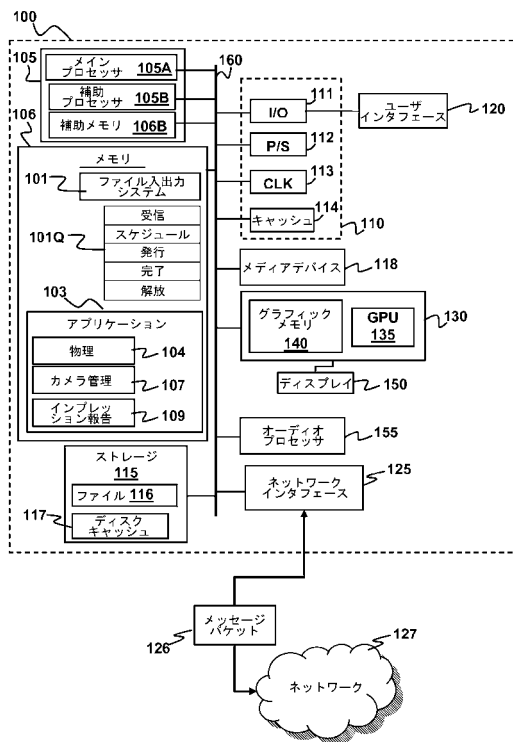
【0110】

本発明の好ましい実施の形態を完全な形で説明してきたが、いろいろな代替物、変形、等価物を用いることができる。したがって、本発明の範囲は、上記の説明を参照して決められるものではなく、請求項により決められるべきであり、均等物の全範囲も含まれる。ここで述べた特徴はいずれも、好ましいかどうかを問わず、他の特徴と組み合わせてもよい。請求項において、明示的に断らない限り、各項目は1またはそれ以上の数量である。請求項において「~のための手段」のような語句を用いて明示的に記載する場合を除いて

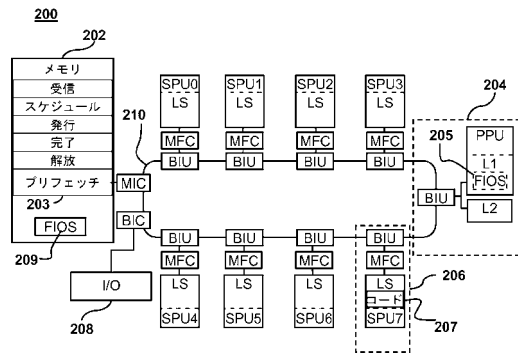
50

、請求項がミーンズ・プラス・ファンクションの限定を含むものと解してはならない。

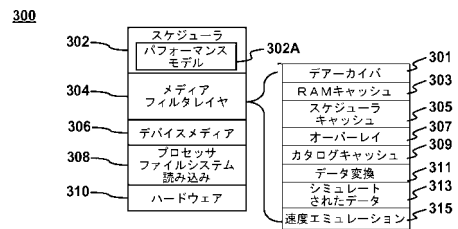
【図1】



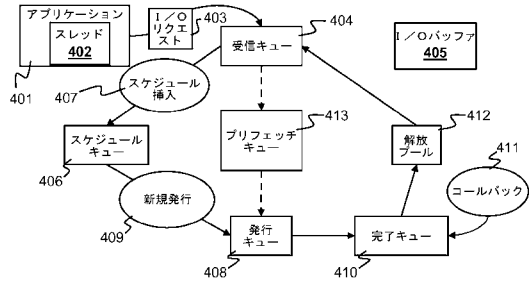
【図2】



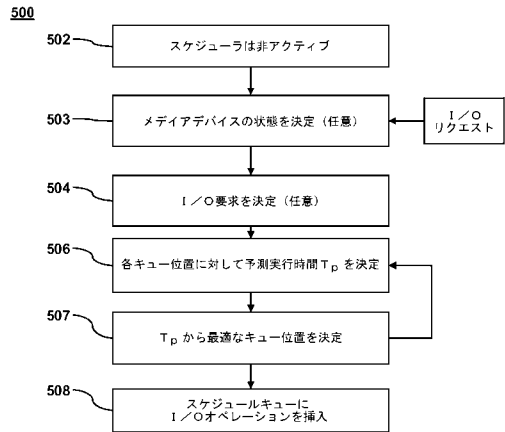
【図3】



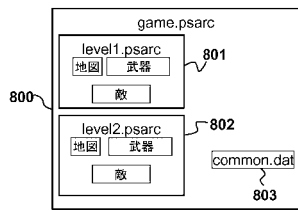
【図4】



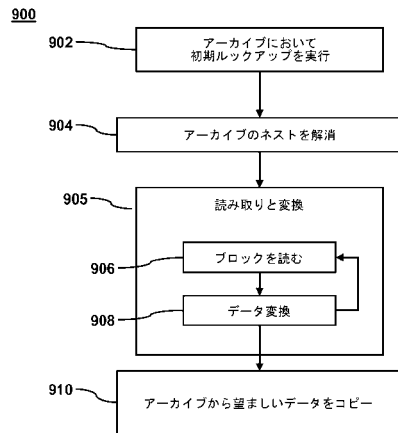
【図5】



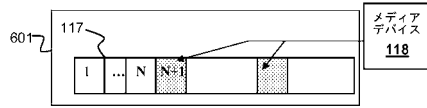
【図8】



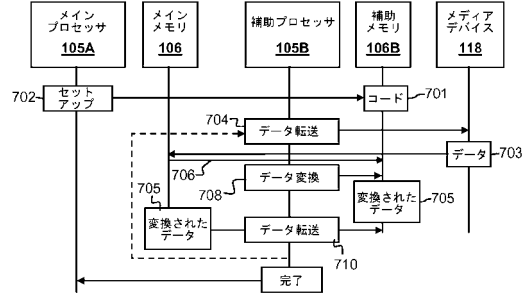
【図9】



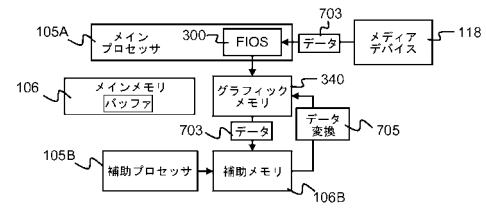
【図6】



【図7A】



【図7B】



---

フロントページの続き

- (72)発明者 ラーナー、エドワード、アダム  
アメリカ合衆国、カリフォルニア州 94404、フォスター・シティー、セカンド・フロアー、  
イースト・ヒルスデイル・ブルバード 919
- (72)発明者 ミカル、ロバート、ジェイ .  
アメリカ合衆国、カリフォルニア州 94404、フォスター・シティー、セカンド・フロアー、  
イースト・ヒルスデイル・ブルバード 919

審査官 篠塚 隆

- (56)参考文献 特開2008-65478(JP,A)  
特開平8-212090(JP,A)  
特開2001-344073(JP,A)

- (58)調査した分野(Int.Cl., DB名)  
G06F9/48  
3/06  
13/10