

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6734768号
(P6734768)

(45) 発行日 令和2年8月5日 (2020. 8. 5)

(24) 登録日 令和2年7月14日 (2020. 7. 14)

(51) Int. Cl. F I
GO 6 F 16/18 (2019. 01) GO 6 F 16/18
GO 6 F 12/16 (2006. 01) GO 6 F 12/16
GO 6 F 16/11 (2019. 01) GO 6 F 16/11

請求項の数 21 (全 23 頁)

(21) 出願番号	特願2016-240058 (P2016-240058)	(73) 特許権者	390019839
(22) 出願日	平成28年12月12日 (2016. 12. 12)		三星電子株式会社
(65) 公開番号	特開2017-120626 (P2017-120626A)		S a m s u n g E l e c t r o n i c s
(43) 公開日	平成29年7月6日 (2017. 7. 6)		C o . , L t d .
審査請求日	令和1年12月11日 (2019. 12. 11)		大韓民国京畿道水原市靈通区三星路129
(31) 優先権主張番号	62/273, 323		129, S a m s u n g - r o , Y e o n
(32) 優先日	平成27年12月30日 (2015. 12. 30)		g t o n g - g u , S u w o n - s i , G
(33) 優先権主張国・地域又は機関	米国 (US)		y e o n g g i - d o , R e p u b l i c
(31) 優先権主張番号	15/089, 237	(74) 代理人	110000051
(32) 優先日	平成28年4月1日 (2016. 4. 1)		特許業務法人共生国際特許事務所
(33) 優先権主張国・地域又は機関	米国 (US)	(72) 発明者	ヒュオウ, ジアンジアン
早期審査対象出願			アメリカ合衆国, カリフォルニア州 95
			120, サンジョゼ, オールド・ミル・コ
			ート, 6602
			最終頁に続く

(54) 【発明の名称】 二重書込みを遂行するストレージ装置を含むシステム、装置、及びその方法

(57) 【特許請求の範囲】

【請求項 1】

ストレージ（格納）装置（d e v i c e）を含むシステムであって、
 プロセッサ及びメモリを含むコンピュータと、
 前記ストレージ装置と、
 前記プロセッサ上で実行され、ジャーナル（j o u r n a l）書込み要請を前記ストレージ装置に伝送し、データ書込み要請をデータストレージシステムに伝送するように動作するアプリケーションと、

第2ジャーナル書込み要請を前記ストレージ装置に伝送し、第2データ書込み要請を前記ストレージ装置に伝送するように動作する前記データストレージシステムと、

前記ストレージ装置上に位置し、前記ストレージ装置が第1ストリーム（s t r e a m）に割り当てられた第1ブロック（b l o c k）にジャーナル情報を書き込み、第2ストリームに割り当てられた第2ブロックにデータを書き込み、第3ストリームに割り当てられた第3ブロックに第2ジャーナル情報を書き込むように指示するコントローラと、を含み、

前記ジャーナル書込み要請は、前記ジャーナル情報を含み、前記第1ストリームに割り当てられ、

前記データ書込み要請は、前記データを含み、前記第2ストリームに割り当てられ、

前記アプリケーションは、前記データ書込み要請を前記プロセッサ上で実行される前記データストレージシステムに伝送するように動作し、

10

20

前記第 2 ジャーナル書込み要請は、前記第 2 ジャーナル情報を含み、第 3 ストリームに割り当てられ、

前記第 2 データ書込み要請は、前記データを含み、前記第 2 ストリームに割り当てられ、

前記第 1 ストリーム、第 2 ストリーム、第 3 ストリームは、データ特性によって定義されることを特徴とするシステム。

【請求項 2】

前記コントローラは、前記データ書込み要請が完遂された以後に前記ジャーナル情報を削除するために無効化 (i n v a l i d a t e) 要請を受信するように動作する、ことを特徴とする請求項 1 に記載のシステム。

10

【請求項 3】

前記アプリケーションは、前記無効化要請を伝送するように動作する、ことを特徴とする請求項 2 に記載のシステム。

【請求項 4】

前記アプリケーションは、前記データ書込み要請が完遂されたという信号を前記アプリケーションが受信したことに応答して、前記無効化要請を伝送するように動作する、ことを特徴とする請求項 3 に記載のシステム。

【請求項 5】

前記データストレージシステムは、第 2 無効化要請を伝送するようにさらに動作する、ことを特徴とする請求項 1 に記載のシステム。

20

【請求項 6】

前記データストレージシステムは、前記第 2 データ書込み要請が完了した後、前記ジャーナル情報を削除するために前記第 2 無効化要請を伝送するようにさらに動作する、ことを特徴とする請求項 5 に記載のシステム。

【請求項 7】

前記ジャーナル書込み要請は、直接入出力 (I / O) 要請として伝送され、

前記データ書込み要請は、バッファリングされた入出力 (I / O) 要請として伝送され、

前記ジャーナル書込み要請は、前記データ書込み要請の前記データを前記ストレージ装置に確実に書き込むために用いられることを特徴とする請求項 1 に記載のシステム。

30

【請求項 8】

前記第 2 データ書込み要請に含まれるデータは、前記データ書込み要請に含まれるデータであることを特徴とする請求項 1 に記載のシステム。

【請求項 9】

前記ジャーナル書込み要請は、前記データ書込み要請の前記データを前記ストレージ装置に確実に書き込むために用いられることを特徴とする請求項 1 に記載のシステム。

【請求項 10】

前記第 2 ジャーナル情報は、前記ジャーナル情報と同一であることを特徴とする請求項 1 に記載のシステム。

【請求項 11】

40

前記ジャーナル書込み要請は、前記アプリケーションによって前記第 1 ストリームに割り当てられ、

前記データ書込み要請は、前記アプリケーションによって前記第 2 ストリームに割り当てられることを特徴とする請求項 1 に記載のシステム。

【請求項 12】

前記第 1 ブロック及び前記第 2 ブロックは、単一の媒体 (m e d i a) タイプであることを特徴とする請求項 1 に記載のシステム。

【請求項 13】

データストレージシステムが、ジャーナル書込み及びデータ書込みの両方を遂行するアプリケーションから書き込まれるデータを識別する段階と、

50

前記データストレージシステムが、前記アプリケーションから無効データに対するガーベッジコレクション (garbage collection) を遂行するストレージ装置に、ジャーナル書込み要請を送送する段階と、

前記データストレージシステムが、前記アプリケーションから前記ストレージ装置にデータ書込み要請を送送する段階と、

前記データストレージシステムが、前記アプリケーションから前記データストレージシステムに前記データ書込み要請を送送する段階と、

前記データストレージシステムが、前記データストレージシステムから前記ストレージ装置に第2ジャーナル書込み要請を送送する段階と、

前記データストレージシステムが、前記データストレージシステムから前記ストレージ装置に第2データ書込み要請を送送する段階と、を有し、

前記ジャーナル書込み要請は、第1ストリームに割り当てられ、直接入出力 (I/O) 要請として送られ、

前記アプリケーションから前記ストレージ装置に送られたか、前記アプリケーションから前記データストレージシステムに送られた前記データ書込み要請は、前記アプリケーションから書き込まれる前記データを含み、第2ストリームに割り当てられ、バッファリングされた入出力 (I/O) 要請として送られ、

前記第2ジャーナル書込み要請は、第3ストリームに割り当てられ、

前記第2データ書込み要請は、前記アプリケーションから書き込まれる前記データを含み、前記第2ストリームに割り当てられ、

前記ジャーナル書込み要請及び前記第2ジャーナル書込み要請は、前記アプリケーションから前記ストレージ装置に送られたか、前記アプリケーションから前記データストレージシステムに送られた前記データ書込み要請の前記アプリケーションから書き込まれる前記データを前記ストレージ装置に確実に書き込むために用いられることを特徴とする二重書込み方法。

【請求項14】

前記データストレージシステムが、前記データ書込み要請が前記ストレージ装置に書き込まれた以後に、ジャーナル情報を削除するために無効化要請を前記ストレージ装置に送送する段階をさらに含む、ことを特徴とする請求項13に記載の二重書込み方法。

【請求項15】

前記無効化要請を前記ストレージ装置に送送する段階は、前記無効化要請を前記アプリケーションが前記ストレージ装置に送送する段階を含む、ことを特徴とする請求項14に記載の二重書込み方法。

【請求項16】

前記無効化要請を前記アプリケーションが前記ストレージ装置に送送する段階は、前記ストレージ装置上の前記データ書込み要請が完遂されたという信号を前記アプリケーションで受信する段階を含む、ことを特徴とする請求項15に記載の二重書込み方法。

【請求項17】

前記第2データ書込み要請が前記ストレージ装置に書き込まれた以後に、前記データストレージシステムが、前記第2ジャーナル書込み要請によって書き込まれた前記データを削除するために無効化要請を前記ストレージ装置に送送する段階をさらに含む、ことを特徴とする請求項13に記載の二重書込み方法。

【請求項18】

非一時的な (Non-Transitory) 命令を格納した有形 (tangible) のストレージ媒体 (tangible storage medium) を含む装置 (article) において、

前記非一時的な命令がマシンによって実行される時、

ジャーナル書込み及びデータ書込みの両方を遂行するアプリケーションから書き込まれるデータを識別する段階と、

無効データに対するガーベッジコレクションを遂行するストレージ装置に、前記アプリ

10

20

30

40

50

ケーションからジャーナル書込み要請を伝送する段階と、

前記アプリケーションからデータストレージシステムにデータ書込み要請を伝送する段階と、

前記データストレージシステムから前記ストレージ装置に、第2ジャーナル書込み要請を伝送する段階と、

前記データストレージシステムから前記ストレージ装置に、第2データ書込み要請を伝送する段階と、を遂行し、

前記ジャーナル書込み要請は、第1ストリームに割り当てられ、直接入出力(I/O)要請として伝送され、

前記データ書込み要請は、第2ストリームに割り当てられ、バッファリングされた入出力(I/O)要請として伝送され、

前記第2ジャーナル書込み要請は、第3ストリームに割り当てられ、

前記第2データ書込み要請は、前記データを含み、前記第2ストリームに割り当てられ、

前記ジャーナル書込み要請及び前記第2ジャーナル書込み要請は、前記データ書込み要請の前記データを前記ストレージ装置に確実に書き込むために用いられることを特徴とする装置。

【請求項19】

前記非一時的な命令を格納した有形のストレージ媒体を含む装置において、前記非一時的な命令がマシンによって実行される時、

前記データ書込み要請が前記ストレージ装置に書き込まれた以後に、ジャーナル情報を削除するために無効化要請を前記ストレージ装置に伝送する段階をさらに遂行する、ことを特徴とする請求項18に記載の装置。

【請求項20】

前記無効化要請を前記ストレージ装置に伝送する段階は、前記無効化要請を前記アプリケーションから前記ストレージ装置に伝送する段階を含む、ことを特徴とする請求項19に記載の装置。

【請求項21】

前記非一時的な命令を格納した有形のストレージ媒体を含む装置において、前記非一時的な命令がマシンによって実行される時、

前記第2データ書込み要請が前記ストレージ装置に書き込まれた以後に、前記第2ジャーナル書込み要請によって書き込まれた前記データを削除するために無効化要請を前記ストレージ装置に伝送する段階をさらに実行する、ことを特徴とする請求項18に記載の装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明はストレージ(格納)装置(device)を含むシステムに関し、さらに具体的には、二重書込みを遂行するストレージ装置を含むシステム、装置(article)及びその方法に関する。

【背景技術】

【0002】

NANDフラッシュメモリベースのSSD(Solid-State Drive)は、全ての種類のストレージ(格納)装置を含むシステムの速度を向上させるために事業用(enterprise)サーバとデータセンタ内で広く使用されて来た。SSD内のフラッシュメモリは独特の特性を有する。従って、従来の磁気ディスク(magnetic disk)をSSDにより直接置換してもSSDに最大能力を発揮させて活用することにならない。

その重要な理由の1つは、SSDがフリー(free)フラッシュメモリブロックのみ

10

20

30

40

50

に書き込むことができ、再使用するためには、ガーベッジコレクション (garbage collection) を遂行して無効とされたフラッシュメモリブロックを回復する、という点である。従来のOS (オペレーティングシステム) 及びアプリケーションがホット (hot) データとコールド (cold) データとを区別しないので、相異なる寿命 (lifespan) を有するデータが通常、混合されている場合、フラッシュメモリを管理し、リクレイム (reclaim、再利用) するためのガーベッジコレクションの遂行をさらに難しくする。これはSSDの性能及び寿命の全てに影響を及ぼす。

【0003】

最近の多くのデータストレージシステムは、データ耐久性及び性能のためにジャーナリング (journaling) を使用する。データストレージシステムはオブジェクト (object) ストレージシステム (例、Ceph等)、ブロックストレージシステム (例、キャッシュ及び他のキャッシングシステム等)、そしてファイルストレージシステム (例、IBM JFS/JFS2、Linux (登録商標) _xfs、及びLinux (登録商標) _ext4等) を含む。このようなシステムは2つのデータコピーを格納する。1つのコピーはジャーナルセクション (journal section) に格納され、他の1つのコピーはデータセクションに格納される。

10

【0004】

このようなシステムが純粋なSSD環境内に展開される場合、性能及び費用上の理由から、システムは一般的にジャーナルと実際データを同一のSSD上に格納する。データが書き込まれるために受信されると、データストレージシステムはまず1つのデータコピーをディスクにフラッシュ (flush、一斉書き込み) されるジャーナルに格納し、2番目のデータコピーをメモリのファイルシステムページキャッシュ内に格納する。そうすると、データストレージシステムはユーザアプリケーションにサクセス (後継動作、success) を返却 (return) する。続いて後ほど、バックグラウンド (background) 内で、データストレージシステムは、ファイルシステムページキャッシュ内のデータ記録をディスクにフラッシュ (flush、一斉書き込み) し、ディスク上のジャーナル内の同一のデータ記録を除去する。このような過程はデータ書き込みの各々に対して反復され、またジャーナルがメタデータのみに対して使用される場合にも生じられる。

20

【0005】

このような二重書き込み (double-write) 方式 (approach) は、SSDと共に使用される場合、全てのSSDブロック内のフラッシュメモリの内部の断片化 (fragmentation) という問題を惹起する。このような内部断片化問題はガーベッジコレクションの増大を引き起こして、ストレージシステムの性能を低下し、読み出し/書き込みレイテンシ (latency、待ち時間) を増加し、SSDの寿命を短縮する。

30

【0006】

そこで、SSD内のフラッシュメモリの断片化を回避できるか、或いは少なくとも最少化できるような、SSDに関する二重書き込みアプローチの利用方法が要求される。

【先行技術文献】

40

【特許文献】

【0007】

【特許文献1】米国特許第7,610,442号公報

【特許文献2】米国特許第8,738,882号公報

【特許文献3】米国特許第8,793,531号公報

【特許文献4】米国特許公開第2014/0297918号明細書

【特許文献5】国際特許公開第WO2015/126518号明細書

【発明の概要】

【発明が解決しようとする課題】

【0008】

50

従って本発明の目的は、SSD内のフラッシュメモリの断片化を回避、乃至最少化した二重書込みを遂行するストレージ装置を含むシステム、装置及びその方法を提供することにある。

【課題を解決するための手段】

【0009】

上記目的を達成するためになされた本発明によるシステムは、ストレージ（格納）装置（device）を含むシステムであって、プロセッサ及びメモリを含むコンピュータと、ストレージ装置と、前記プロセッサ上で実行され、ジャーナル（journal）書込み要請をストレージ装置に伝送し、データ書込み要請をデータストレージシステムに伝送するように動作するアプリケーションと、第2ジャーナル書込み要請を前記ストレージ装置に伝送し、第2データ書込み要請を前記ストレージ装置に伝送するように動作する前記データストレージシステムと、前記ストレージ装置上に位置し、前記ストレージ装置が第1ストリーム（stream）に割り当てられた第1ブロック（block）にジャーナル情報を書き込み、第2ストリームに割り当てられた第2ブロックにデータを書き込み、第3ストリームに割り当てられた第3ブロックに第2ジャーナル情報を書き込むように指示するコントローラと、を含み、前記ジャーナル書込み要請は、前記ジャーナル情報を含み、前記第1ストリームに割り当てられ、前記データ書込み要請は、前記データを含み、前記第2ストリームに割り当てられ、前記アプリケーションは、前記データ書込み要請を前記プロセッサ上で実行されるデータストレージシステムに伝送するように動作し、前記第2ジャーナル書込み要請は、第2ジャーナル情報を含み、第3ストリームに割り当てられ、前記第2データ書込み要請は、前記データを含み、前記第2ストリームに割り当てられ、前記第1ストリーム、第2ストリーム、第3ストリームは、データ特性によって定義されることを特徴とする。

【0010】

上記目的を達成するためになされた本発明による二重書込み方法は、ジャーナル書込み及びデータ書込みの両方を遂行するアプリケーションから書き込まれるデータを識別する段階と、前記アプリケーションから無効データに対するガーベッジコレクション（garbage collection）を遂行するストレージ装置に、ジャーナル書込み要請を伝送する段階と、前記アプリケーションから前記ストレージ装置にデータ書込み要請を伝送する段階と、前記アプリケーションからデータストレージシステムにデータ書込み要請を伝送する段階と、前記データストレージシステムから前記ストレージ装置に第2ジャーナル書込み要請を伝送する段階と、前記データストレージシステムから前記ストレージ装置に第2データ書込み要請を伝送する段階と、を有し、前記ジャーナル書込み要請は、第1ストリームに割り当てられ、直接入出力（I/O）要請として伝送され、前記データ書込み要請は、前記データを含み、前記第2ストリームに割り当てられ、バッファリングされた入出力（I/O）要請として伝送され、前記第2ジャーナル書込み要請は、第3ストリームに割り当てられ、前記第2データ書込み要請は、前記データを含み、前記第2ストリームに割り当てられ、前記ジャーナル書込み要請及び前記第2ジャーナル書込み要請は、前記データ書込み要請の前記データを前記ストレージ装置に確実に書き込むために用いられることを特徴とする。

【0011】

上記目的を達成するためになされた本発明による装置は、非一時的な（Non-Volatile）命令を格納した有形（tangible）のストレージ媒体（tangible storage medium）を含む装置（article）において、前記非一時的な命令がマシンによって実行される時、ジャーナル書込み及びデータ書込みの両方を遂行するアプリケーションから書き込まれるデータを識別する段階と、無効データに対するガーベッジコレクションを遂行するストレージ装置に、前記アプリケーションからジャーナル書込み要請を伝送する段階と、前記アプリケーションからデータストレージシステムにデータ書込み要請を伝送する段階と、前記データストレージシステムから前記ストレージ装置に、第2ジャーナル書込み要請を伝送する段階と、前記データストレージ

10

20

30

40

50

システムから前記ストレージ装置に、第２データ書込み要請を伝送する段階と、を遂行し、前記ジャーナル書込み要請は、第１ストリームに割り当てられ、直接入出力（Ｉ／Ｏ）要請として伝送され、前記データ書込み要請は、第２ストリームに割り当てられ、バッファリングされた入出力（Ｉ／Ｏ）要請として伝送され、前記第２ジャーナル書込み要請は、第３ストリームに割り当てられ、前記第２データ書込み要請は、前記データを含み、前記第２ストリームに割り当てられ、前記ジャーナル書込み要請及び前記第２ジャーナル書込み要請は、前記データ書込み要請の前記データを前記ストレージ装置に確実に書き込むために用いられることを特徴とする。

【発明の効果】

10

【００１２】

本発明の実施形態によれば、ストレージ装置内の断片化されたブロックの誘発を防止しながら、二重書込み動作が可能である。これに、ストレージ装置の動作効率が増加する。これに、ストレージ装置の動作効率が増加する。

【図面の簡単な説明】

【００１３】

【図１】本発明の実施形態によるジャーナリングを備えたデータストレージを使用するサーバを示すブロック図である。

【図２】図１のサーバの詳細に示す図面である。

【図３】従来技術に係り、ジャーナリング及びデータ書込みの両方を遂行するために図１のストレージ装置と通信する図１のアプリケーションを示す図面である。

20

【図４】ジャーナル及びデータを格納するためにマルチストリーミングを使用する図１のストレージ装置を示す図面である。

【図５】従来のアプローチを使用する図１のストレージ装置の使用法を示す図面である。

【図６】本発明の実施形態によるアプローチを使用する図１のストレージ装置の使用法を示す図面である。

【図７】ジャーナリング及びデータ書込みの両方を遂行するために図１のストレージ装置と通信する図１のアプリケーション及び図１のデータストレージシステムを示す図面である。

30

【図８】ジャーナリング及びデータ書込みの両方を遂行するために図１のストレージ装置と通信する図１のアプリケーション及び図１のデータストレージシステムを示す図面である。

【図９】本発明の実施形態によってジャーナリング及びデータ書込みの両方を遂行し、図１のストレージ装置と通信するための図１のアプリケーション及び図１のデータストレージシステムに対する手順を例示的に示すフローチャートである。

【図１０】本発明の実施形態によってジャーナリング及びデータ書込みの両方を遂行し、図１のストレージ装置と通信するための図１のアプリケーション及び図１のデータストレージシステムに対する手順を例示的に示すフローチャートである。

【発明を実施するための形態】

40

【００１４】

以下では、本発明の技術分野で通常の知識を有する者が本発明を容易に実施できる程度に、本発明の実施形態を明確に、且つ詳細に記述する。

【００１５】

本明細書では第１番目、第２番目等のような用語が多様な構成要素を表現するために使用されたが、このような構成要素がこのような用語によって限定されないことは容易に理解されよう。このような用語は単に、１つの構成要素を他の構成要素と区分するために使用されており、例えば、本発明の範囲を逸脱せずに、第１番目のモジュールは第２番目のモジュールと称され得るし、同様に、第２番目のモジュールは第１番目のモジュールと称され得る。

50

【0016】

図1は本発明の実施形態によるジャーナリングを備えたデータストレージシステムを使用するサーバを示すブロック図である。図1に示されたサーバ105は任意の種類のサーバである。サーバ105はプロセッサ110及びメモリ115を含む。プロセッサ110は任意の種類のプロセッサである。また、メモリ115は任意の種類のメモリである。

【0017】

プロセッサ110上で動作できるのはデータストレージシステム120である。データストレージシステム120は二重書込み（即ち、ジャーナリング及びデータ書込み）を遂行する任意の種類のシステムである。データストレージシステム120は‘Ceph’（登録商標）のようなオブジェクトストレージシステム及びファイルストレージシステムのみでなく、アプリケーションが二重書込みを遂行する他のオペレーティングシステム上で動作するアプリケーションを含む。（‘Ceph’は、Inktank Storage、Inc.の米国内で登録された商標である。）

【0018】

データストレージシステム120に加えて、アプリケーション125はデータストレージシステム120のさらに上（top）で動作する。本発明の一部の実施形態では、アプリケーション125自身が内部理由によって二重書込みを遂行できる。例えば、アプリケーション125が実時間シミュレーションプログラムの場合である。このようなプログラムは動作が遂行される時間に大きく依存する。仮にデータが、バッファリングされている（buffered）が書き込まれていない（unwritten）状態で、実時間シミュレーションプログラムが中断される場合、シミュレーションの結果は捨てられる（wasted）。従って、データがデータストレージシステム120を通じてストレージ装置130に書き込まれない場合であっても、シミュレーションプログラムはデータがジャーナリングを通じて格納されることを保障させようとする。

【0019】

ストレージ装置130は、無効データのガーベッジコレクションを遂行できる多様な形態のストレージ装置の内の任意の一つである。例えば、ストレージ装置130はフラッシュベースのSSD（Solid State Drive）である。ストレージ装置130はストレージ装置130の動作の管理を担当するコントローラ135を含む。例えば、他に色々の機能がある中で、コントローラ135はデータ読出し及び書込みを管理し、ロジックブロック住所をストレージ装置305上の物理ブロック住所にマッピング（mapping）する。

【0020】

他の構成要素の中で、コントローラ135はコントローラ135を直接又は間接的にサーバ105に連結する物理インタフェース、ストレージ装置130の動作を制御するプロセッサ、フラッシュストレージ内に格納されたデータに対してエラーを感知し、訂正する機能を提供するECC（Error Correction Code）回路、ストレージ装置130内のDRAM（Dynamic Random Access Memory）を管理するDRAMコントローラ、そしてフラッシュストレージを管理する1つ又は1つ以上のフラッシュコントローラを含み得る。

【0021】

コントローラ135は、またマルチストリーミングコントローラを含み得る。マルチストリーミングコントローラは、何のデータが何のブロック（下で説明されるように、各ブロックは互に異なるストリームに関連付けされている）に書き込まれるかを管理する。一部の実施形態では、コントローラ135はこのような構成要素の機能をプログラムするのに適合する単一チップからなる。他の実施形態では、コントローラ135はこのような構成要素の全部又は一部を別々の構成要素（例えば、複数のチップ）として含む。

【0022】

図2は図1のサーバの詳細に示す図面である。図2を参照すれば、一般的に単数又は複数のサーバ105は単数又は複数のプロセッサ110を含む。プロセッサ110はメモリ

10

20

30

40

50

コントローラ 205 と、サーバ 105 の構成要素の動作を調整するのに使用されるクロック 210 を含む。また、プロセッサ 110 はメモリ 115 に連結される。例えば、メモリ 115 は RAM (Random Access Memory)、ROM (Read-Only Memory)、又は他の状態保持媒体 (state preserving media) を含む。

また、プロセッサ 110 はストレージ装置 130 及びネットワークコネクタ 215 に連結される。例えば、ネットワークコネクタ 215 はイーサネット (登録商標) (ethernet) コネクタである。また、他に色々の構成要素がある中で、プロセッサ 110 はバス 220 に連結される。バス 220 はユーザインタフェイス 225 と連結され、I/O (Input/Output) エンジン 230 を利用して管理される I/O インタフェイスポートに連結される。

10

【0023】

図 3 は従来技術に係り、ジャーナリング及びデータ書込みの両方を遂行するために図 1 のストレージ装置 130 と通信する図 1 のアプリケーション 125 を示す図面である。図 3 を参照すれば、アプリケーション 125 は、図 1 のデータストレージシステム 120 の利便を介すること無くストレージ装置 130 と通信する。

ただ 1 つのソフトウェア要素 (アプリケーション 125、図 3 参照) がストレージ装置 130 と通信する本発明の一実施形態において、該ソフトウェア要素はジャーナリングを遂行する任意のソフトウェア要素である場合、図 3 において、アプリケーション 125 は適用性 (applicability) を喪失すること無く、図 1 のデータストレージシステム 120 により置換できる。(以下の図 7 及び図 8 においては、図 1 のデータストレージシステム 120 及びアプリケーション 125 の両方がストレージ装置 130 とジャーナリングを遂行する本発明の別の実施形態が説明される。

20

【0024】

図 3 を参照すれば、アプリケーション 125 はジャーナル書込み要請 310 をストレージ装置 130 に伝送する。ジャーナル書込み要請 310 はジャーナル情報 305 及び第 1 ストリーム識別子 315 を含む。第 1 ストリーム識別子 315 はジャーナル情報 305 が割当てられる特定ストリームを指定する。以下で図 4 を参照して説明されるように、相異なるストリームはストレージ装置 130 上の相異なるブロック又はスーパーブロックと関連される。ストレージ装置 130 は単数又は複数の特性 (例えば、期待寿命又は他の何らかの分割基準等) に基づいてデータを分割 (partition) する。

30

ジャーナル情報 305 が直ちにストレージ装置 130 に書き込まれることを保障するために、ジャーナル書込み要請 310 は直接 I/O コマンドを使用して伝送される。又は別の選択肢として、同じくジャーナル情報 305 が直ちにストレージ装置 130 に書き込まれることを保障するために、ジャーナル書込み要請 310 は直ちに (一杯に書き込まれない場合にも) フラッシュ (flush) されるバッファに伝送される。

【0025】

また、アプリケーション 125 はデータ書込み要請 320 をストレージ装置 130 に伝送する。データ書込み要請 320 はデータ 325 及び第 2 ストリーム識別子 330 を含む。データ 325 がジャーナル情報 305 とは異なるストリームに書き込まれるように (結果的に相異なるブロック又はスーパーブロックに書き込まれる) するために、第 2 ストリーム識別子 330 は第 1 ストリーム識別子 315 と相異なるストリームを識別する。データ書込み要請 320 はバッファリングされた (buffered) 書込み要請として伝送される。その際、ジャーナル書込み要請 310 が直ちにストレージ装置 130 に書き込まれているので、データ 325 は、必ず直ぐではないが、何れはストレージ装置 130 に書き込まれる。

40

【0026】

最終的に、ストレージ装置 130 はデータ書き込み完了信号 335 をアプリケーション 125 に伝送する。信号 335 はデータ書込み要請 320 が完了され、データ 325 がストレージ装置 130 に書き込まれたことを示す。このような観点で、ジャーナル情報 30

50

5 はデータがストレージ装置 130 上のどこかに書き込まれた (w r i t t e n) ことをそれ以上保障する必要がなくなる。ガーベッジコレクションを遂行するストレージ装置上で、データは一般的にデータがガーベッジコレクションを通じて削除される前に無効化 (i n v a l i d a t e d) される。

このようにして、アプリケーション 125 は無効化要請 340 をストレージ装置 130 に伝送して、ジャーナル情報 305 がストレージ装置 130 から削除されるように要請する。ジャーナル情報 305 がデータ 325 とは異なるブロック又はスーパーブロックに書き込まれたので、無効化要請 340 はストレージ装置 130 内に有効か無効かにより断片化された (v a l i d i t y - f r a g m e n t e d) ブロック又はスーパーブロックを誘発しない。

10

【0027】

図4はジャーナル及びデータを格納するためにマルチストリーミングを使用する図1のストレージ装置130を示す図面である。図4を参照すれば、複数のブロックA、B、C、D(405、410、415、420)に区分されたストレージ装置130が示された。複数のブロック405、410、415、420の各々は順にページに分割される。例えば、ブロック405は複数のページ425、430、435、440を含み、ブロック410は複数のページ445、450、455、460を含む。図4で、複数のブロック405、410、415、420の各々が4つのページを含むことと図示されたが、これは単なる例示的なものであり、複数のブロック405、410、415、420の各々は任意の望む数のページを含み得る。

20

【0028】

S S Dにおいて生起するように、ページはストレージ装置130に書き込まれるか、或いは読み出される最小単位 of データを示す。反面に、一部の実施形態で、ブロックはガーベッジコレクションが遂行される最小単位 of データを示す。図4には図示しないが、一部の実施形態で、ストレージ装置130の複数のブロックはスーパーブロックであると称されるさらに大きいブロックに組織化されることができる。スーパーブロックはガーベッジコレクションが遂行される最小単位 of データである。ガーベッジコレクションがブロック又はスーパーブロック上で遂行されるかに拘わらず、ガーベッジコレクションの最小単位はページより大きい。

このような不一致はガーベッジコレクションがストレージ装置130の動作に否定的な影響を及ぼす理由を説明する。即ち、仮にガーベッジコレクションの対象であるブロックの単数又は複数のページ内に有効なデータがある場合、該ブロックがガーベッジコレクション動作の対象になる前にこのような有効データは他のブロックにコピーされなければならない。例えば、ページ445がフリー(f r e e)である場合、仮にページ425が有効データを含む場合、ブロック405がガーベッジコレクション動作の対象になる前に該データは、例えば別のブロック410内のページ445にコピーされなければならない。

30

【0029】

一部の実施形態で、ページは複数のブロックに組織される。そして、ガーベッジコレクションが複数のブロックに別々に遂行されるよりは、複数のブロックがスーパーブロックに組織化されてガーベッジコレクションがスーパーブロック上で遂行される。

40

但し、スーパーブロックの概念がストレージ装置130内のガーベッジコレクションの具現上に影響を及ぼすが、理論的な観点からは、スーパーブロックはガーベッジコレクション目的のためのブロックサイズの再調整に過ぎない。複数のブロックに関する全ての論議はスーパーブロックにも同様に適用されることは容易に理解されよう。

【0030】

図3について上述したように、個別ブロックはストリームに割当てられる。例えば、複数のブロック405、410は第1ストリーム315に割当てられる。また、複数のブロック415、420は第2ストリーム330に割当てられる。実施形態として、ホット(h o t)データ及びコールド(c o l d)データが相異なるストリームに区分できる状態で、ストリーム割当てはガーベッジコレクション動作に問題を惹起するジャーナル書込み

50

及びデータ書込みの混在を防止することができる。

【 0 0 3 1 】

図 5 及び図 6 は従来のアプローチ方法と本発明の実施形態によるアプローチ方法とを使用する図 1 のストレージ装置 1 3 0 の使用方法を比較して示す図面である。上述したように、従来のシステムでは、図 5 のブロック A (5 0 5) に示すように、ジャーナル書込み及びデータ書込みは図 1 のストレージ装置 1 3 0 内の同一のブロックに書き込まれる。図 5 を参照すれば、ブロック A (5 0 5) はジャーナル書込みを含む複数のページ 5 1 0、5 1 5、5 2 0 とデータ書込みを含む複数のページ 5 2 5、5 3 0、5 4 0 とを含む。

【 0 0 3 2 】

ジャーナル書込みは対応するデータ書込み動作が完了されれば、削除可能になる、即ち、ジャーナル書込みは短い寿命を有する傾向が存在するので、ジャーナル書込み及びデータ書込みの混在はブロック A (5 3 5) に示すように断片化されたブロックを残す。ここで、ブロック A (5 3 5) は、ブロック A (5 0 5) と物理的には同一のブロックであるが、ジャーナル書込みが無効化された後のブロック A を表す。仮にその後にブロック A (5 3 5) がガーベッジコレクションの対象になれば、複数のページ 5 2 5、5 3 0、5 4 0 は先ず他のブロックにコピーされなければならない。このようなコピー動作は時間を消耗し、ストレージ装置 1 3 0 上の他の読出し及び書込み動作を遅延させる。

【 0 0 3 3 】

しかし、本発明の実施形態のように、仮にジャーナル書込み及びデータ書込みが相異なるストリームに伝送されれば、ジャーナル書込みを消去することは断片化されたブロックを残さない。図 6 はこのような状況を示す。図 6 を参照すれば、ジャーナル書込みはブロック A (5 4 5) に伝送される。他方、データ書込みはブロック C (5 5 0) に伝送される。複数のジャーナル書込み 5 1 0、5 1 5、5 2 0 が無効化する場合、ブロック 5 4 5 は他のブロックにコピーされなければならない何らのデータも格納しない。複数のデータ書込み 5 2 5、5 3 0、5 4 0 は予め、別のブロック C (5 5 0) に格納されているからである。

上述した記述はこのような状況を単純化して記述したものである。但し、ジャーナル書込みは同一の時間に全て消去される必要がないので、一般的にジャーナル書込みはそれらが書き込まれた以後に、特にデータ書込みの寿命と比較して相対的に短い時間内に削除される。それ故、ブロック A (5 4 5) 内の全てのデータは最後のジャーナル書込みがブロック A (5 4 5) に書き込まれた後の短時間内に無効化される筈である。その際、何らのデータも他のブロックにコピーされる必要がなく、全ページを含む該ブロックがガーベッジコレクションの対象になる。

【 0 0 3 4 】

図 7 及び図 8 はジャーナリング及びデータ書込みの両方を遂行するために図 1 のストレージ装置 1 3 0 と通信する図 1 のアプリケーション 1 2 5 及び図 1 のデータストレージシステム 1 2 0 を示す図面である。アプリケーション 1 2 5 がデータストレージシステム 1 2 0 の利便を介する事無くストレージ装置 1 3 0 と通信する図 3 と対照的に、図 7 及び図 8 は一連のイベント進行手順においてデータストレージシステム 1 2 0 を含む。

【 0 0 3 5 】

図 7 を参照すれば、アプリケーション 1 2 5 はジャーナル書込み要請 3 1 0 を図 3 の場合と同じようにジャーナル情報 3 0 5 及び第 1 ストリーム識別子 3 1 5 と共に伝送する。ジャーナル書込み要請 3 1 0 が直接 I / O コマンドであるので、図 7 で、アプリケーション 1 2 5 がジャーナル書込み要請 3 1 0 をストレージ装置 1 3 0 に伝送し、その際、データストレージシステム 1 2 0 をバイパス (b y p a s s) することを示した。しかし、他の実施形態では、アプリケーション 1 2 5 はジャーナル書込み要請 3 1 0 をデータストレージシステム 1 2 0 に命令と共に伝送する。その際、該命令はデータストレージシステム 1 2 0 がジャーナル書込み要請 3 1 0 を完了するために直接 I / O コマンドを遂行させる。

しかし、図 3 と対照的に、アプリケーション 1 2 5 はデータ 3 2 5 及び第 2 ストリーム

10

20

30

40

50

識別子 330 と共にデータ書き込み要請 320 をストレージ装置 130 ではないデータストレージシステム 120 に伝送する。そうすると、データストレージシステム 120 はストレージ装置 130 に対するデータ 325 の書き込み動作の監督を担当する。

【0036】

データストレージシステム 120 自身が自分のデータ及び／又はメタデータに対するジャーナリングを遂行できるので、データストレージシステム 120 は第 2 ジャーナル書き込み要請 605 をストレージ装置 130 に伝送できる。第 2 ジャーナル書き込み要請 605 は第 2 ジャーナル情報 610 及び第 3 ストリーム識別子 615 を含む。これは単一データユニットの書き込みが複数のジャーナルを含み、従って複数のストリームを含み得ることを示す。

10

【0037】

図 8 を参照すれば、データストレージシステム 120 自身が自分の第 2 データ書き込み要請 620 をストレージ装置 130 に伝送する。第 2 データ書き込み要請 620 はデータ 325 及び第 2 ストリーム識別子 330 を含む。ここで、第 2 データ書き込み要請 620 内のデータ 325 及びストリーム識別子 330 がデータ書き込み要請 320 内のデータ 325 及びストリーム識別子 330 と同一であることが分かる。これで、同一のデータがデータストレージシステム 120 を通じて単に迂回して書き込まれることが理解されよう。最終的に、ストレージ装置 130 は第 2 データ書き込み要請 620 を完遂したことをデータストレージシステム 120 に通報する第 2 データ書き込み完了信号 625 を伝送し、以後データストレージシステム 120 は第 2 ジャーナル情報 610 を削除するために第 2 無効化要請 630 を伝送する。

20

【0038】

また、仮にストレージ装置 130 がアプリケーション 125 の存在を分かる場合、ストレージ装置 130 はデータ書き込み完了信号 335 を再びアプリケーション 125 に伝送し、これに対してアプリケーション 125 は自身の無効化要請 340 を伝送する。

しかし、他の実施形態では、データストレージシステム 120 は第 2 ジャーナル情報 610、ジャーナル情報 305 の全てを削除するために第 2 無効化要請 630 を使用する。

また、他の実施形態で、アプリケーション 125 にアプリケーション 125 自身が無効化要請 340 を伝送できることを知らせるために、データストレージシステム 120 はデータ書き込み完了信号 335 をアプリケーション 125 に伝送する。

30

【0039】

上述した実施形態で、アプリケーション 125 及び／又はデータストレージシステム 120 はジャーナル情報 305、第 2 ジャーナル情報 610 の全て又は一部の消去を担当することができる。従って、ジャーナル情報 305、第 2 ジャーナル情報 610 の全て又は一部の消去が何時可能になるかを知らるために、アプリケーション 125 及び／又はデータストレージシステム 120 はデータ書き込み完了信号 335、第 2 データ書き込み完了信号 625 を全て又は一部を受信する必要がある。

しかし、他の実施形態で、ストレージ装置 130 はデータ 325 のソースを知り、アプリケーション 125 及び／又はデータストレージシステム 120 に対する無効化要請 340、第 2 無効化要請 630 の全て又は一部を伝送する必要を事前に防止して、データ書き込み要請 320、第 2 データ書き込み要請 620 の全て又は一部が完遂されれば、自動的にジャーナル情報 305、第 2 ジャーナル情報 610 の全て又は一部を削除する。

40

【0040】

一部の実施形態では、多重個 (multiple_instances) のデータストレージシステム 120 が単一ストレージ装置 130 上に共存する。例えば、ストレージ装置 130 は複数のジャーナリングファイルシステムパーティション (partitions) を含む。又は、ストレージ装置 130 は多重個のオブジェクト (object) ストレージを維持する。このような実施形態で、多重個のデータストレージシステム 120 の各々は自分の第 2 ジャーナル書き込み要請 605 を同一のストレージ装置 130 に伝送する。各々の個別の第 2 ジャーナル書き込み要請 605 は自分の第 2 ジャーナル情報 610 及び

50

第3ストリーム識別子615を含む。多重個のデータストレージシステム120の各々は自分の第2データ書き込み要請620を同一のストレージ装置130に伝送する。各々の個別の第2データ書き込み要請620は自分のデータ325及び第2ストリーム識別子330を含む。

【0041】

多重個のデータストレージシステム120の各々は、自分の第2ジャーナル情報610及びデータ325を相異なるストリームに載せ、これに対してストレージ装置130は多様な第2ジャーナル情報610及びデータ325を相異なるブロック又はスーパーブロックに格納する。

斯くして、多重個のデータストレージシステム120の各々に対する第2ジャーナル情報610及びデータ325が、相異なるブロック又はスーパーブロックに存在するのみならず、相異なるデータストレージシステム120からの相異なる第2ジャーナル情報610は相異なるブロック又はスーパーブロック内に格納され、且つ、相異なるデータストレージシステム120からの相異なるデータ325は相異なるブロック又はスーパーブロック内に格納される。

【0042】

図9及び図10は本発明の実施形態によってジャーナリング及びデータ書き込みの両方を遂行するために、図1のストレージ装置130と通信する図1のアプリケーション125及び図1のデータストレージシステム120に対する手順を例示的に示すフローチャートである。

【0043】

図9を参照すれば、S705段階で、図1のアプリケーション125は、図1のストレージ装置130に書き込まれるべき図3のデータ325を識別する。S710段階で、図3のストリーム識別子315を図3のジャーナル情報305に対して使用するためのストリームとして指定するために、図1のアプリケーション125は図3のジャーナル書き込み要請310を直接I/Oコマンドとして図1のストレージ装置130に伝送する。

【0044】

この地点で、手順は相異なる経路に従って進行され得る。

一部の実施形態では、S715段階で、図1のアプリケーション125は図3のデータ書き込み要請320をバッファリングされたI/Oコマンドとして図3のストレージ装置130に伝送する。続いて、S720段階で、図1のアプリケーション125は図1のストレージ装置130からデータ書き込み完了信号335を受信する。データ書き込み完了信号335は図3のデータ書き込み要請320が完遂されたことを示す。最後に、S725段階で、図3のジャーナル情報305を削除するために、図1のアプリケーション125は図3の無効化要請340を図1のストレージ装置130に伝送する。

【0045】

又は、他の実施形態として、図1のアプリケーション125は、図3のデータ書き込み要請320を図1のストレージ装置130に直接伝送せず、その代わりに、S730段階で、図1のアプリケーション125は図3のデータ書き込み要請320を図1のデータストレージシステム120に伝送する。S735段階で、図1のデータストレージシステム120は図7の第2番目のジャーナル書き込み要請605を図1のストレージ装置130に伝送する。S740段階で、図1のデータストレージシステム120は図8の第2番目のデータ書き込み要請620を図1のストレージ装置130に伝送する。

【0046】

S745段階で、図1のデータストレージシステム120は図1のストレージ装置130から図8の第2データ書き込み完了信号625を受信する。第2データ書き込み完了信号625は図8の第2データ書き込み要請620が完遂されたことを示す。また、S750段階で、図7の第2ジャーナル情報610を削除するために、図1のデータストレージシステム120は図8の第2無効化要請630を図1のストレージ装置130に伝送する。続いて、手順は図9のS720段階に進行される。

【0047】

図9及び図10において、本発明の一部の実施形態が示された。しかし、当業者は上述した段階の順序を変更するか、一部段階を省略するか、或いは図面に示されない連結を含む、相異なる実施形態が具現できることは容易に理解できよう。明示的に説明されたか否かに関係無く、このような順序の変形は本発明の実施形態として看做される。

【0048】

以下に、本発明の概念の特定形態が具現できる適切な単数又は複数のマシンに対する簡略であり、一般的な説明が提供される。単数又は複数のマシンは、キーボード、マウス等のような従来の入力装置からの入力によって少なくとも一部が制御され、同様に他のマシン、VR (Virtual Reality、仮想) 環境との相互作用、バイオメトリック (biometric) フィードバック (feedback)、又は他の入力信号から受信された指示によっても制御され得る。

10

ここで、使用される「マシン」という用語は幅広く、単体マシン、仮想マシン、又は、通信可能に結合された、マシン、仮想マシン、又は協同動作する装置からなるシステムを含む。例示的なマシンはパーソナル (Personal) コンピュータ、ワークステーション、サーバ、ポータブル (Portable) コンピュータ、ポケット用 (Handheld) 装置、携帯電話、タブレット等を含む。また、例示的なシステムは自動車、列車、タクシーなどの、個人用又は公共用交通装置を含む。

【0049】

マシンは組み込み (embedded) コントローラ等を含む。例えば、組み込みコントローラはプログラムが可能であるか、或いは不可能なロジック装置又はアレイ (Array)、ASIC (Application Specific Integrated Circuits)、組み込みコンピュータ、スマートカード等を含む。マシンは単数又は複数の遠隔マシンとの単数又は複の連結を使用できる。例えば、このような連結はネットワークインタフェース、モデム、又は他の擬似伝達連結を通じてなされ得る。

20

【0050】

マシンはイントラネット、インターネット、LAN (Local Area Network)、WAN (Wide Area Network) 等の物理的及び/又は論理的ネットワーク方法によって相互連結され得る。当業者はネットワーク通信が多様な有線及び/又は無線近距離又は長距離キャリア、及びプロトコルを利用できることを容易に理解できよう。例えば、キャリア及びプロトコルは、RF (Radio Frequency)、衛星 (Satellite)、マイクロウェーブ (Microwave)、IEEE (Institute of Electrical and Electronics Engineers) 802.11、Bluetooth (登録商標)、可視光線、赤外線、ケーブル、レーザー等を含む。

30

【0051】

本発明の実施形態は機能、段階、データ構造、アプリケーションプログラムを含む関連データを参照して説明される。これらの機能、段階、データ構造、アプリケーションプログラムは、マシンによってアクセスされる場合にマシンがタスクを遂行するか、或いは抽象的なデータタイプ又はローレベル (Low-Level) のハードウェアコンテキスト (Context) を定義するようにする。例えば、上述した関連データはRAM、ROMのような揮発性及び/又は不揮発性メモリに格納される。また、関連データは他のストレージ装置及びそれらの関連ストレージ媒体に格納される。例えば、関連ストレージ媒体はハードドライブ、フロッピーディスク (Floppy Disks)、光学ストレージ (Optical Storage)、テープ (Tapes)、フラッシュメモリ (Flash Memory)、メモリスティック (Memory Sticks)、デジタルビデオディスク (Digital Video Disks)、生体ストレージ (Biological Storage) 等を含む。

40

【0052】

関連データは、パケット、シリアル (Serial) データ、パラレル (Parallel

50

e1) データ、電波信号等の形態で、物理的及び／又は論理的ネットワークを含む通信環境を通じて伝送される。また、関連データは圧縮されるか、或いは暗号化された形態で利用される。関連データは分散環境で利用され得、マシンアクセスに対して局所的に、及び／又は、遠隔に格納される。

【0053】

本発明の実施形態は、有形 (tangible) 且つ非一時的な (Non-transitory) マシンリーダブル (Readable) 媒体を含み得る。マシンリーダブル媒体は1つ又は1つ以上のプロセッサによって遂行される命令を含み、該命令はここで記述された本発明の要素を遂行する命令を含む。

【0054】

本発明の実施形態は、以下の例 (statement) に何ら制限されることなく拡張され得る。

【0055】

第1例で、本発明の実施形態によるシステムは、プロセッサ及びメモリを含むコンピュータと、ストレージ装置と、プロセッサ上で実行され、ジャーナル書込み要請及びデータ書込み要請の両方を前記ストレージ装置に伝送するように動作するアプリケーションと、前記ストレージ装置上に位置し、前記ストレージ装置が第1ストリームと関連された第1ブロックにジャーナル情報を書き込み、第2ストリームと関連された第2ブロックにデータを

書き込むように指示するコントローラと、を含み、前記ジャーナル書込み要請は前記ジャーナル情報を含み、前記第1ストリームと関連され、前記データ書込み要請は前記データを含み、前記第2ストリームと関連される。

【0056】

第2例で、本発明の実施形態による第1例のシステムにおいて、前記ストレージ装置がSSD (Solid State Drive) を含む。

【0057】

第3例で、本発明の実施形態による第1例のシステムにおいて、前記アプリケーションは直接 (direct) 入出力コマンドを利用して前記ストレージ装置に前記ジャーナル書込み要請を伝送するように動作する。

【0058】

第4例で、本発明の実施形態による第1例のシステムにおいて、前記アプリケーションはバッファリングされた (buffered) 書込みコマンドを利用して前記ストレージ装置に前記データ書込み要請を伝送するように動作する。

【0059】

第5例で、本発明の実施形態による第1例のシステムにおいて、前記コントローラはデータ書込み要請が完了された後に前記ジャーナル情報を削除するための無効化要請を受信するように動作する。

【0060】

第6例で、本発明の実施形態による第5例のシステムにおいて、前記アプリケーションは前記無効化要請を伝送するように動作する。

【0061】

第7例で、本発明の実施形態による第6例のシステムにおいて、前記データ書込み要請が完了したことの信号を前記アプリケーションが受信したことに反応して、前記アプリケーションが前記無効化要請を伝送するように動作する。

【0062】

第8例で、本発明の実施形態による第7例のシステムにおいて、前記データ書込み要請が完了したことの信号を前記コントローラから前記アプリケーションが受信したことに反応して、前記アプリケーションが前記無効化要請を伝送するように動作する。

10

20

30

40

50

【 0 0 6 3 】

第 9 例で、本発明の実施形態による第 1 例のシステムにおいて、前記アプリケーションは前記プロセッサ上で実行されるデータストレージシステムに前記データ書込み要請を送送するように動作し、

前記データストレージシステムは第 2 データ書込み要請を前記ストレージ装置に伝送するように動作し、

前記第 2 データ書込み要請は前記データを含み、前記第 2 ストリームと関連される。

【 0 0 6 4 】

第 1 0 例で、本発明の実施形態による第 9 例のシステムにおいて、前記データストレージシステムは第 2 ジャーナル書込み要請を前記ストレージ装置に伝送するように動作し、前記第 2 ジャーナル書込み要請は第 2 ジャーナル情報を含み、第 3 ストリームと連関される。

【 0 0 6 5 】

第 1 1 例で、本発明の実施形態による第 1 0 例のシステムにおいて、前記データストレージシステムは第 2 無効化要請を送送するように動作する。

【 0 0 6 6 】

第 1 2 例で、本発明の実施形態による第 1 0 例のシステムにおいて、前記第 2 データ書込み要請が完了したことの第 2 信号を前記データストレージシステムが受信したことに反応して、前記データストレージシステムが前記第 2 無効化要請を送送するように動作する。

【 0 0 6 7 】

第 1 3 例で、本発明の実施形態による第 1 2 例のシステムにおいて、前記第 2 データ書込み要請が完了したことの第 2 信号を前記コントローラから前記データストレージシステムが受信したことに反応して、前記データストレージシステムが前記第 2 無効化要請を送送するように動作する。

【 0 0 6 8 】

第 1 4 例で、本発明の他の実施形態による方法は、ジャーナル書込み及びデータ書込み両方を実行するアプリケーションから書き込まれるデータを識別する段階と、

前記アプリケーションから無効データにガーベッジコレクションを遂行するストレージ装置にジャーナル書込み要請を送送する段階と、

前記アプリケーションから前記ストレージ装置にデータ書込み要請を送送する段階と、を含み、

前記ジャーナル書込み要請は第 1 ストリームに割当され、前記データ書込み要請は第 2 ストリームに割当される。

【 0 0 6 9 】

第 1 5 例で、本発明の実施形態による第 1 4 例の方法において、前記アプリケーションから前記ストレージ装置に前記ジャーナル書込み要請を送送する段階は前記アプリケーションから S S D (S o l i d S t a t e D r i v e) に前記ジャーナル書込み要請を送送する段階を含み、

前記アプリケーションから前記ストレージ装置に前記データ書込み要請を送送する段階は前記アプリケーションから前記 S S D に前記データ書込み要請を送送する段階を含む。

【 0 0 7 0 】

第 1 6 例で、本発明の実施形態による第 1 4 例の方法において、前記アプリケーションから前記ストレージ装置に前記ジャーナル書込み要請を送送する段階は直接 (d i r e c t) 入出力コマンドを利用して前記アプリケーションから前記ストレージ装置に前記ジャーナル書込み要請を送送する段階を含み、

前記アプリケーションから前記ストレージ装置に前記データ書込み要請を送送する段階はバッファリングされた (b u f f e r e d) 書込みコマンドを利用して前記アプリケーションから前記ストレージ装置に前記データ書込み要請を送送する段階を含む。

【 0 0 7 1 】

第 1 7 例で、本発明の実施形態による第 1 4 例の方法において、前記データ書込み要請が前記ストレージ装置に書き込まれた以後に前記ジャーナル情報を削除するために無効化要請を前記ストレージ装置に伝送する段階をさらに含む。

【 0 0 7 2 】

第 1 8 例で、本発明の実施形態による第 1 7 例の方法において、前記無効化要請を前記ストレージ装置に伝送する段階は前記無効化要請を前記アプリケーションから前記ストレージ装置に伝送する段階を含む。

【 0 0 7 3 】

第 1 9 例で、本発明の実施形態による第 1 8 例の方法において、前記無効化要請を前記アプリケーションから前記ストレージ装置に伝送する段階は前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記アプリケーションで受信する段階を含む。

10

【 0 0 7 4 】

第 2 0 例で、本発明の実施形態による第 1 9 例の方法において、前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記アプリケーションで受信する段階は前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記アプリケーションで前記ストレージ装置から受信する段階を含む。

【 0 0 7 5 】

第 2 1 例で、本発明の実施形態による第 1 4 例の方法において、前記アプリケーションから前記ストレージ装置にデータ書込み要請を伝送する段階は、前記アプリケーションからデータストレージシステムに前記データ書込み要請を伝送する段階と、前記データストレージシステムから前記ストレージ装置に第 2 データ書込み要請を伝送する段階と、を含む。

20

【 0 0 7 6 】

第 2 2 例で、本発明の実施形態による第 2 1 例の方法において、前記データストレージシステムから前記ストレージ装置に前記第 2 データ書込み要請を伝送する段階は前記データストレージシステムから前記ストレージ装置に第 2 ジャーナル書込み要請を伝送する段階を含む。

30

【 0 0 7 7 】

第 2 3 例で、本発明の実施形態による第 2 2 例の方法において、前記第 2 データ書込み要請が前記ストレージ装置に書き込まれた以後に前記第 2 ジャーナル書込み要請によって書き込まれた前記データを削除するために無効化要請を前記ストレージ装置に伝送する段階をさらに含む。

【 0 0 7 8 】

第 2 4 例で、本発明の実施形態による第 2 3 例の方法において、前記無効化要請を前記ストレージ装置に伝送する段階は前記無効化要請を前記データストレージシステムから前記ストレージ装置に伝送する段階を含む。

【 0 0 7 9 】

第 2 5 例で、本発明の実施形態による第 2 4 例の方法において、前記無効化要請を前記データストレージシステムから前記ストレージ装置に伝送する段階は前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記データストレージシステムで受信する段階を含む。

40

【 0 0 8 0 】

第 2 6 例で、本発明の実施形態による第 2 5 例の方法において、前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記データストレージシステムで受信する段階は前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記データストレージシステムで前記ストレージ装置から受信する段階を含む。

【 0 0 8 1 】

50

第 27 例で、本発明の実施形態による非一時的な (Non - Transitory) 命令を格納するタイプのストレージ媒体を含む装置において、
前記非一時的な命令はマシンによって実行される時、
ジャーナル書込み及びデータ書込み両方を実行するアプリケーションから書き込まれるデータを識別する段階と、
前記アプリケーションから無効データにガーベッジコレクションを遂行するストレージ装置にジャーナル書込み要請を伝送する段階と、
前記アプリケーションから前記ストレージ装置にデータ書込み要請を伝送する段階と、を
遂行し、
前記ジャーナル書込み要請は第 1 ストリームに割当され、前記データ書込み要請は第 2 ス
トリームに割当される。

10

【0082】

第 28 例で、本発明の実施形態による第 27 例の装置において、
前記アプリケーションから前記ストレージ装置に前記ジャーナル書込み要請を伝送する段
階は前記アプリケーションから SSD (Solid State Drive) に前記ジ
ャーナル書込み要請を伝送する段階を含み、
前記アプリケーションから前記ストレージ装置に前記データ書込み要請を伝送する段階は
前記アプリケーションから前記 SSD に前記データ書込み要請を伝送する段階を含む。

【0083】

第 29 例で、本発明の実施形態による第 27 例の装置において、
前記アプリケーションから前記ストレージ装置に前記ジャーナル書込み要請を伝送する段
階は直接 (direct) 入出力コマンドを利用して前記アプリケーションから前記スト
レージ装置に前記ジャーナル書込み要請を伝送する段階を含み、
前記アプリケーションから前記ストレージ装置に前記データ書込み要請を伝送する段階は
バッファリングされた (buffered) 書込みコマンドを利用して前記アプリケー
ションから前記ストレージ装置に前記データ書込み要請を伝送する段階を含む。

20

【0084】

第 30 例で、本発明の実施形態による第 27 例の装置において、
前記非一時的な命令はマシンによって実行される時、
前記データ書込み要請が前記ストレージ装置に書き込まれた以後に前記ジャーナル情報を
削除するために無効化要請を前記ストレージ装置に伝送する段階をさらに遂行する。

30

【0085】

第 31 例で、本発明の実施形態による第 30 例の装置において、前記無効化要請を前記
ストレージ装置に伝送する段階は前記無効化要請を前記アプリケーションから前記スト
レージ装置に伝送する段階を含む。

【0086】

第 32 例で、本発明の実施形態による第 31 例の装置において、前記無効化要請を前記
アプリケーションから前記ストレージ装置に伝送する段階は前記ストレージ装置上の前記
データ書込み要請が完了したことの信号を前記アプリケーションで受信する段階を含む
。

40

【0087】

第 33 例で、本発明の実施形態による第 32 例の装置において、前記ストレージ装置上
の前記データ書込み要請が完了したことの信号を前記アプリケーションで受信する段階
は前記ストレージ装置上の前記データ書込み要請が完了したことの信号を前記アプリ
ケーションで前記ストレージ装置から受信する段階を含む。

【0088】

第 34 例で、本発明の実施形態による第 27 例の装置において、前記アプリケーション
から前記ストレージ装置にデータ書込み要請を伝送する段階は、
前記アプリケーションからデータストレージシステムに前記データ書込み要請を伝送する
段階と、

50

前記データストレージシステムから前記ストレージ装置に第２データ書き込み要請を送送する段階と、を含む。

【００８９】

第３５例で、本発明の実施形態による第３４例の装置において、前記データストレージシステムから前記ストレージ装置に前記第２データ書き込み要請を送送する段階は前記データストレージシステムから前記ストレージ装置に第２ジャーナル書き込み要請を送送する段階を含む。

【００９０】

第３６例で、本発明の実施形態による第３５例の装置において、
前記非一時的な命令はマシンによって実行される時、
前記第２データ書き込み要請が前記ストレージ装置に書き込まれた以後に前記第２ジャーナル書き込み要請によって書き込まれた前記データを削除するために無効化要請を前記ストレージ装置に伝送する段階をさらに実行する。

10

【００９１】

第３７例で、本発明の実施形態による第３６例の装置において、前記無効化要請を前記ストレージ装置に伝送する段階は前記無効化要請を前記データストレージシステムから前記ストレージ装置に伝送する段階を含む。

【００９２】

第３８例で、本発明の実施形態による第３７例の装置において、前記無効化要請を前記データストレージシステムから前記ストレージ装置に伝送する段階は前記ストレージ装置上の前記データ書き込み要請が完了したことの信号を前記データストレージシステムで受信する段階を含む。

20

【００９３】

上記に説明した内容は本発明を実施するための具体的な例である。本発明には上記に説明した実施形態のみだけでなく、単純に設計変更するか、或いは容易に変更できる実施形態も含まれる。また、本発明には上述した実施形態を利用して将来に容易に変形して実施できる技術も含まれる。

【符号の説明】

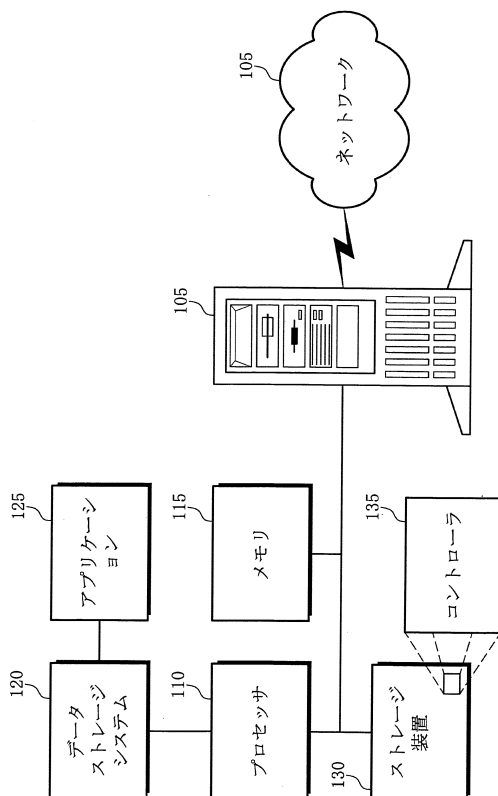
【００９４】

１０５	サーバ	30
１１０	プロセッサ	
１１５	メモリ	
１２０	データストレージシステム	
１２５	アプリケーション	
１３０	ストレージ装置	
１３５	コントローラ	
２０５	メモリコントローラ	
２１０	クロック	
２１５	ネットワークコネクタ	
２２０	バス	40
２２５	ユーザインタフェース	
２３０	Ｉ／Ｏエンジン	
３０５	ジャーナル情報	
３１０	ジャーナル書き込み要請	
３１５	第１ストリーム識別子、第１ストリーム	
３２０	データ書き込み要請	
３２５	データ	
３３０	第２ストリーム識別子、第２ストリーム	
３３５	データ書き込み完了信号	
３４０	無効化要請	50

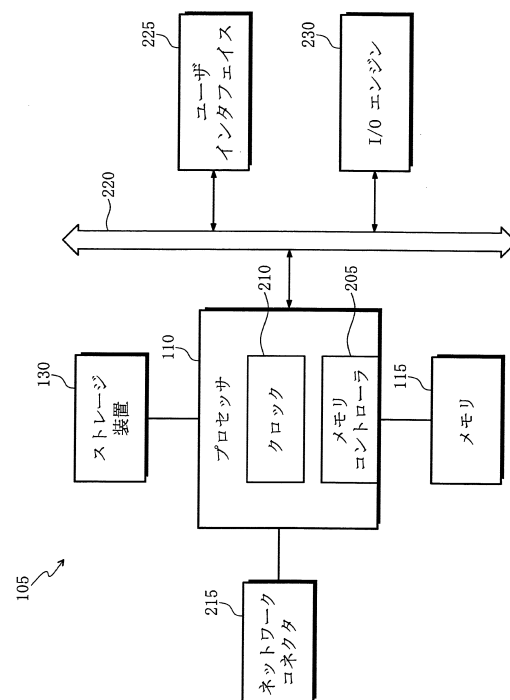
4 0 5、4 1 0、4 1 5、4 2 0	ブロック A、B、C、D
4 2 5 ~ 4 4 0、4 4 5 ~ 4 6 0	ページ
5 0 5、5 3 5	ブロック A、A
5 1 0、5 1 5、5 2 0	ジャーナル書き込みを含むページ、ジャーナル書き込み
5 2 5、5 3 0、5 4 0	データ書き込みを含むページ、データ書き込み
5 4 5、5 5 0	ブロック A、C
6 0 5	第 2 ジャーナル書き込み要請
6 1 0	第 2 ジャーナル情報
6 1 5	第 3 ストリーム識別子、第 3 ストリーム
6 2 0	第 2 データ書き込み要請
6 2 5	第 2 データ書き込み完了信号
6 3 0	第 2 無効化要請

10

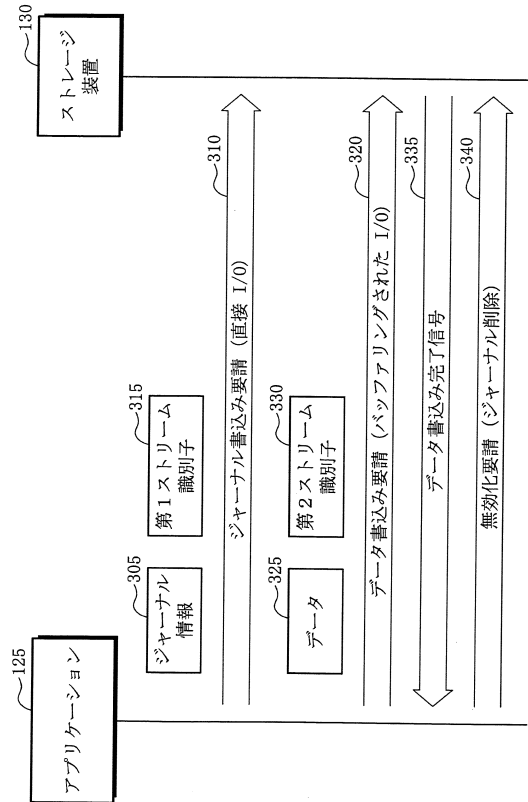
【図 1】



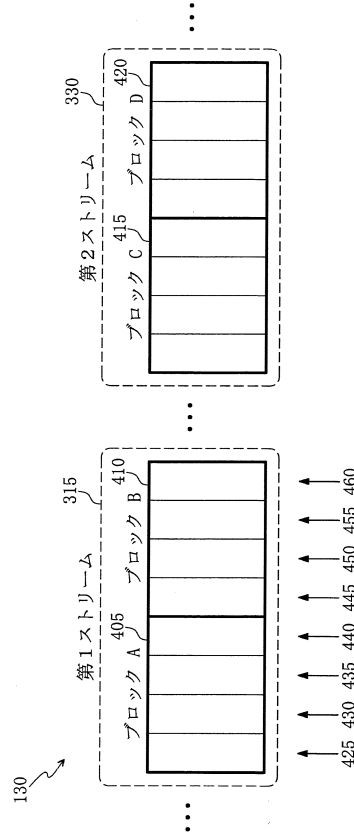
【図 2】



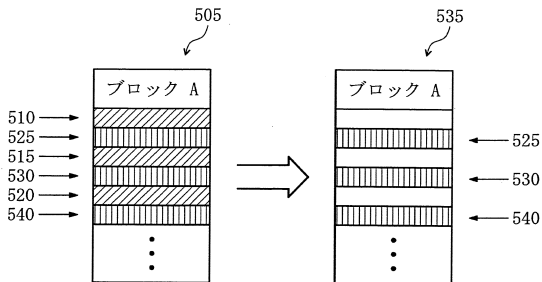
【図 3】



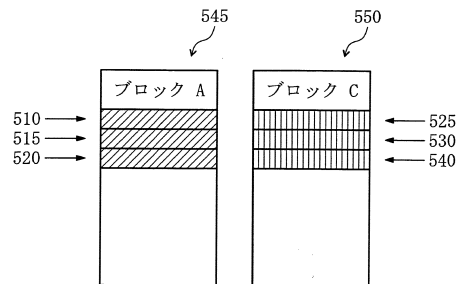
【図 4】



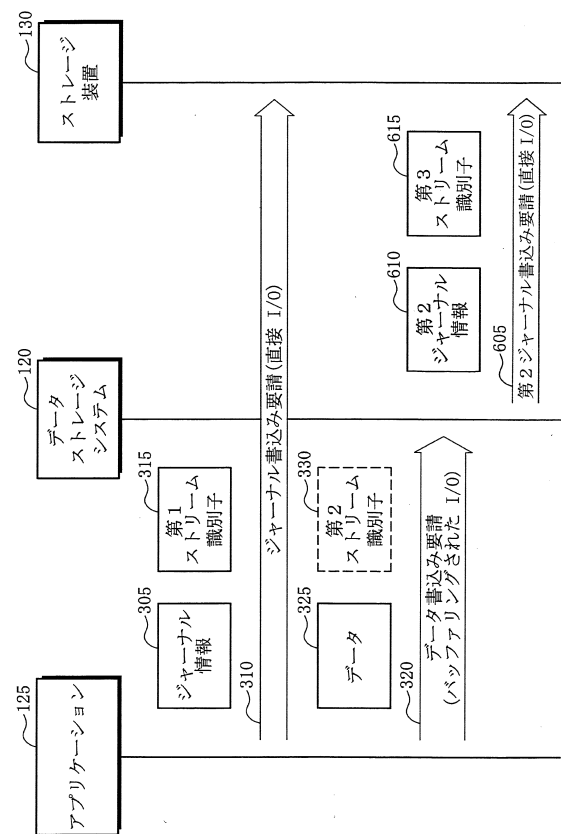
【図 5】



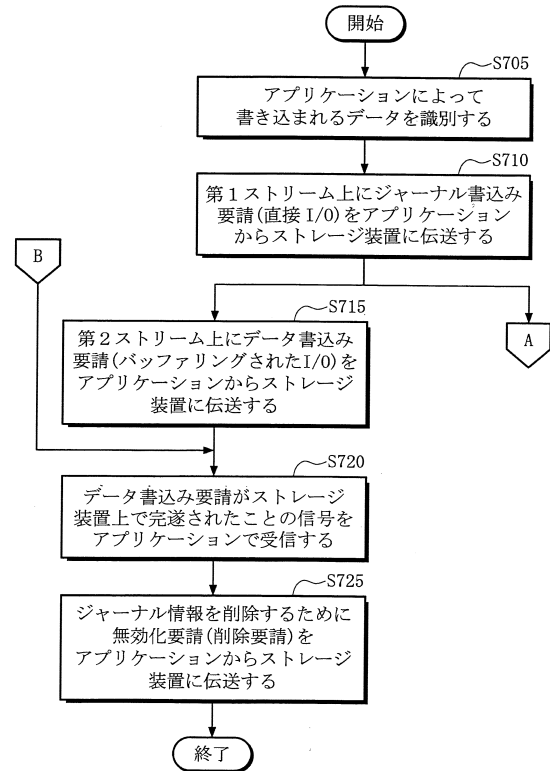
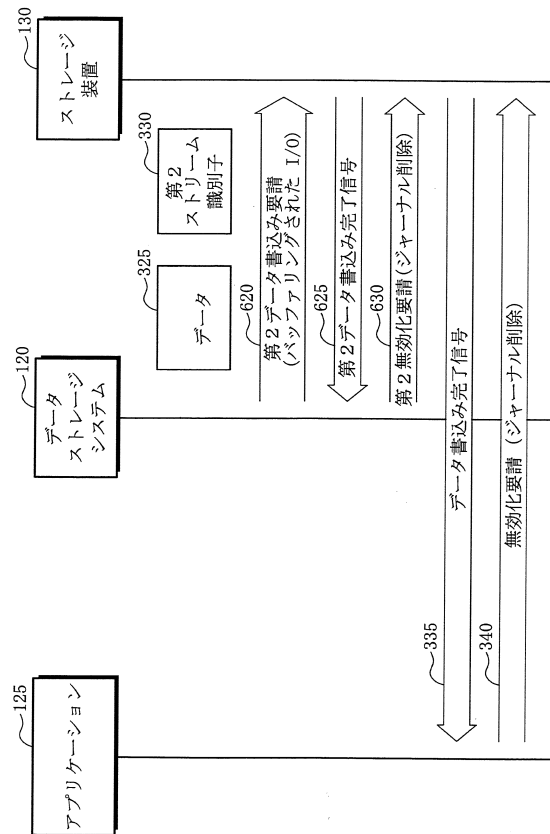
【図 6】



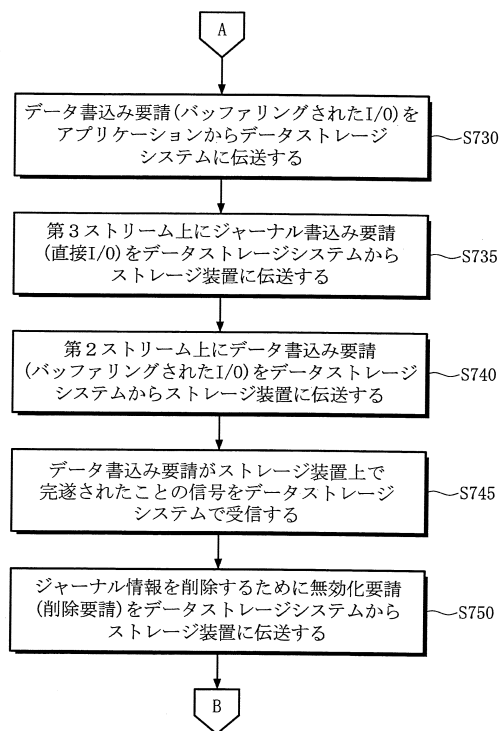
【図 7】



【 図 9 】



【 図 1 0 】



フロントページの続き

(72)発明者 チェ, チャンホウ

アメリカ合衆国, カリフォルニア州 95120, サンジョゼ, マウント・ハリ・ドライブ, 66
22

審査官 松尾 真人

(56)参考文献 米国特許出願公開第2014/0281172(US, A1)

米国特許第06128630(US, A)

米国特許出願公開第2014/0149473(US, A1)

米国特許出願公開第2014/0208001(US, A1)

米国特許出願公開第2014/0337562(US, A1)

(58)調査した分野(Int.Cl., DB名)

G06F 16/00 - 16/958

G06F 12/16