

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5695040号
(P5695040)

(45) 発行日 平成27年4月1日(2015.4.1)

(24) 登録日 平成27年2月13日(2015.2.13)

(51) Int. Cl.	F I
G 0 6 F 1 2 / 0 0 (2 0 0 6 . 0 1)	G 0 6 F 1 2 / 0 0 5 2 0 J

請求項の数 36 (全 52 頁)

(21) 出願番号	特願2012-517787 (P2012-517787)	(73) 特許権者	511314120
(86) (22) 出願日	平成22年6月25日 (2010.6.25)		シンプリヴィティ・コーポレーション
(65) 公表番号	特表2012-531675 (P2012-531675A)		アメリカ合衆国マサチューセッツ州015
(43) 公表日	平成24年12月10日 (2012.12.10)		81-1756, ウェストボロー, テクノ
(86) 国際出願番号	PCT/US2010/040058		ロジー・ドライブ 8
(87) 国際公開番号	W02010/151813	(74) 代理人	100140109
(87) 国際公開日	平成22年12月29日 (2010.12.29)		弁理士 小野 新次郎
審査請求日	平成25年6月25日 (2013.6.25)	(74) 代理人	100075270
(31) 優先権主張番号	61/269, 633		弁理士 小林 泰
(32) 優先日	平成21年6月26日 (2009.6.26)	(74) 代理人	100080137
(33) 優先権主張国	米国 (US)		弁理士 千葉 昭男
		(74) 代理人	100096013
			弁理士 富田 博行
		(74) 代理人	100120112
			弁理士 中西 基晴

最終頁に続く

(54) 【発明の名称】 ファイルシステム

(57) 【特許請求の範囲】

【請求項1】

1つ又は複数のコンピュータストレージデバイスでのファイルの命名及び記憶のためのコンピュータファイルシステムであって、

前記ファイルシステムは、オブジェクトストアにアクセスする名前空間ファイルシステムを含み、前記名前空間ファイルシステムは、メモリと、該メモリと通信するハードウェアプロセッサとを含み、前記プロセッサは、オブジェクトフィンガープリントを用いて前記オブジェクトストアにアクセスするためにプログラム命令を実行し、前記オブジェクトストアはファイル、データ及びメタデータをオブジェクトとして保持し、各オブジェクトは、該オブジェクトのコンテンツから導出され、前記オブジェクトストアにアクセスするために用いられるグローバルに一意的なオブジェクトフィンガープリントを有し、

各ファイルのオブジェクトは、該ファイルの前記データのオブジェクト又はメタデータのオブジェクトに対するオブジェクトフィンガープリントのマッピングを含み、前記ファイルのオブジェクトは、該ファイル内の前記オブジェクトの前記フィンガープリントから導出された前記ファイルのオブジェクト自体のオブジェクトフィンガープリントを有し、前記オブジェクトストアは、更に、

ファイルシステムのiノード番号とオブジェクトフィンガープリントとのマッピングを含むiノードマップのオブジェクトであって、該オブジェクトフィンガープリントは、前記iノード番号を一定とすることを可能にしつつ、前記ファイルのコンテンツに変更が生じると変更される、iノードマップのオブジェクトと、

10

20

各々が i ノード番号とファイル名とのマッピングを含む複数のディレクトリのオブジェクトと、

を含み、前記 i ノードマップのオブジェクトとディレクトリのオブジェクトとの各々は、それぞれのオブジェクトのコンテンツから導出した自分自身のオブジェクトフィンガープリントを有することを特徴とするファイルシステム。

【請求項 2】

請求項 1 に記載のファイルシステムであって、オブジェクト参照は、前記オブジェクトフィンガープリントによってマッピングされることを特徴とするファイルシステム。

【請求項 3】

請求項 1 又は 2 に記載のファイルシステムであって、前記ファイルのオブジェクトのマッピングは、線形リスト、ツリー構造又は間接指定テーブルを含むことを特徴とするファイルシステム。

10

【請求項 4】

請求項 1 から 3 の何れか一項に記載のファイルシステムであって、前記ファイルのオブジェクトはルートオブジェクトを含み、前記ルートオブジェクトは、前記ファイルシステム内の前記オブジェクトの全てから導出された該ルートオブジェクト自体のオブジェクトフィンガープリントを有し、それにより、前記ファイルシステム内のあらゆるオブジェクトは前記ルートオブジェクトを通じてアクセス可能であることを特徴とするファイルシステム。

【請求項 5】

20

請求項 1 から 4 の何れか一項に記載のファイルシステムであって、前記名前空間ファイルシステムは、仮想ファイルシステム層とブロックストレージ抽象化層の間の、ストレージスタックにおける層として提供されることを特徴とするファイルシステム。

【請求項 6】

請求項 2 から 5 の何れか一項に記載のファイルシステムにおいて、前記オブジェクトストアは、オブジェクトフィンガープリント、オブジェクト位置及びオブジェクト参照カウンタのインデックスを含むことを特徴とするファイルシステム。

【請求項 7】

請求項 6 に記載のファイルシステムであって、前記オブジェクトストアのインデックスは、不揮発性メモリに記憶されたことを特徴とするファイルシステム。

30

【請求項 8】

請求項 1 から 7 の何れか一項に記載のファイルシステムであって、前記フィンガープリントは、前記オブジェクトのコンテンツの暗号ハッシュダイジェストであることを特徴とするファイルシステム。

【請求項 9】

請求項 1 から 8 の何れか一項に記載のファイルシステムであって、前記オブジェクトのサイズは可変であることを特徴とするファイルシステム。

【請求項 10】

請求項 1 から 9 の何れか一項に記載のファイルシステムであって、前記ファイルシステムは P O S I X 規格準拠ファイルシステムであることを特徴とするファイルシステム。

40

【請求項 11】

請求項 1 から 10 の何れか一項に記載のファイルシステムであって、前記オブジェクトストアは、更に、

各オブジェクトに対する参照カウントを含む位置インデックスを含み、概参照カウントは前記オブジェクトが参照される回数を示すことを特徴とするファイルシステム。

【請求項 12】

請求項 11 に記載のファイルシステムであって、

読出し、書込み、削除及び参照カウント更新を含むオブジェクトアクティビティのトランザクションログ

を含むファイルシステム。

50

【請求項 13】

請求項 4 又は請求項 4 に従属する場合の請求項 5 から 12 の何れか一項に記載のファイルシステムであって、ファイルシステムの何れかのオブジェクトのコンテンツの変更は前記ルートオブジェクトの変更を生じ、前記ルートオブジェクトにおける変更を追跡することは、ファイルシステムアクティビティの履歴を提供することを特徴とするファイルシステム。

【請求項 14】

請求項 1 から 13 の何れか一項に記載のファイルシステムであって、

前記オブジェクトの生成及び記憶の 1 つ又は複数を実施するためのハードウェアアクセラレータ

10

を含むファイルシステム。

【請求項 15】

請求項 1 から 14 の何れか一項に記載のファイルシステムであって、

スタック

を含み、前記オブジェクトストアは前記スタックの下部を構成し、前記ファイルシステムは前記スタックの上部を構成することを特徴とするファイルシステム。

【請求項 16】

請求項 7 又は請求項 7 に従属する場合の請求項 8 から 15 の何れか一項に記載のファイルシステムであって、前記不揮発性メモリはソリッドステートディスクを含むことを特徴とするファイルシステム。

20

【請求項 17】

請求項 16 に記載のファイルシステムであって、前記メモリはフラッシュメモリを含むことを特徴とするファイルシステム。

【請求項 18】

請求項 1 から 17 の何れか一項に記載のファイルシステムであって、前記名前空間ファイルシステム及び前記オブジェクトストアは、デジタル電子回路、コンピュータハードウェア、ファームウェア、非一時的マシン読み取り可能ストレージデバイス内のコンピュータプログラム又はそれらの組み合わせのうちの 1 つ又は複数において実装されたことを特徴とするファイルシステム。

【請求項 19】

30

名前空間ファイルシステムによりオブジェクトストアにアクセスするステップであって、該オブジェクトストアは、ファイル、データ又はメタデータをオブジェクトとして保持し、各オブジェクトは、グローバルに一意であり、該オブジェクトのコンテンツから導出され、前記オブジェクトストアにアクセスするために用いられるオブジェクトフィンガープリントを有し、各ファイルのオブジェクトは、オブジェクトフィンガープリントの、該ファイルの前記データのオブジェクト又はメタデータのオブジェクトに対するマッピングを含み、前記ファイルのオブジェクトは、該ファイル内の前記オブジェクトの前記フィンガープリントから導出された該ファイルのオブジェクト自体のオブジェクトフィンガープリントを有する、ステップと、

ファイルシステムの *i* ノード番号とオブジェクトフィンガープリントとのマッピングを含む *i* ノードマップのオブジェクトを前記オブジェクトストアに維持するステップであって、該オブジェクトフィンガープリントは、前記 *i* ノード番号を一定とすることを可能にしつつ、前記ファイルのコンテンツに変更が生じると変更される、ステップと、

40

前記オブジェクトストアにおいてディレクトリのオブジェクトを維持するステップであって、各ディレクトリのオブジェクトは *i* ノード番号とファイル名とのマッピングを含む、ステップと

を含み、前記 *i* ノードマップのオブジェクト及びディレクトリのオブジェクトの各々は、それぞれの前記オブジェクトのコンテンツから導出した自分自身のオブジェクトフィンガープリントを有することを特徴とする方法。

【請求項 20】

50

請求項 19 に記載の方法であって、オブジェクトフィンガープリントと前記オブジェクトの物理的位置とをマッピングするための位置インデックスを維持するステップを含む方法。

【請求項 21】

請求項 20 に記載の方法であって、前記位置インデックスは、前記オブジェクトに対する参照カウントを含むことを特徴とする方法。

【請求項 22】

請求項 21 に記載の方法であって、前記フィンガープリント、位置インデックス、i ノードマップのオブジェクト、ファイルのオブジェクト及びディレクトリのオブジェクトは、ファイルシステムを構成することを特徴とする方法。

10

【請求項 23】

請求項 19 から 22 の何れか一項に記載の方法であって、前記ファイルシステムのスナップショットを前記 i ノードマップのオブジェクトの前記フィンガープリントから生成するステップを含む方法。

【請求項 24】

請求項 19 から 23 の何れか一項に記載の方法であって、前記 i ノードマップのフィンガープリントを別のオブジェクトストアに発行するステップを含む方法。

20

【請求項 25】

請求項 19 から 24 の何れか一項に記載の方法であって、前記 i ノードマップのフィンガープリントを前記ファイルシステムのスナップショットとして用いて障害回復を行うステップを含む方法。

【請求項 26】

請求項 19 から 25 の何れか一項に記載の方法であって、前記 i ノードマップのオブジェクトは以前の i ノードマップのフィンガープリントを含むことを特徴とする方法。

【請求項 27】

請求項 26 に記載の方法であって、前記以前の i ノードマップのフィンガープリントは、前記ファイルシステムのスナップショットの履歴を構成することを特徴とする方法。

30

【請求項 28】

請求項 19 から 27 の何れか一項に記載の方法であって、前記オブジェクトストアにおいてオブジェクト名と物理的オブジェクト位置との位置インデックスを維持するステップを含む方法。

【請求項 29】

請求項 21 又は請求項 21 に従属する場合の請求項 22 から 28 の何れか一項に記載の方法であって、

前記ファイルシステムに対する変更の際に、前記 i ノードマップのオブジェクトの下のあらゆるオブジェクトの前記参照カウントを調節するステップを含む方法。

40

【請求項 30】

請求項 29 に記載の方法であって、調節する前記ステップは、あらゆる入力/出力トランザクションに関して行われて、継続的データ保護が提供されることを特徴とする方法。

【請求項 31】

請求項 29 に記載の方法であって、調節する前記ステップは、周期的に、オンデマンドで、又は、特定のイベントにおいて行われて、スナップショットが生成されることを特徴とする方法。

【請求項 32】

50

請求項 19 から 31 の何れか一項に記載の方法であって、前記ファイルのオブジェクトのマッピングは、前記ファイルの前記コンテンツへのオフセットによってインデックス付けされ、線形リスト、ツリー構造又は間接指定テーブルを含むことを特徴とする方法。

【請求項 33】

請求項 21 から 32 の何れか一項に記載の方法であって、読出し、書込み、削除及び参照カウント更新を含む全てのオブジェクトアクティビティのトランザクションログを生成するステップを含む方法。

【請求項 34】

請求項 19 から 33 の何れか一項に記載の方法であって、
前記ファイルのオブジェクトを追加、修正又は削除し、ファイルの新しいオブジェクトフィンガープリントを生成するステップを含む方法。

10

【請求項 35】

請求項 34 に記載の方法であって、
ハードウェアアクセラレーションを利用してオブジェクトの生成、圧縮及び暗号化のうちの 1 つ又は複数を行うステップを含む方法。

【請求項 36】

プロセッサに、請求項 19 から 35 の何れか一項に記載の方法の前記ステップを実行させるプログラムコードを含む、非一時的マシン読み取り可能ストレージデバイスに記憶されたコンピュータプログラム。

20

【発明の詳細な説明】

【技術分野】

【0001】

[001]本発明は、コンピュータファイルシステムデータ構造、並びに、ファイルの命名及び格納のための方法及び装置に関する。

【背景技術】

【0002】

[002]十分な機能を有するストレージソリューションは、生のディスク、ファイルシステム、スナップショット、ファイルバージョンニング、圧縮、暗号化、内蔵容量最適化（例えば、データ重複排除）、監査及びタンパーレジスタンス等のセキュリティ機能、障害回復のためのオフサイト位置への効率的なレプリケーション等を含む場合がある。これらの機能の多くは、別々の機器において達成され、これらの機器は次いで、経験豊富な技術者によって接続されなければならない。

30

【0003】

[003]数テラバイト（TB）のデータ向けの、今日の技術によるそのようなストレージソリューションの構築は、結果として、10万ドルものコストを容易に上回る可能性のあるマルチボックスソリューションとなり、そのような十分な機能を有するストレージソリューションが、多数の企業及び顧客にとっては使用できないものとなることが多い。

40

【0004】

[004]このマルチボックス、アドホックなソリューションは、ストレージの基本的な側面ではなく、むしろ、ファイルシステムのアーキテクチャ及び実装が、他の技術開発についてきていないということである。例えば、大抵のファイルシステムアーキテクチャは、より高速なコンピュータ処理装置（CPU）、フラッシュメモリ、及び、ネットワーク帯域幅、ディスク密度及びディスクアクセスレートの間の異なるバランスを十分に活用するほどに、発展してきていない。

【発明の概要】

【発明が解決しようとする課題】

【0005】

50

[005]データアクセシビリティを、アクセス帯域幅とアドレス指定可能ストレージの比として定義する場合、データのアクセシビリティは低下しつつある。ストレージ密度は、ディスクへのアクセスより速く増大しつつあるので、所与のデータセットサイズでは、データにアクセスするために必要とされる時間は増大しつつある（且つ、よって、アクセシビリティの低減を引き起こす）。ストレージアーキテクチャへの影響は、以下の通りであり、即ち、データを格納した後、絶対に必要でない限り、データを移動するべきではない。この単純な観察事実は、データが読み込まれ、再び書き出されることが絶えず行われている現在のストレージアーキテクチャでは、何度も妨げられる。この結果、著しい余分な支出（例えば、10チャンネル、CPU、電力、時間、管理）となる。

【課題を解決するための手段】

【0006】

[006]本発明の一実施形態によれば、ファイルシステム提供され、このファイルシステムは、

デジタルに署名されたファイルシステムであって、該システムにおいてデータ、メタデータ及びファイルはオブジェクトであり、各オブジェクトは、グローバルに一意でありコンテンツ導出のフィンガープリントを有し、オブジェクト参照はフィンガープリントによってマップされる、デジタルに署名されたファイルシステムを含み、

ファイルシステムは、ファイルシステム内の全てのオブジェクトのフィンガープリントのマッピングを含むルートオブジェクトを有し、

ファイルシステムに対する変更は、ルートオブジェクト内の変更を生じ、ルートオブジェクト内の変更を追跡することは、ファイルシステムアクティビティの履歴を提供する。

【0007】

一実施形態では、

ファイルシステムは、ファイルオブジェクトのフィンガープリントへのiノード番号のマッピングを含むiノードマップオブジェクトを含み、iノードマップオブジェクトのフィンガープリントは、ファイルシステムのスナップショットを含む。

【0008】

本発明の別の実施形態によれば、格納されたオブジェクトをインデックス付けする方法のための実行可能プログラム命令を含むコンピュータ可読媒体が提供され、上記方法は、

データ、メタデータ及びファイルをオブジェクトとして提供するステップと、

オブジェクトごとにフィンガープリントを提供するステップであって、該フィンガープリントはグローバルに一意でありオブジェクトのコンテンツから導出される、ステップとを含み、

ファイルシステムルートオブジェクトが提供され、該オブジェクトは、ファイルシステム内の全てのオブジェクトのフィンガープリントのマッピングを含み、ファイルシステムに対する変更がルートオブジェクト内の変更を生じ、ルートオブジェクト内の変更を追跡することはファイルシステムアクティビティの履歴を提供する。

【0009】

一実施形態では、この方法は、

ファイルシステムiノードマップオブジェクトを提供するステップであって、該オブジェクトはファイルオブジェクトのフィンガープリントへのiノード番号のマッピングを含む、ステップを含み、

iノードマップオブジェクトのフィンガープリントは、ファイルシステムのスナップショットを含む。

【0010】

一実施形態では、この方法は、

iノードマップフィンガープリントを、別個のオブジェクトストア上の別のコンピュータシステムに発行するステップを含む。

【0011】

一実施形態では、この方法は、

10

20

30

40

50

i ノードマップフィンガープリントを、障害回復のために、ファイルシステムのスナップショットとして使用するステップを含む。

【0012】

—実施形態では、

i ノードマップオブジェクトは、以前の i ノードマップのフィンガープリントを含む。

—実施形態では、

以前の i ノードマップフィンガープリントは、ファイルシステムのスナップショットの履歴を含む。

【0013】

—実施形態では、

オブジェクトは、参照カウントを有し、

ファイルシステムに対する変更があると、i ノードマップオブジェクトの下のあらゆるオブジェクトのオブジェクト参照カウントを調節するステップである。

【0014】

—実施形態では、

調節するステップは、あらゆる I/O トランザクションにおいて行われて、継続的データ保護が提供される。

【0015】

—実施形態では、

調節するステップは、周期的に、オンデマンドで又は特定のイベントにおいて行われて、スナップショットが生成される。

【0016】

—実施形態では、

オブジェクトは、参照カウントを有し、

参照カウントに対する調節は、新しいデータコンテンツのみが格納されるように、データ重複排除のために利用される。

【0017】

本発明の別の実施形態によれば、1つ又は複数のコンピュータストレージデバイス上でのファイルの命名及び格納のためのコンピュータファイルシステムが提供され、このシステムは、

名前空間ファイルシステムを含み、ファイル、データ及びメタデータはオブジェクトであり、各オブジェクトは、オブジェクトのコンテンツから導出されたグローバルに一意的なフィンガープリントを有し、各ファイルオブジェクトは、ファイルのデータオブジェクト及びメタデータオブジェクトのうち的一方又は双方のためのオブジェクトフィンガープリントのマッピングを含み、ファイルオブジェクトは、ファイル内のオブジェクトのフィンガープリントから導出された該ファイルオブジェクト自体のオブジェクトフィンガープリントを有し、このシステムは、ファイルオブジェクトのフィンガープリントへの i ノード番号のマッピングを含む。

【0018】

—実施形態では、

オブジェクト参照は、オブジェクトフィンガープリントによって定まる。

—実施形態では、

ファイルオブジェクトのマッピングは、線形リスト、ツリー構造又は間接指定テーブルを含む。

【0019】

—実施形態では、

ファイルオブジェクトはルートオブジェクトを含み、このルートオブジェクトは、ファイルシステム内のオブジェクトの全てから導出された該ルートオブジェクト自体のオブジェクトフィンガープリントを有し、ファイルシステム内のあらゆるオブジェクトはルートオブジェクトを通じてアクセス可能である。

10

20

30

40

50

【 0 0 2 0 】

－実施形態では、

名前空間ファイルシステムは、仮想ファイルシステム層とブロックストレージ抽象化層の間の、ストレージスタックにおける層として提供される。

【 0 0 2 1 】

－実施形態では、このシステムは、

オブジェクトフィンガープリント、オブジェクト位置及びオブジェクト参照カウン트의インデックスを含むオブジェクトストアを更に含む。

【 0 0 2 2 】

－実施形態では、

オブジェクトストアインデックスは、不揮発性メモリに格納される。

－実施形態では、

フィンガープリントは、オブジェクトコンテンツの暗号ハッシュダイジェストである。

【 0 0 2 3 】

－実施形態では、

オブジェクトサイズは可変である。

－実施形態では、

ファイルシステムは P O S I X 準拠ファイルシステムである。

【 0 0 2 4 】

本発明の別の実施形態によればある方法が提供され、この方法は、

ファイルシステム内のデータオブジェクトのためのオブジェクトフィンガープリントを生成するステップであって、データオブジェクトはデータ又はメタデータを含み、オブジェクトフィンガープリントは、データオブジェクトのコンテンツから導出されたグローバルに一意的なフィンガープリントを含む、ステップと、

ファイルオブジェクトのためのオブジェクトフィンガープリントを生成するステップであって、各ファイルオブジェクトは、ファイル内の複数のデータオブジェクトのフィンガープリントを含み、ファイルオブジェクトのフィンガープリントは、ファイルオブジェクトのコンテンツから導出されたグローバルに一意的なフィンガープリントを含む、ステップと、

ファイルシステム内のオブジェクトフィンガープリントの全てのマッピングを含むルートオブジェクトを生成するステップとを含む。

【 0 0 2 5 】

－実施形態では、この方法は、

オブジェクトごとに参照カウントを維持し、オブジェクトへの参照が追加又は削除されたときに、オブジェクトの参照カウントを更新するステップを含む。

【 0 0 2 6 】

－実施形態では、この方法は、

読出し、書込み、削除及び参照カウント更新を含むオブジェクトアクティビティのトランザクションログを生成するステップを含む。

【 0 0 2 7 】

－実施形態では、この方法は、

ファイル内のデータオブジェクトを追加、修正又は削除し、新しいファイルオブジェクトのフィンガープリントを生成するステップを含む。

【 0 0 2 8 】

－実施形態では、

ファイルオブジェクト又はデータオブジェクトのコンテンツが変更されるとき、変更をルートオブジェクトまで伝えるステップである。

【 0 0 2 9 】

－実施形態では、この方法は、

10

20

30

40

50

伝えるステップを、
I/Oトランザクションごとに、
周期的に、
オンデマンドで、
特定のイベントで、
のうち1つで行うステップを含む。

【0030】

本発明の別の実施形態によればある方法が提供され、この方法は、
複数のデータオブジェクトを提供するステップであって、各データオブジェクトはデー
タ又はメタデータを含み、各データオブジェクトは、グローバルに一意であり該データオ
ブジェクトのコンテンツから導出されたフィンガープリントを有する、ステップと、

10

複数の関連付けられたデータオブジェクトのための複数のデータオブジェクトのフィン
ガープリントを含むファイルオブジェクトを生成し、グローバルに一意でありファイルオ
ブジェクトのコンテンツから導出されるファイルオブジェクトのフィンガープリントを生
成するステップと、

ファイルオブジェクトのフィンガープリントへのiノード番号のインデックスを維持す
るステップと
を含む。

【0031】

—実施形態では、この方法は、
オブジェクトフィンガープリント及びオブジェクトの物理位置をマッピングするための
位置インデックスを維持するステップを含む。

20

【0032】

—実施形態では、
位置インデックスは、オブジェクトのための参照カウントを含む。
—実施形態では、
フィンガープリント及びインデックスは、ファイルシステムを含む。

【0033】

—実施形態によれば、
プロセスによって実行されるとき、方法の請求項27のステップを行うプログラムコー
ド手段を含むコンピュータプログラム製品が提供される。

30

【0034】

本発明の別の実施形態によれば、
格納されたオブジェクトをインデックス付けする方法のための実行可能プログラム命令
を含むコンピュータ可読媒体が提供され、上記方法は、

グローバルに一意であり、データ及びメタデータのオブジェクトのコンテンツから導出
されたフィンガープリントを生成するステップと、

データ及びメタデータのうちの一方又は双方のオブジェクトの複数のフィンガープリン
トを含むファイルオブジェクトを生成し、ファイルオブジェクトのフィンガープリントを
生成するステップであって、ファイルオブジェクトのフィンガープリントは、グローバ
ルに一意でありファイルオブジェクトのコンテンツについて導出された、ステップと、

40

データ、メタデータ及びファイルのオブジェクトのフィンガープリントの全てのマッピ
ングを含むルートオブジェクトを生成するステップと
を含む。

【0035】

本発明の別の実施形態によれば、
データ、メタデータ及びファイルへのアクセスを提供する物理プロセッサ及びストレ
ージデバイスが提供され、

データ、メタデータ及びファイルはオブジェクトであり、各オブジェクトは、グロー
バルに一意でありコンテンツ導出のフィンガープリントを有し、オブジェクト参照はフィン

50

ガープリントによってインデックス付けされ、

インデックス付けは、ファイルオブジェクトのフィンガープリントへの i ノード番号のマッピングを含む。

【0036】

本発明の別の実施形態によれば、

データオブジェクト、及び、ファイルオブジェクトを含むデータオブジェクトの集まりを命名し格納するための処理及びストレージの装置であって、各データオブジェクトはデータ又はメタデータを含み、各オブジェクトは、コンテンツベースのグローバルに一意的なフィンガープリントを該オブジェクトのオブジェクト名として有し、ファイルオブジェクトは、データオブジェクト名の集まりであり、該ファイルオブジェクト自体のコンテンツベースのグローバルに一意的なフィンガープリントを該ファイルオブジェクトのファイルオブジェクト名として有する、処理及びストレージの装置と、

2つの層を有するファイルシステムであって、これら2つの層は、

オブジェクト名及び物理オブジェクト位置のマッピングインデックスを含むオブジェクトストア層と、

ファイルオブジェクトごとのデータオブジェクト名のマッピングインデックスを含む名前空間層と

を含む、2つの層を有するファイルシステムと

が提供される。

【0037】

一実施形態では、

名前空間層は、ファイルオブジェクト名への i ノード番号のマッピングインデックスを含む。

【0038】

一実施形態では、

オブジェクトストア層は、オブジェクトごとの参照カウントを含む。

一実施形態では、

オブジェクト名は、オブジェクトのコンテンツの暗号ハッシュダイジェストである。

【0039】

一実施形態では、

システムは、オブジェクトの命名、圧縮及び暗号化のうち1つ又は複数のために機能するハードウェアアクセラレーション装置を含む。

【0040】

一実施形態では、

オブジェクトストア層は、ファイルシステム内の全てのオブジェクトのグローバルインデックスを含み、グローバルなオブジェクトのインデックスのための主キーはオブジェクト名であり、オブジェクト名は、オブジェクトのコンテンツの暗号ハッシュダイジェストである。

【0041】

[007]本発明は、詳細な説明を以下の図面と共に参照することによって、より十分に理解されよう。

【図面の簡単な説明】

【0042】

【図1】[008]オペレーティングシステムカーネル空間に統合された本発明の一実施形態を例示する概略ブロック図である。

【図2】[009]オブジェクトストアが様々な物理媒体上でホストされることを可能とする、オブジェクトストアの一実施形態の主なコンポーネントの概略ブロック図である。

【図3】[0010]重要な機能性を、オブジェクトストア設計に影響を与えることなくこの機能が様々な方法で実装されることを可能としながら抽象化することが可能な、オブジェクトストアの一実施形態の概略ブロック図であり、実装は、純粋なソフトウェアソリューション

10

20

30

40

50

ションから、ハードウェアアクセラレーションを使用するものまで及ぶ場合がある。

【図4】[0011]統合ファイルシステムの基本構成要素として、コンストラクト(「hノード」)に共にグループ化されたオブジェクトのセットの一実施形態の概略ブロック図である。

【図5】[0012]ファイル、ディレクトリ及びiマップ等、ファイルシステムによって必要とされるような他のデータ構造に特化されうる、hノードの一実施形態の概略ブロック図である。

【図6】[0013]ファイルシステムに対する変更がどのように追跡され、経時的に維持されるか、並びに、使用された技術がどのように自然に空間効率、不変性及びセキュリティの結果となるかを例示する、一実施形態の概略ブロック図である。

10

【図7】[0014]オブジェクトのためのグローバル一意名を提供しながら、圧縮、暗号化及び位置独立を透過的に扱うことができるオブジェクトの、一実施形態の概略ブロック図である。

【図8】[0015]FUSE、ユーザ空間ファイルシステムにより、ユーザ空間内で実装された本発明の代替実施形態の概略ブロック図であり、FUSEは、ユーザ空間内でファイルシステムの構築を可能とする、ライブラリ及びカーネルモジュールのオープンソースセットである。

【図9】[0016]インデックス付けの発明の一実施形態により行われる、様々なインデックス付けオペレーションを例示する概略ブロック図である。

【図10A】[0017]この発明で使用されうるデータ構造の様々な実施形態を例示する図である。

20

【図10B】この発明で使用されうるデータ構造の様々な実施形態を例示する図である。

【図10C】この発明で使用されうるデータ構造の様々な実施形態を例示する図である。

【図10D】この発明で使用されうるデータ構造の様々な実施形態を例示する図である。

【図11】[0018]この発明の一実施形態によるルックアップオペレーションを例示する概略ブロック図である。

【図12】[0019]この発明の一実施形態による挿入オペレーションを例示する概略ブロック図である。

【図13】[0020]この発明の一実施形態による削除オペレーションの概略ブロック図である。

30

【図14】[0021]この発明の一実施形態による更新オペレーションの概略ブロック図である。

【図15A】[0022]この発明の一実施形態によるフリー消去ブロックを生成するためのランダム読み出し処理を例示する概略ブロック図である。

【図15B】この発明の一実施形態によるフリー消去ブロックを生成するためのランダム読み出し処理を例示する概略ブロック図である。

【図16A】[0023]スカベンジング処理に従ってフリー消去ブロックを生成する別の方法を例示する概略ブロック図である。

【図16B】スカベンジング処理に従ってフリー消去ブロックを生成する別の方法を例示する概略ブロック図である。

40

【図17】[0024]この発明の実装を例示するための6層のビュー又はスタックを例示する概略ブロック図である。

【図18】[0025]この発明の一実施形態で使用されるようなレコードエントリの概略図である。

【図19】[0026]図19Aは、この発明の一実施形態によるカッコウハッシング(cuckoo hashing)の実装を概略的に例示する図である。図19Bは、この発明の一実施形態によるカッコウハッシングの実装を概略的に例示する図である。図19Cは、この発明の一実施形態によるカッコウハッシングの実装を概略的に例示する図である。図19Dは、この発明の一実施形態によるカッコウハッシングの実装を概略的に例示する図である。図19Eは、この発明の一実施形態によるカッコウハッシングの実装を概略的に

50

例示する図である。

【図 2 0】[0027]この発明の一実施形態による、各バケットが複数のレコードを保持する、複数のバケットの概略的例示の図である。

【図 2 1】[0028]この発明の一実施形態によるバケットのコンテンツの概略図である。

【図 2 2】[0029]この発明の一実施形態による、複数のダイ、消去ブロック、ページ及びバケットを有する物理フラッシュチップの一例を例示する概略ブロック図である。

【図 2 3】[0030]図 2 3 A は、この発明の一実施形態によるデバイス管理層のあるコンポーネントを例示する図である。図 2 3 B は、この発明の一実施形態によるデバイス管理層のあるコンポーネントを例示する図である。

【発明を実施するための形態】

【0043】

A．従来技術の（レガシー）ファイルシステムの従来のファイルシステムデータ構造及び制限。

[0031]従来のファイルシステムは、幾つかの基本データ構造を有する。ユーザが見ることができるディレクトリ及びファイルに加えて、内部構造は、スーパーブロック、i ノード、アロケーションマップ、及びトランザクションログを含む。

【0044】

[0032]アロケーションマップは、ディスク上のどのブロックが使用中であるか、或いは使用中でないかを示すデータ構造である。これらのデータ構造は、ビットマップと同じくらい簡素になり、或いは、B ツリーと同じくらい複雑になりうる。アロケーションマップは大きくなり、ほとんどメモリに合わない可能性がある。新しいブロックの単純なアロケーションは、低いディスク性能という結果となるが、前述のメモリ制限を考えると、最適な配置には高度なアロケーションアルゴリズムを要する。

【0045】

[0033]ディレクトリは、ファイル及び他のディレクトリの名前のリストであり、多数のファイルシステムでは、単に異なるように解釈される別のファイルタイプとして扱われる。内部的には、ディレクトリは、ファイル名/i ノード番号の対のリストである。ファイルシステムは、ファイル名へのアクセスを望むとき、ディレクトリ内のファイル名、及び対応する i ノード番号を発見しなければならない。

【0046】

[0034]ファイルは、命名されたデータの集まりである。ファイル名は、それが参照する i ノードと共に、ディレクトリ構造に格納にされる。多数のファイルシステムは、異なるファイル名が同じデータ（i ノード）をポイントすることができる、リンクの概念をサポートする。

【0047】

[0035]トランザクションログは、原子性、一貫性、独立性及び持続性（ACID）特性に従って、ファイルシステムを一貫性のあるものに保つために使用される。多数のファイルシステムは、メタデータ一貫性を保証するが、データについての異なるサービスレベルアグリーメント（SLA）を有するようになる。

【0048】

[0036]スーパーブロックは、ディスク又は永続媒体上の既知の位置に存在する小さいデータ構造である。スーパーブロックから、i ノードテーブルのサイズ及び位置、アロケーションマップ、ルートディレクトリ等、ファイルシステムに関連する全ての他のデータ構造が発見されうる。ファイルシステムが搭載されるとき、最初にアクセスされるものがスーパーブロックである。安全性の理由のため、スーパーブロックはしばしば、ディスク上の様々なポイントでレプリケートされる。

【0049】

[0037]恐らく、最も基本的なデータ構造は、i ノード（「インデックスノード」）である。多数のファイルシステムに共通して、i ノードは、ファイル等、コンテンツのための基本コンテナであるデータ構造である。i ノード自体は、ファイル名を含まず、ファイル

10

20

30

40

50

名はディレクトリに格納される。i ノードは、ディスク常駐データ構造 (i ノードテーブル) へのインデックスを示す整数によって識別される。テーブル内の各 i ノードエントリは、このファイルについて、ディスク上のどこでコンテンツが発見されうるかを記述する。この「マップ」は、線形リスト、間接指定テーブル (i n d i r e c t i o n t a b l e)、様々なツリータイプを含む、様々な形式を取ることができ、その各々は様々な速度 / 空間トレードオフを有する。重要なことは、マップが、論理ブロック番号 (L B N) 等、物理又は論理アドレス指定を使用することである。L B N がどのディスクのために意図されるかを知る場合にのみ、L B N は意味をなす。

【 0 0 5 0 】

[0038] 上記から、レガシーファイルシステムが、何か (コンテンツ) 及びどこか (データの配置) についての厳重な制御を有することは、明らかになるべきである。この何か及びどこかが互いに混ざり合うこと、概して履歴の産物は、結果として、現代のストレージのニーズに届くことが困難であるアーキテクチャをもたらす。

【 0 0 5 1 】

B . 本発明のファイルシステムの新規のデータ構造及び特徴。

[0039] 本発明の様々な実施形態によれば、新しいデータ構造が、新しいタイプのファイルシステムを埋め込むために提供される。このファイルシステムは、他のファイルシステムと並んで存在し、働くことができ、レガシーファイルシステム及び既知のユーザレベルユーティリティと互換性がある。しかし、本発明の新しいデータ構造は、レガシーファイルシステムでは達成できない利益を提供する。これらの利益は、限定されないが、以下のうち 1 つ又は複数を含む。

【 0 0 5 2 】

- ・物理又は論理ブロックアドレス指定に頼らない、ファイルの命名及び格納のための抽象化のレベルを提供すること、

- ・データオブジェクトのコンテンツから導出されたグローバル一意フィンガープリントをオブジェクト名として利用することであって、各データオブジェクトはデータ又はメタデータを備えること、

- ・名前空間ファイルシステム内でオブジェクトフィンガープリントを利用することであって、全ての内部データ構造はオブジェクトであり、全てのオブジェクト間参照がオブジェクトフィンガープリントによって定義されることを可能とすること、

- ・「h ノード」構造と呼ばれる新しいデータ構造を提供することであって、ファイル h ノードは、ファイル内の全てのデータオブジェクトフィンガープリントのマッピング構造であり、それ自体は、ファイルオブジェクトのコンテンツから導出されたグローバル一意フィンガープリントを有するオブジェクトであること、

- ・同様に、ルート h ノード (オブジェクト) は、ファイルシステム内の全てのオブジェクトフィンガープリントのマッピング構造であり、ファイルシステムに対するいかなる変更もルートオブジェクト内の変更を生じ、且つ、ルートオブジェクトに対する変更を追跡することがファイルシステムアクティビティの履歴を提供するようにすること、

- ・ファイルオブジェクトフィンガープリントへの i ノード番号のマッピングを備える、i ノードマップオブジェクト (i マップ) を提供し、ファイルコンテンツが変化するときにはオブジェクト名 (フィンガープリント) が変化しながら、i ノード番号が一定であることを可能にすることであって、i ノードマップオブジェクトのフィンガープリントは、ファイルシステムのスナップショットを備えること。

【 0 0 5 3 】

[0040] 開示された実施形態では、オブジェクト名、即ち、オブジェクトフィンガープリントは、オブジェクトのコンテンツの暗号ハッシュダイジェストである。これは、オブジェクト名がオブジェクトコンテンツのフィンガープリントとしてグローバル一意且つ識別可能となることを可能とする。フィンガープリントは、オブジェクトよりかなり、例えば、100 倍、1000 倍以上も小さく、よって、フィンガープリントを操作することはしばしば、下にあるコンテンツを操作することよりも高速且つ容易である。

【 0 0 5 4 】

[0041]データオブジェクトの組み合わせ又は集まりを、同じくオブジェクトフィンガープリントであるオブジェクト名を有するオブジェクトである、hノードとして提供することによって、hノードはグローバル一意であり、hノード内に含まれたデータオブジェクトのコンテンツから導出される。いかなる変更（例えば、追加、削除、メタデータ変更、読出し）も、ファイルシステムhノードフィンガープリントが変更される結果となる。iマップに対する変更を追跡することによって、全てのファイルシステムアクティビティの完全な履歴が提供される。

【 0 0 5 5 】

[0042]本発明に固有のものは、iノードマップオブジェクト（別名iマップ）であり、iノード番号をオブジェクトフィンガープリントに変換するものである。これは、名前空間ファイルシステムがiノード番号に対処することを可能とし、多数のユーザレベルアクティビティがiノード番号を参照するので、中心のものである。iノード番号へのフィンガープリント（オブジェクト名）のhノードマッピングは、従来の静的なiノードテーブルを越えて、追加の間接指定（*i n d i r e c t i o n*）（又は仮想化）の層を提供する。この間接指定テーブルを使用することによって、iノード番号は一定であることができるが、iノードに対応するファイルが変化するとき、関連付けられたオブジェクト名（フィンガープリント）は変化することができる。iマップ自体がオブジェクトであるので、ファイルシステムが修正されるとき、その名前もまた変化ようになる。iマップのフィンガープリントは本質的に、ファイルシステムの完全な「スナップショット」である。スナップショットフィンガープリントを有した後、ファイルシステム（書込み可能なスナップ）に取り組むことを継続し、且つ、今後の使用のため（例えば、障害回復のため）にそれを記憶していることが可能である。スナップショットフィンガープリントを、別個のオブジェクトストア上に位置する別のシステムに発行することもできる。他のオブジェクトストアは、スナップショットデータ（オブジェクト）の全てを十分にホストしない場合があるが、記載された機構はなお十分に一貫性があり、使用可能である。

【 0 0 5 6 】

[0043]本発明のこれら及び他の利益は、本発明の様々な実施形態を参照して、より詳細に後述される。

[0044]カーネル空間及び次いでユーザ空間の両方で実装される、新しいファイルシステムの具体例を記載する前に、本実施形態で利用される様々なコンポーネントのより全体的な説明が定義される。

【 0 0 5 7 】

オブジェクトストア

[0045]オブジェクトストアは、本実施形態では、不透明（*o p a q u e*）データ（オブジェクト）のフラットな集まりである。各オブジェクトは固有であり、参照カウント（名前空間ファイルシステムによって参照される回数）を有する。オブジェクトの名前は、オブジェクトのコンテンツの暗号ハッシュであり、即ち、コンテンツを変更すると、名前は変化しなければならない。

【 0 0 5 8 】

[0046]任意の十分に強力な暗号ハッシュが、オブジェクト名（フィンガープリント）を生成するために受け入れ可能である。例として、セキュアハッシュアルゴリズム（*S H A*）ハッシュ関数は、国家安全保障局（*N S A*）によって設計され、*N I S T*によって米国連邦情報処理規格として発行された、暗号ハッシュ関数のセットである。*S H A - 1*は、既存の*S H A*ハッシュ関数の最もよく確立されたものであり、幾つかの広く使用されたセキュリティアプリケーション及びプロトコルで採用される。

【 0 0 5 9 】

[0047]実際には、オブジェクトサイズは典型的には2の累乗であり、512バイト（ 2^9 ）から最大1MB（ 2^{20} ）以上に及ぶが、オブジェクトのサイズに対するアーキテクチャ上の制約はない。

【 0 0 6 0 】

[0048]典型的なオブジェクトサイズは、2 KB (2^{11} バイト) である。8 TB (2^{43} バイト) のファイルシステムでは、 2^{32} ものオブジェクトであり、即ち、およそ 20 億ものオブジェクトである。インデックス内の各オブジェクトのエントリは約 $32 \times (2^5)$ バイトであり、そのため、オブジェクトインデックスは、密に詰まっていると仮定すると、 2^{37} 、即ち 128 GB、即ち全ファイルシステム空間の約 2% である。他のオブジェクトサイズは、適用性又は一般性を失わずに使用されうる。

【 0 0 6 1 】

[0049]オブジェクトは、オブジェクトのユーザに対して透過的に圧縮且つ暗号化される。オブジェクト名は、クリーンな非圧縮データ (及び任意選択のソルト) に基づく。オブジェクトに実際に格納されるものは、(クリーン)、(クリーン圧縮)、(クリーン、圧縮暗号化) 又は (クリーン暗号化) データのうち 1 つである。

10

【 0 0 6 2 】

[0050]オブジェクトは典型的には、クリーンデータのみで読み / 書きされ、圧縮 / 暗号化は、オブジェクトストアに内部で起こる。

[0051]強力な暗号ダイジェストを使用することは、オブジェクトがグローバル一意且つ一貫性のある名前を有することを可能とする。同じ名前を有する 2 つのオブジェクトは、実際上は、同じコンテンツを有するようになる。

【 0 0 6 3 】

名前空間

20

[0052]名前空間ファイルシステムは、本実施形態では、ファイル、ディレクトリ構造、リンク、スーパーブロック等を有する。

【 0 0 6 4 】

[0053]名前空間ファイルシステムは、データを直接的には含まず、代わりに全てのデータがオブジェクトに格納される。オブジェクトは比較的小さく、しばしばより大きいデータ構造が必要とされる。オブジェクトを集約する構造は、h ノードと呼ばれる。

【 0 0 6 5 】

[0054]実際の方法としては、Unix (登録商標) 又は Linux (登録商標) 環境にプラグインするファイルシステムは、i ノード番号をエクスポートする必要がある。i ノードは、ファイルを一意に識別する数である。

30

【 0 0 6 6 】

h ノード

[0055]h ノードは、本実施形態では、ファイル等、コンテンツを共に結び付けるデータ構造である。時として、コンテンツは大変大きくなる可能性があり (多 GB)、ディスク又は永続媒体上で連続的に適合しない。そのコンテンツは分割され、複数の個別のユニットとして格納される。従来のファイルシステムの場合、これはディスク上のブロックとなる。本発明では、これらがオブジェクト名である。h ノードは、マッピング構造内のオブジェクト名の全てのリストを保持する。線形リストは、そのようなマッピング構造の一例であるが、より複雑な間接指定テーブルもまた可能である。

【 0 0 6 7 】

40

[0056]h ノードと i ノードの間には、2 つの主な違いがある。第 1 に、h ノードはオブジェクトのコンテンツを識別するオブジェクト名 (フィンガープリント) を使用するのに対し、i ノードは物理又は論理ブロックアドレス指定を使用することである。第 2 に、h ノードは、明確に定義されたグローバル一意名 (そのコンテンツのハッシュ) を有することである。後述の好ましい実施形態では、h ノード名は、オブジェクトコンテンツのハッシュ及びソルトである。

【 0 0 6 8 】

i ノードマップオブジェクト (i マップ)

[0057]本発明に固有のものは、i マップであり、i ノード番号をオブジェクトフィンガープリント (名前) に変換するものである。このフィンガープリントは典型的には h ノー

50

ドであり、hノードは、コンテキストに応じて様々な方法で解釈される。これは、名前空間ファイルシステムの残りがiノード番号に対処することを可能とし、多数のユーザレベルユーティリティがそのようなコンストラクトを見ることを必要とするので、必須のものである。ある意味で、これは、従来の静的なiノードテーブルを越えて、追加の間接指定（又は仮想化）の層を提供する。

【0069】

[0058]この間接指定テーブルを使用することによって、iノード番号は一定であることができるが、iノードに対応するファイルが変化するとき、関連付けられたオブジェクト名（フィンガープリント）は変化することができる。iマップ自体がオブジェクトであるので、ファイルシステムが修正されるとき、その名前もまた変化するようになる。

10

【0070】

[0059]従来のファイルシステムでは、ルートディレクトリは既知のiノード番号にあり、iマップの場合もまた、そうである。

[0060]iマップのフィンガープリントを有する場合、本質的に、ファイルシステムの完全な「スナップ」を有する。このフィンガープリントの下のあらゆる可視のオブジェクトの参照カウントをバンプすることは、スナップをロックし、他のファイルシステムアクティビティにかかわらず、スナップが削除されることを妨げる。

【0071】

[0061]スナップフィンガープリントを有した後、ファイルシステム（書込み可能なスナップ）に取り組むことを継続し、且つ、今後の使用のために（恐らく、障害回復の目的で）それを記憶していることが可能である。スナップフィンガープリントを、別個のオブジェクトストア上に位置する別のシステムに発行することもできる。他のオブジェクトストアに気付く範囲内で、あるオブジェクトストアが特定のフィンガープリントの読み出し要求を解決することができない場合、そのオブジェクトストアは、その要求をそれらのストアに転送してもよい。よって、スナップのフィンガープリントは、そのオブジェクトストアがスナップのデータ（オブジェクト）の全てを十分にホストしない場合のあるシステムに移動することがあるが、上記の機構を介して、なお十分に一貫性があり、使用可能である。

20

【0072】

スーパーブロック

[0062]スーパーブロックは、本実施形態では、オブジェクトストアが永続媒体上で存続するとき、使用されるデータ構造である。スーパーブロックは、既知の（1つ又は複数の）位置内で存続する。スーパーブロックは、アロケーションマップ、iマップ、オブジェクトプール、インデックス及び他の構造が媒体上のどこで存続するかを記述する。オブジェクトストアは、常にグローバル意識別子（GUID）を有し、GUIDは、オブジェクトストアのその固有インスタンスを表現する。

30

【0073】

[0063]オブジェクトストアが大きいオブジェクトプールに参加する場合、スーパーブロックもまた、より大きいプールのGUID、及び、全てのメンバのGUID、及び、メンバの関係（取り除かれた、レプリケートされた、消去符号化された等）を含む。

40

【0074】

ファイル

[0064]ファイルコンストラクトは、本実施形態では、hノードから導出される。ファイルコンストラクトは、読み出し、書込み、開く、閉じる等、ファイルに関する通常（例えば、POSIX（登録商標））のセマンティクスの全てを有する。

【0075】

ディレクトリ

[0065]ディレクトリは、本実施形態では、hノードの特化されたバージョンである。ディレクトリは、（iノード番号、オブジェクト名）の対のマップを含む。線形リスト、ベクトル、又は他のより複雑な構造は、実装の例である。マップは少なくとも、hノードに

50

対してそれとして持続するために、直列化可能且つ非直列化可能でなければならない。マッピング構造に応じて、ランダムアクセスもまた可能である。

【0076】

追跡

[0066]ファイルシステムは、通常の手書き、削除及び読出しにより修正される（読出しがアクセス時間を変化させることを認める）ので、そのファイルシステムを構成するオブジェクト及びhノードもまた変化する。これは、結果としてルートハッシュの履歴となり、ルートハッシュの履歴は、大変細かい粒度では、継続的データ保護（CDP）と呼ばれ、より粗い粒度では、スナップと呼ばれる。違いは、どのくらいの頻度でルートハッシュが取り込まれるか、のみにある。

10

【0077】

[0067]システム内のあらゆるオブジェクトは、少なくとも1つのルートハッシュを通じてアクセス可能でなければならない。

[0068]本実施形態では、hノードとしてHが書き込まれ、新しいhノードH'が作成され、より多くの変化が生じる場合、場合によってはH"である。これらの変化は蓄積してもよいが、ある時点で、最後の変化がルートまで遡って伝わる。この保留中の入力/出力（IO）は、ファイルシステムが変化を蓄積し、且つ、変化ごとにルートまで伝わらないことを可能とする。どのくらいの頻度でこれが発生するかは、ポリシーに基づく。変化リストH -> H' -> H"の途中のオブジェクトのための参照カウントは、それに応じて対処され、ダングリング参照又は到達不可能オブジェクトがないようにしなければならない。

20

【0078】

C. 本発明の実施例（実装）

[0069]ここで図1を参照すると、図示されるものは、オペレーティングシステムカーネル101内の様々なストレージコンポーネントである。Linux（登録商標）環境をもとに描かれるが、この図は十分に一般的であるので、Windows（登録商標）、Solaris（登録商標）及び他のUnix（登録商標）クラスのオペレーティングシステム等、他のオペレーティングシステムに適用される。

【0079】

[0070]POSIX（登録商標）104のスタイルのファイルシステムの一実施例が示され、但し、POSIX（登録商標）は、一般性を失わずに、ResierFs、Extfs、btrfs及びzfs等、任意の数のファイルシステムのうち任意のものでありうる。仮想ファイルシステム（VFS）層103は、ファイルシステムの多数の共通機能を抽象化するために使用され、一貫性のあるインタフェース160をユーザ空間100及び他のコンポーネントに提供する。VFS103はまた、いかなるファイルシステムも使用しなければならない（VFS103層によって認識されることを期待する場合）、明確に定義された「より低いエッジ」のインタフェース150をも有する。実際には、典型的には、並行して働く多数のファイルシステムがある。

30

【0080】

[0071]ファイルシステムは、ブロックドライバ105によって実装された、ブロックストレージ抽象化の上に正常に位置する。ブロックストレージは、論理装置番号LUNローカルストレージデバイス109上にあってもよく、或いは、iSCSIプロトコルを使用してリモートLUN上にあってもよい。ブロックドライバ105もまた、オペレーティングシステム内で明確に定義されたインタフェースを有する。

40

【0081】

[0072]この実施形態では、新しいファイルシステムは、カーネル内の他のファイルシステムと並んで働く。新しいファイルシステムは、軽量オブジェクトファイルシステム108の上にスタックされる、名前空間ファイルシステム107からなる。2つのコンポーネント間のインタフェース152は、ANSI T-10オブジェクト標準等、様々な業界標準オブジェクトインタフェースの任意のものであってもよい。

50

【 0 0 8 2 】

[0073]オブジェクトファイルシステム（オブジェクトストア）108は次には、パーティションで区切られ、共通して使用される機能のライブラリである、ダイジェスト、インデックス付け、圧縮、暗号化（D I C E）ライブラリ310が抜き出されるようになる。ライブラリ310は、ソフトウェアにおいて完全に実現されてもよく、或いは、様々なハードウェアアクセラレーション113の技術をうまく利用してもよく、その1つが例示される。

【 0 0 8 3 】

[0074]オブジェクトファイルシステム108は、生のLUN、ディスク上のパーティション、又は大きいファイルの上に位置してもよい、オブジェクトコンテナを作成する。オブジェクトファイルシステム108はまた、i S C S I又は他のリモートアクセスブロックプロトコル（F C o Eは、もう1つの例である）等、プロトコルを使用して、ネットワークスタック106を介して、コンテナを参照してもよい。ネットワークファイルシステム（N F S）102は、（インタフェース154を介して）ネットワークスタック106の上に位置し、N F Sは、V F S 103に接続される。ネットワークスタック106は、インタフェース160を介してLUN109に、インタフェース159を介してクラウド110に接続される。

【 0 0 8 4 】

[0075]図2を参照すると、オブジェクトストア108は更に分解される。オブジェクトストア108は、バイナリの不透明オブジェクトを含み、その例は、P201、Q202及びR203である。オブジェクトは様々なサイズであってもよいが、好ましい実装では、2の累乗である。オブジェクトは、コンテナ内であるオフセットにて存在し、このオフセットは、バイトオフセット、又は、最小オブジェクトサイズを法とするオフセット（即ち、最小オブジェクトが512バイトである場合、バイトオフセットを得るために、オフセットに512が乗じられる）であってもよい。

【 0 0 8 5 】

[0076]各オブジェクトは、名前（フィンガープリント）を有し、名前（フィンガープリント）は、オブジェクトのコンテンツ全体の暗号ダイジェスト（ハッシュ）に、あるサイト特有のソルトを加えたものである。図2では、オブジェクト名は、H（P）、H（q）及びH（r）によって示される。

【 0 0 8 6 】

[0077]インデックス構造204は、オブジェクト名、オブジェクト位置、及びオブジェクト参照を追跡し続ける。オブジェクトの参照は、オブジェクトが書き込まれるたびに増分される。名前空間ファイルシステム107は、同じオブジェクトの多数のコピーであると思うものを生成してもよく、オブジェクトストア108は、1つのみを格納するが、名前空間がいくつを有すると実際に考えるかを追跡し続ける。

【 0 0 8 7 】

[0078]オブジェクトストア108は、幾つかのインタフェースクラスを有する。読出し、書込み、削除インタフェース152aは、オブジェクトに対してまさにそれを行う。これに関連するオブジェクト削除は、実際にはオブジェクトの参照カウンタの減分である。オブジェクトストア内部のオブジェクトのためのストレージは、参照カウンタが0になるときにのみ、解放される。

【 0 0 8 8 】

[0079]インデックス付けオペレーション152bは、名前によるオブジェクトの列挙、参照カウンタ調節、及び、名前によるオブジェクトのルックアップを可能とする。

[0080]オブジェクトストア108は、トランザクションのセマンティクス（A C I D特性）を有し、トランザクション境界は、トランザクションオペレーション152cを通じて管理される。これは、保留中のトランザクションをリストすることに加えて、トランザクションの開始、コミット及び中断を含む。

【 0 0 8 9 】

10

20

30

40

50

[0081]プロビジョニングインタフェース152dは、オブジェクトストアが作成され、削除され、マージされ、分割され、集約されることを可能とする。

[0082]インデックス204はマップであり、その主キーはオブジェクト名である。他のところで論じたように、インデックスは大変大きくなりうる。システム内のあらゆるオブジェクトのためのインデックスエントリがある。各エントリは、以下を含む。

【0090】

a) オブジェクトのコンテンツのフィンガープリント。フィンガープリントは、コンテンツ上の暗号ダイジェストによって生成され、少量の追加のコンテンツ(「ソルト」)が付加される。ソルトは、オブジェクトストア内の全てのオブジェクトに共通である。

【0091】

b) 何回オブジェクトが参照されるかを示す参照カウント。参照カウントは、空間を節約するために飽和演算(saturating arithmetic)を使用してもよい。例えば、参照を追跡するために8ビットのみを使用してもよく、参照カウントは追加且つ減分されうるが、255に等しくなるか、或いは255を上回る場合、カウントは「飽和」し、更なる減分は許可されない。

【0092】

c) 物理ロケータ。オブジェクトが物理ディスク上にある場合、これは論理ブロック番号LBNであってもよい。オブジェクトがホスト側プロバイダ(例えば、Amazon S3)によってホストされる場合、クラウドオブジェクトへの参照となる。

【0093】

d) 様々な使用のためのフラグ。あるフラグは、オブジェクトが格納され、圧縮されるかどうかを示し、別のフラグは、暗号化されるかどうかを示す。他のフラグは使用可能であるが、具体的な使用に割り振られない。

【0094】

[0083]アロケーションマップ220は、オブジェクトコンテナ206上で割り振られたブロックに対して使用された通常のビットマップである。

[0084]オブジェクトコンテナ206は、ランダムにアドレス指定可能な永続ストレージ抽象化である。例は、生のLUN、ファイル、ディスク上のパーティション、又は、広域ネットワークWANに渡るiSCSIデバイスを含む。

【0095】

[0085]オブジェクトコンテナ206は、幾つかのコンポーネント207~211(縮尺通りに図示されない)を有する。既知のオフセットにて存続するコンテナ記述子ブロック207を除いて、他のコンポーネントの順序は重要ではない。

【0096】

[0086]インデックス208は、Bツリー等、コンテナ常駐部分、若しくはメモリ204内の部分、又は両方を有してもよい。アロケーションマップ210もまた、部分的にディスク上且つメモリ220内であってもよい。2つの間のマイグレーションは、ページング技術により達成されうる。

【0097】

[0087]オブジェクトストアが修正されるとき、トランザクションログ211が永続ストレージ上で保持される。ログは、読出し、書込み、削除、参照調節等を含む、全てのオブジェクトアクティビティを追跡する。ログは、時間順に保持され、周期的にメインインデックス208にロールインされる。オブジェクトアクティビティは、メインインデックスを検索する前に、最初にログ上で「ヒット」しなければならない。各ログエントリは、オペレーションタイプ152a、152b、152c、152d、フィンガープリント、参照カウント、トランザクションID又はエポック番号、及びプール位置からなる。ログエントリは構造的に、トランザクションIDを追加したインデックスエントリに類似する。

【0098】

[0088]グローバルオブジェクト命名は、一貫性のある命名及びアクセスをなお保ちながら、オブジェクトストアがオブジェクトをあちこちに移動させることを可能とする。オブ

10

20

30

40

50

ジェクトを移動させるための理由は、以下を含む。

【0099】

a) 性能の理由のため、物理ディスク上で関連オブジェクトを互いの近くに移動させること。

b) 故障境界 (*fault boundaries*) に渡ってオブジェクトをレプリケートすること。これは、2つの別々のローカルディスク、ローカルディスク及びリモートディスク、又は、それらの任意の複数のものに渡る可能性がある。レプリケーションはまた、読出し性能利益を与えることもできる。レプリケーションはまた、消去コードによる等、オブジェクトを分割することを含むこともできる。

【0100】

c) 圧縮、解凍、暗号化、暗号解読等、オブジェクトにおけるバックグラウンドオペレーション。

d) 温度、即ち、それらの使用の頻度又は期待される頻度に基づいて、オブジェクトを移動させること。

【0101】

[0089] 図3は、オブジェクトストア108の、DICEライブラリ310との関係を例示する。ライブラリ310は、ダイジェスト153a、インデックス付け153b、圧縮153c及び暗号化153d等、オブジェクトストアの共通機能を抽象化する。

【0102】

[0090] 一貫性のあるインタフェースを提供しながら、内部的には、ライブラリは、様々な技術を使用してサービスを達成してもよい。実装技術は、ソフトウェアのみ、部分的ハードウェア支援 (例えば、Intel Quick Assist (登録商標))、又は、大量のインデックスを格納することができるカスタムハードウェア実装、又は、上記の任意の組み合わせを含む。

【0103】

[0091] ハードウェアアクセラレータ113を使用する場合、そのアクセラレータは2つの幅広いサービスのクラスを有してもよい。即ち、1つは、計算集中的なオペレーション用111 (圧縮、暗号化、フィンガープリンティング)、もう1つは、インデックス等、メモリ集中的なオペレーション用112である。ハードウェア実装は、一方若しくは他方又は両方を有してもよい。

【0104】

[0092] 図4は、本実施形態におけるhノード構造401の重要なコンポーネントを例示する。hノードは、レガシーiノードが使用する物理/論理ブロックアドレス指定ではなく、オブジェクト識別子 (フィンガープリント) を使用して、コンテンツを識別する。

【0105】

[0093] hノードは、ファイルのような、コンテンツのシーケンスであり、ランダムに読出し、書込み、付加、作成、削除及び切り捨て可能である。コンテンツは、任意のバイト境界上で、任意の範囲でアクセス可能である。どのようにコンテンツが解釈されるかは、コンテキストによって決まる。

【0106】

[0094] hノード401は、stat構造420、例えば、ファイルメタデータに対して使用されるPOSIX (登録商標) 構造を有してもよい。その構造の一部は、ファイル、即ち、この場合はhノードのバイト長を含んでもよい。データシーケンスは、個別のオブジェクト、例えば、図4のS410、T411及びU412に分割される。各オブジェクトの名前は、マッピングテーブル402に格納され、マッピングテーブル402はS、T及びUの各々のフィンガープリントを記録する。オブジェクトは、必ずしも同じ長さでなくてもよい。

【0107】

[0095] マッピングテーブル402は、一般性を失わずに、線形リスト、ツリー構造又は間接指定構造 (*an indirection structure*) を含む、様々な表

10

20

30

40

50

現を有してもよい。マッピングテーブル402は、標準のUnix(登録商標)のiノード間接指定テーブルが動く方法に類似の方法で、どの(1つ又は複数の)オブジェクトが参照されるべきかを決定するために、コンテンツ(シーケンスS、T及びU)へのオフセットによって、インデックス付けされる。

【0108】

[0096] hノード自体はオブジェクトであり、よって、固有の名前を有する。stat構造420、マッピングテーブル402、及び、参照されたオブジェクトの何れかのうち、何れか1つ又は複数が変化するとき、hノードの名前(フィンガープリント)もまた変化するようになる。

【0109】

[0097] hノードは、読出し、書込み及び付加の両方のために、ランダムにアクセスされる。hノードは、書き込まれていないデータが既知の値(典型的には0)を返す、疎空間(sparse space)をサポートする。

【0110】

[0098] hノードの名前はそのコンテンツの関数であるので、hノードに対するいかなる変更も、結果として新しいhノードとなる。元のhノードは、ファイルシステムポリシーに応じて、参照解除されても、(参照カウントを増すことによって)保持されてもよい。

【0111】

[0099] hノード401は、例えば、標準のUnix(登録商標)「stat」構造420に加えて、追加の構造を有してもよい。

[0100] 図5に示すように、hノード401は、ファイルに類似した、ランダムにアドレス指定可能なコンテンツのシーケンスである。どのようにそのコンテンツが解釈されるかは、コンテキストによって決まる。本発明の本実施形態では、hノードは更に、ファイル、ディレクトリ及びiマップに特化される。オブジェクト指向プログラミングの専門用語では、クラスファイル、ディレクトリ及びiマップは、ベースクラスhノードから派生される。

【0112】

[0101] ファイル504は、hノードが、開き、閉じ、読出し、書込み等が可能な通常のPOSIX(登録商標)ファイルのように見えるようにする、薄いラッパであってもよい。

【0113】

[0102] ディレクトリ505は、hノード401の別の解釈である。ディレクトリ505は、ファイル名(文字列)へのiノード番号(整数)のマッピング501である。マッピングは、限定されないが、線形リスト、Bツリー及びハッシュマップを含む、様々な形式を取ることができる。マップ501が完全にメモリ内にある場合、マップが直列化且つ非直列化されることが要件である。

【0114】

[0103] iマップ(「iノードマップ」)502は、(ディレクトリ501からの)iノード番号をオブジェクトダイジェスト(フィンガープリント)に変換する。オブジェクトは、hノード(及び、従って拡張して、ファイル、ディレクトリ又は他のiマップ)、スーパーブロック等の構造、又は他のデータを表現してもよい。

【0115】

[0104] iマップは、周知のオブジェクトのために、インデックス0、インデックス1等、位置を予約済みであってもよい。例は、(1つ又は複数の)以前のiマップ、ファイルシステムスーパーブロック等を含む。

【0116】

[0105] 図6は、コンテンツが追加される時、ファイルコンテンツ及びメタデータが最初の時間 T_0 から時間 T_1 から T_2 までにどのように変化するかを例示する。コンテンツの削除は、類似の経路を辿る。

【0117】

10

20

30

40

50

[0106]この図は、インタフェース152によって分離された、オブジェクトストア108コンポーネント、及び、名前空間107コンポーネントを共に示す。

[0107]時間 T_0 610で、 $Root_0$ ディレクトリ $Root_0$ 640は、2つのファイル FOO 641及び BAR 642を有する。ファイル FOO 641は、オブジェクト P 652及び Q 655に分割されたコンテンツからなる。 P 652及び Q 655のためのオブジェクト名は、 FOO 641の、以前に例示された(図4)マッピングテーブルに格納される。同様に、ファイル BAR 642は、コンテンツ Q 655を有する。ルートディレクトリ640もまた、 $Root_0$ 653によって示されたオブジェクトである。同様に、ファイル(hノード) FOO 641及び BAR 642は、それぞれオブジェクト651及び654において表現される。オブジェクト $Root_0$ 653を有するルートディレクトリ $Root_0$ 640がそうであるように、最初の $iMap_0$ 502aもまた、オブジェクト $iMap_0$ 650において表現される。

10

【0118】

[0108]オブジェクト Q 655は、ファイル FOO 641及び BAR 642の両方に共通であるので、参照カウントが2であるのに対して、オブジェクト P 652は、時間 T_0 610で参照カウントが1のみである。

【0119】

[0109]ルートディレクトリ640は、2つのエントリを含み、それぞれ FOO 及び BAR のものである。 FOO のエントリは、 i ノードインデックスが4であり、 BAR の i ノードインデックスは9である。

20

【0120】

[0110] $iMap_0$ 502aはhノードであり、オブジェクト650等として格納される。図面を複雑にすることを避けるため、 i マップはhノードであり、hノードは多数のオブジェクト上にマップしてもよいが、ここでは1つのオブジェクトとして図示される。

【0121】

[0111]慣例により、ルートディレクトリのダイジェストは常に、 i マップインデックス2に格納される。 i マップのダイジェストは、ファイルシステムへのフルアクセスを可能とする。 i マップに関連付けられたオブジェクトを読み出すことによって、ルートディレクトリが得られ、そこから、任意の後続のディレクトリ及びファイルのうち的一方又は双方が得られる。更に、 i マップのダイジェストは、ダウンストリームのファイルシステム全体のコンテンツを、正確に、且つ、曖昧でなく定義する。

30

【0122】

[0112]不変性：例えば、オブジェクト Q が変化する場合、名前が変化する(オブジェクトの名前は、そのコンテンツの関数である)。修正された Q にポイントするいかなるマッピングテーブルも、このときはポイントせず、従って、修正された Q は「可視」ではない。同様の引数は、マップのダイジェストによって参照可能であるいかなるオブジェクトにも適用される。

【0123】

[0113]時間 T_1 611で、ファイル BAR 642にはコンテンツ S 658が付加されているので、新しいファイル BAR 644が作成される。新しいファイル BAR は、ダイジェスト及びオブジェクト名が一貫性のあるものになるように、作成されなければならない。新しいコンテンツ S 658が追加されるので、それを参照するあらゆるものもまた更新され、新しいバージョンが作成される。これは、より新しいバージョンの BAR 644、ルートディレクトリ643、及び最も重要なことには、新しい i マップテーブル502bに当てはまる。時間 T_0 610でのオブジェクト参照カウント614は、コンテンツが追加/除去されるときに調節されるので、時間 T_1 で、 T_1 オブジェクト参照カウント615は、 T_0 に固有である、 T_1 に固有であるコンテンツ、及び、共通であるコンテンツを表現する。

40

【0124】

[0114]時間 T_1 611で、本質的に、多くの共通コンテンツを有する2つのファイル

50

システムがある。これらの2つのファイルシステムは、それらのそれぞれのiマップ、i map₀、502a及びi map₁、502bのダイジェストによって十分に特定される。例えば、時間T₀、610で、オブジェクトQ655は、経路(640a、641b)、(640b、642a)、(643a、641b)及び(643b、644a)を通じて参照されうる。

【0125】

[0115]ファイルのコンテンツが修正(追加、削除、修正)されるとき、ファイルのマッピングテーブルもまた変更される。次には、ファイルマッピングを含むオブジェクト、hノードもまた変化する。様々な理由(性能、管理インタフェース)のため、あらゆる変化をツリー中ずっと、ルートディレクトリまで、且つiマップへと伝わることは、適切でない場合がある。しかし、あらゆるI/Oトランザクション上で行われる場合、システムはCDPを暗黙的に実装し、CDPではiマップのあらゆるダイジェストが特定のI/Oトランザクションを表現する。周期的に(例えば、毎時等)、オンデマンドで、或いは、特定のイベント(ファイル閉じ)において行われる場合、ビヘイビアはファイルシステムスナップショットに類似する。

10

【0126】

[0116]同一のオブジェクトがある程度まで、オブジェクトが参照カウントを有するとき、重複排除はシステムに固有である。ファイルシステムが、修正の結果として変化するとき、大部分では、変化のみが、新しいコンテンツがストレージプールに追加される結果となる。

20

【0127】

[0117]図7では、オブジェクト701は、不透明のバイナリデータのシーケンスであり、オブジェクト名703を有する。オブジェクトサイズは任意であるが、好ましい実装では、割振り及び永続ストレージ管理をより容易にするために、2の累乗のサイズであってもよい。

【0128】

[0118]オブジェクトのユーザに対して、コンテンツは常に、クリーンオブジェクトコンテンツ710として読出し、書込み且つアクセスされる。オブジェクトストアは、任意選択の圧縮711及び暗号化712のうち的一方又は双方を含んでもよい形式で、オブジェクトを内部で格納する。よって、ユーザにとって2048バイトのオブジェクトのように見える可能性のあるものは、(4:1の圧縮比を仮定すると)512バイトのデータとして内部で格納され、このデータは更に暗号化される。オブジェクトストアは暗号化ドメインであり、全てのオブジェクトが暗号化に関して同様に扱われることを意味する。これは、オブジェクトの呼出元が使用する場合のあるいかなる暗号化とも別個である。

30

【0129】

[0119]図8では、ユーザ空間100で実装される、ファイルシステムの代替実施形態が例示される。オープンソースFUSEフレームワークを使用して、名前空間ファイルシステム107は、ユーザモードFUSEライブラリ802に対してリンクされる。名前空間ファイルシステムは、オブジェクトストア108への同じプライベートインタフェース152を有する。オブジェクトストア108もまた、DICEライブラリ310への同じインタフェース153を有する。DICEライブラリは、ハードウェア支援113を任意選択で使用してもよい。垂直ファイルシステム103は、カーネル101内に存在し、それはハードウェア支援113及び(VFS103及びFUSEライブラリ802の両方に接続された)FUSEモジュール801と同様である。

40

【0130】

D. 代替実施形態

[0120]機能の組み合わせを統合するファイルシステムを、従来のシステムのコストの何分の1かで構築する新規の方法が上述された。様々な修正は、代替実施形態の構築時に当業者には明らかになるであろう。

【0131】

50

[0121]新しいファイルシステムは、任意の他のファイルシステムとしてコンピュータ上で実行する、純粋なソフトウェア形式で実現されうる。更に、統合ファイルシステムの編成は、レガシーファイルシステムでは可能ではない、固有のハードウェアアクセラレーション技術に役立つ。ハードウェアアクセラレーションは、所与のコストに対してより高い性能、又は、所与の性能レベルに対してより低い所有コスト総額を可能とする。

【0132】

[0122]上記の実施形態では、このファイルシステムは、統合された機能セットを提供する。このファイルシステムは、2つの別個のファイルシステムである、オブジェクトファイルシステム及び名前空間ファイルシステムを含む、スタックとして実装される。スタックは完全にPOSIX（登録商標）準拠であり、第2の拡張ファイルシステム（EXT2）10、第3の拡張ファイルシステム（EXT3）、ReiserFs等、POSIX（登録商標）準拠ファイルシステムが求められる場合は常に使用されうる。

【0133】

[0123]スタックのより下位の部分は、オブジェクトファイルシステムである。オブジェクトベースのファイルシステムは、オブジェクトの形式でデータをホストするために使用される。オブジェクトは、不透明のバイナリデータのシーケンスである。オブジェクトは、生データであっても、メタデータ（例えば、生データの作成のレコード、及び、生データに対する任意の変更）であってもよい。オブジェクトサイズは様々でありうるが、典型的には、数キロバイト（KB）の範囲に制限される。しかし、これは、本発明の訂正オペレーションには必要とされない。オブジェクトの名前（本明細書では、フィンガープリントとも呼ばれる）は、例えば、強力な暗号ハッシュを使用して、オブジェクトのコンテンツから導出される。これは、オブジェクト名がグローバル一意且つ識別可能であること、20、即ち、コンテンツのフィンガープリントを可能とする。オブジェクトファイルシステムは、主としてマシン指向である。

【0134】

[0124]等しい2つのフィンガープリントは、どこでそれらのフィンガープリントが計算されたかにかかわらず、実際には同じコンテンツを表現するようになる。反対に、異なる2つのフィンガープリントは、異なるコンテンツを表現する。フィンガープリントは、オブジェクトよりかなり（例えば、100倍、1000倍以上も）小さいので、フィンガープリントを操作することはしばしば、下にあるコンテンツを操作することよりも高速且つ30、容易である。

【0135】

[0125]上記の実施形態に記載されたオブジェクトファイルシステムは、軽量且つフラットであり、ANSI T-IO仕様に記載されたもの等、重量のオブジェクトファイルシステム、又は、市販のEMC Centera（登録商標）、若しくはHitachiの製品（Archivasを介して取得）等、コンテンツアドレス指定可能ファイルシステムとは異なる。本明細書で使用されるようなオブジェクトは、C++及びJava（登録商標）等、プログラミング言語で使用されるようなオブジェクトと混同されるべきではない。

【0136】

[0126]オブジェクトファイルシステムは、オブジェクトの全てを追跡する「インデックス」を有する。そのようなインデックスの構築及び管理は、何百万又は数十億ものエントリさえもがインデックス内にある可能性がある場合、オブジェクトファイルシステムにとって大きな課題となりうる。

【0137】

[0127]記載された実施形態によれば、ストレージスタックの上に、ファイル、ディレクトリ等を有する名前空間ファイルシステムが提供される。既知のもの（例えば、POSIX（登録商標）ファイルシステム）との違いは、しかし、コンテンツにアクセスするために、論理ブロック番号アドレス指定（LBN）を使用する代わりに、オブジェクトフィンガープリントが使用されることである。更に、名前空間ファイルシステムの全ての内部デ40、

10

20

30

40

50

ータ構造は、それら自体でオブジェクトである。よって、ストレージスタック全体（名前空間及びオブジェクト層）は、オブジェクト参照によって共に「結び付け」られ、ルートを表現するオブジェクトのフィンガープリントを有することは、ファイル構造全体を完全に、且つ曖昧でなく定義することを可能とする。

【0138】

[0128]あらゆる変更（追加、削除、メタデータ変更、読出し）も、ファイルシステムのシグネチャが変更される結果となる。ルートシグネチャを追跡することによって、従って、全てのファイルシステムアクティビティの完全な履歴を得ることができる。

【0139】

[0129]本発明の開示された実施形態によれば、2つの別々のコンポーネント（名前空間107及びオブジェクトストア108）への仕事の分割、及び、これらのコンポーネントがどのように対話するかは、POSIX（登録商標）セマンティクスをなお保ちながら、重複排除、スナップ、書込み可能スナップ、継続的データ保護（CDP）、広域ネットワーク効率、バージョニング、ファイルシステム整合性チェック及び不変性が自然に起こるような方法で行われる。

10

【0140】

[0130]開示された実施形態によれば、ファイルシステムの編成は、ハードウェア支援の適用を可能とする。ハードウェア支援は、2つの形式を取ってもよい。1つの形式は、圧縮、暗号化及び暗号ダイジェスト等、計算アクセラレーション用である。第2の形式は、実際のオブジェクトストアを構築するために使用される、大きいインデックスのための構築及び維持用である。

20

【0141】

[0131]かなりのCPU資源が、暗号ハッシング、圧縮及び暗号化で費やされる。より高速なCPUクロック及びより多くのCPUコアは、これをある程度まで軽減するが、性能要件が増大するにつれて、これらの機能の一部又は全部を専用ハードウェア（アクセラレーション）へオフロードすることが望ましい。これを達成することができる、幾つかの市販のチップセット（例えば、HiFive、Cavium）がある。

【0142】

[0132]オブジェクトストアインデックスは大きくなる可能性があり、実際的なメモリ制限をすぐに上回る場合がある。ランダムに読出し及び書込みされる（そのようなインデックスのための主キーは暗号ハッシュであり、ランダムな分布を有する）、グローバルオブジェクトインデックス（即ち、ストレージの全てのためのインデックス）は、ページング及びキャッシングアルゴリズムを効果のないものにする場合がある。そのようなインデックスを、ソリッドステートディスク（SSD）等、より高速の不揮発性ストレージ上に置くことは、よって、性能利益を提供することになる。

30

【0143】

[0133]SSDは、読出し速度が書込み速度よりかなり速い（即ち、Seagate xxxは、35,000 iops / 読出し且つ3000 iops / 書込みを達成することができる）ように構築される。インデックスアクセスが読出しと書込みの間で均等に分割される場合、SSDの利益の多くは実現されない。

40

【0144】

[0134]FLASH及びFPGAでできたカスタムビルドのインデックス付けソリューションは、インデックス付け帯域幅を更により増大させることができる。310

[0135]ハードウェア支援は、DICEライブラリによって、予め記載されたように、管理されうる。

【0145】

[0136]本発明の実施形態は、デジタル電子回路において、或いは、コンピュータハードウェア、ファームウェア、ソフトウェアにおいて、或いは、それらの組み合わせにおいて実装されうる。本発明の実施形態は、コンピュータプログラム製品として、即ち、情報キャリア内で、例えば、マシン可読ストレージデバイス内で、データ処理装置、例えば、プ

50

ログラマブルプロセッサ、コンピュータ又は複数のコンピュータによる実行のため、或いはそれらのオペレーションを制御するために、有形に実施されたコンピュータプログラムとして実装されうる。コンピュータプログラムは、コンパイラ型又はインタープリタ型言語を含む任意の形式のプログラミング言語で書かれてよく、スタンドアロンのプログラム又はモジュール、コンポーネント、サブルーチン若しくはコンピューティング環境で使用するために適した他のユニットを含む任意の形式で配置されうる。コンピュータプログラムは、1つのコンピュータ又は複数のコンピュータ上で実行されるように配置することができ、これら複数のコンピュータは、1つの場所にあるか又は複数の場所に渡って分散され通信ネットワークによって相互接続される。

【0146】

10

[0137]本発明の実施形態の方法ステップは、1つ又は複数のログラマブルプロセッサがコンピュータプログラムを実行して、入力データに関して動作し、出力を生成することにより、本発明の機能を行うことによって、行われうる。方法ステップはまた、専用論理回路、例えば、FPGA（フィールドプログラマブルゲートアレイ）又はASIC（特定用途向け集積回路）によって行われてもよく、本発明の装置は、専用論理回路、例えば、FPGA又はASICとして実装されうる。

【0147】

[0138]コンピュータプログラムの実行に適したプロセッサは、例として、汎用及び専用マイクロプロセッサ、及び、任意の種類デジタルコンピュータの任意の1つ又は複数のプロセッサを含む。一般に、プロセッサは、命令及びデータを、読出し専用メモリ若しくはランダムアクセスメモリ又は両方から受信するようになる。コンピュータの必須要素は、命令を実行するためのプロセッサ、並びに、命令及びデータを格納するための1つ又は複数のメモリデバイスである。一般に、コンピュータはまた、データを格納するための1つ又は複数の大容量ストレージデバイス、例えば、磁気、光磁気ディスク若しくは光ディスクをも含むか、或いは、それらからデータを受信し、若しくはそれらへデータを転送するように動作可能に結合されるか、或いはその両方となる。コンピュータプログラム命令及びデータを実施するために適した情報キャリアは、例として、例えば、EPROM、EEPROM及びフラッシュメモリデバイス等、半導体メモリデバイスと、例えば、内蔵ハードディスク又はリムーバブルディスク等、磁気ディスクと、光磁気ディスクと、CD-ROM及びDVD-ROMディスクとを含む、全ての形式の不揮発性メモリを含む。プロセッサ及びメモリは、専用論理回路によって補われるか又は専用論理回路に組み込まれる。

20

30

【0148】

E . h ノード（コンテンツにソルトを加えたもの）

[0139]一実施形態では、hノード名は、そのコンテンツのハッシュにソルトを加えたものである。ソルトは、8から32バイト程度の小さい値であり、シグネチャが計算される前に、内部で自動的にあらゆるオブジェクトの前に付加又は後に付加される。オブジェクトが書き出されるとき、ソルトは格納されない。

【0149】

[0140]例えば、ユーザがパスワードを入力し、そこから、暗号鍵生成のために使用される様々な標準技術の任意のものを使用して、ソルトが生成される。ユーザは、任意の他のパスワードのように、このパスワードを保護するようになる。人がその販売を得る場合であっても、元のパスワードを生成することは計算的に可能ではない。

40

【0150】

[0141]ソルトは、主として、不正が行われたデータに対する防御機構であり、この実施例では、以下の通りである。

- あらゆるオブジェクトは、そのフィンガープリントを使用して追跡される、
- フィンガープリントは、インデックス内に格納される、
- フィンガープリントの配布はフラットであり、即ち、全てのフィンガープリントは等しく発生する可能性がある、

50

- インデックス付けアルゴリズムは、鍵（フィンガープリント）の一律の配布を前提とする。

【0151】

[0141]ファイルシステムが、例えば、SHA-1等、特定のフィンガープリントアルゴリズムを使用することを、悪意のあるエンティティが知る場合、このエンティティは、大変狭い範囲に入るフィンガープリントを有するコンテンツを容易に生成することができる。そうするために、このエンティティは、ランダムコンテンツを生成し続け、それをフィンガープリントし、特定の狭い範囲に入るコンテンツのみを保持する。それは、インデックス付けアルゴリズムが大変低い性能を有するようにさせることとなる。

【0152】

[0142]しかし、暗号ハッシュの性質は、オブジェクトのコンテンツのわずか1ビットを変更する場合、フィンガープリントのビットのおよそ50%が変化するようになるものである。元のコンテンツの異なるビットを変更するとき、どの50%であるかもまたランダム化される。

【0153】

[0143]ソルトを追加すること（即ち、比較的小さい変更）は、よって、フィンガープリントをランダム化し、インデックス付けアルゴリズムを「操る」ことを非常に困難にする。

【0154】

G. スケーラブルなインデックス付け（実施形態）

[0144]本発明の方法及び装置は、本願と同日（2010年6月25日）に出願され、2009年6月26日出願の米国仮特許出願第61/269,633号に対する優先権を主張する、本願と同じ発明者であるP. Bowden及びA. J. Beaversonによる、同時係属の、所有者が同じである「Scalable Indexing（スケーラブルなインデックス付け）」という名称の米国特許出願第12/823,452号に記載された、以下のインデックス付けアルゴリズム及びメモリ技術と共に実装されうる。両出願に対する優先権が本明細書で主張され、各々の完全な開示はそれらの全体として参照により本明細書に組み込まれる。

【0155】

[0145]前述及び以下の記載は、本発明の範囲を例示するためのものであり、限定するように意図されないことを理解されたい。

1) 概要

[0146]本発明の一実施形態によれば、ユニフォームアクセスインデックス付け処理によって、ノンユニフォームアクセスメモリに格納されたインデックスにアクセスする方法が提供され、この方法は、

インデックス付け処理によって生成された論理バケット識別子を、メモリの物理バケット位置にマップして、インデックス内の各レコードデータエントリにアクセスするために、変換テーブルを維持するステップと、

インデックスに書き込まれる複数のレコードデータエントリをキャッシュ内に集めるステップであって、該ステップは、メモリの少なくとも1つの物理バケット位置へのエントリの集まりの後続の順次書込みに先立つ、ステップとを含む。

【0156】

一実施形態では、この方法は、レコードデータエントリの集まりを、キャッシュからメモリのバケット位置へ、順次書込みとして書き込むステップと、

集まりのレコードデータエントリのためのバケット位置により、変換テーブルを更新するステップとを含む。

【0157】

－実施形態では、この方法は、
 1つ又は複数の順次レコードデータエントリを、メモリからキャッシュへ読み出すステップと、
 1つ又は複数のエントリが読み出されたメモリ内の物理位置を、フリーとして指定するステップと
 を含む。

【0158】

－実施形態では、この方法は、
 ブロック内の任意の有効エントリをキャッシュに読み出すこと、及び、そのようなエントリが読み出されたメモリ内の物理位置を、フリーとして指定することによって、メモリ内の複数の順次物理バケット位置をフリーブロックとするステップを含む。

10

【0159】

－実施形態では、
 インデックス付け処理が、一律に配布され一意のインデックスキーに基づいて、インデックスへのランダムアクセス要求を生成する。

【0160】

－実施形態では、
 キーは、暗号ハッシュダイジェストを含む。
 －実施形態では、
 インデックス付け処理は、置き換えハッシングの処理を含む。

20

【0161】

－実施形態では、
 置き換えハッシングは、カックウハッシング処理を含む。
 －実施形態では、
 メモリは、フラッシュ、相変化及びソリッドステートディスクのメモリデバイスのうちの1つ又は複数を含む。

【0162】

－実施形態では、
 メモリは、ランダム書込みアクセス時間、ランダム読出し - 修正 - 書込みアクセス時間、順次書込み、アラインメント制約、消去時間、消去ブロック境界及び摩耗のうちの1つ又は複数によって制限される。

30

【0163】

－実施形態では、
 物理バケットのサイズは、メモリの最小書込みサイズを含む。
 －実施形態では、
 物理バケットのサイズは、ページ又は部分ページを含む。

【0164】

－実施形態では、
 メモリは、複数のページを含む消去ブロックを有する。
 －実施形態では、この方法は、
 メモリ内のどのバケット位置が有効であるかを追跡するためのバケット有効テーブルを維持するステップを含む。

40

【0165】

－実施形態では、
 メモリ内のバケットは、1つ又は複数のレコードデータエントリのセット、及び、バケット変換テーブルへのセルフインデックスを含む。

【0166】

－実施形態では、
 バケット内のレコードデータエントリは順序付けされない。
 －実施形態では、この方法は、

50

メモリに順次書き込まれたレコードデータエントリを、キャッシュ内で読出し専用として指定するステップを含む。

【0167】

—実施形態では、

バケット変換テーブルは、永続メモリに格納される。

—実施形態では、この方法は、

消去ブロック内のフリーバケットの数を追跡し、フリーバケットのしきい値が満たされたとき、フリー消去ブロックを生成するための処理を実施するステップを含む。

【0168】

—実施形態では、

インデックス付け処理は、レコードが挿入、削除、ルックアップ及び修正のうちの1以上をされる要求に基づいて、インデックス付けオペレーションを行う。

【0169】

—実施形態では、

インデックス付け処理は、インデックスのレコードを格納する物理バケットを読み出し該物理バケットに書き込むための、論理バケットオペレーションを提示する。

【0170】

—実施形態では、

物理バケットオペレーションは、ランダム読出し及び順次書込みを含む。

—実施形態では、

物理バケットオペレーションは、トリムコマンドを更に含む。

【0171】

—実施形態では、

メモリは、ノンユニフォーム読出し及び書込みアクセス、並びに、サイズ、アラインメント及びタイミングについての不変性によって特徴付けられた物理デバイス層を含む。

【0172】

—実施形態では、

レコードデータエントリは、キー、参照カウント及び物理ブロックアドレスのためのフィールドを含む。

【0173】

—実施形態では、

キーは、データの暗号ハッシュダイジェストを含み、

物理ブロックアドレスフィールドは、ストレージデバイス上に格納されたデータの物理ブロックアドレスへのポインタを含む。

【0174】

—実施形態では、

論理バケット位置は、複数のハッシュ関数によって生成される。

—実施形態では、

メモリは、複数の消去ブロックを含むフラッシュメモリデバイスを含み、各消去ブロックは複数のページを含み、各ページは複数のバケットを含む。

【0175】

本発明の別の実施形態によれば、

プロセッサによって実行されるとき、前述の方法のステップを行うプログラムコード手段を含むコンピュータプログラム製品が提供される。

【0176】

本発明の別の実施形態によればコンピュータ可読媒体が提供され、このコンピュータ可読媒体は、ユニフォームアクセスインデックス付け処理によってノンユニフォームアクセスメモリに格納されたインデックスにアクセスする方法のための実行可能プログラム命令を含み、上記方法は、

インデックス付け処理によって生成された論理バケット識別子を、メモリの物理バケッ

10

20

30

40

50

ト位置にマップして、インデックス内の各レコードデータエントリにアクセスするために、変換テーブルを維持するステップと、

インデックスに書き込まれる複数のレコードデータエントリをキャッシュ内に集めるステップであって、該ステップは、メモリの少なくとも1つの物理バケット位置へのエントリの集まりの後続の順次書込みに先立つ、ステップとを含む。

【0177】

本発明の別の実施形態によれば、物理プロセッサと、コンピュータ可読媒体を含むメモリデバイスとを含むシステムが提供され、

該コンピュータ可読媒体は、ユニフォームアクセスインデックス付け処理によってノンユニフォームアクセスメモリに格納されたインデックスにアクセスする方法のための実行可能プログラム命令を含み、上記方法は、

インデックス付け処理によって生成された論理バケット識別子を、メモリの物理バケット位置にマップして、インデックス内の各レコードデータエントリにアクセスするために、変換テーブルを維持するステップと、

インデックスに書き込まれる複数のレコードデータエントリをキャッシュ内に集めるステップであって、該ステップは、メモリの少なくとも1つの物理バケット位置へのエントリの集まりの後続の順次書込みに先立つ、ステップとを含む。

【0178】

一実施形態では、

インデックスを格納するメモリは、ノンユニフォーム読み出し及び書き込みアクセス、並びに、サイズ、アラインメント及びタイミングについての不変性によって特徴付けられた物理デバイス層を含む。

【0179】

一実施形態では、

インデックスを格納するメモリは、フラッシュ、相変化及びソリッドステートディスクのメモリデバイスのうちの1つ又は複数を含む。

【0180】

一実施形態では、

インデックスを格納するメモリは、複数の消去ブロックを含むフラッシュメモリデバイスを含み、各消去ブロックは複数のページを含み、各ページは複数のバケットを含む。

【0181】

本発明の別の実施形態によれば、ユニフォームアクセスインデックス付け処理によって、ノンユニフォームアクセスメモリに格納されたインデックスにアクセスする方法が提供され、この方法は、

インデックス内のレコードデータエントリごとに論理バケット識別子をメモリの物理バケット位置にマップする変換テーブルに、インデックス付け処理によって生成された論理バケット識別子を提供するステップと、

論理バケット識別子にマップされた物理バケット位置にアクセスするステップと、

インデックスに書き込まれるレコードデータエントリをキャッシュ内に集めるステップと、

続いて、レコードデータエントリの集まりを、キャッシュから、メモリの少なくとも1つの新しい物理バケット位置におけるインデックスに順次書き込むステップと、

変換テーブルを更新して、少なくとも1つの新しい物理バケット位置を論理バケット識別子に関連付けるステップと

を含む。

【0182】

本発明の別の実施形態によればあるコンピュータシステムが提供され、このコンピュータシステムは、

ノンユニフォームアクセスメモリであって、メモリの物理バケット位置においてレコードデータエントリを含むインデックスが格納される、ノンユニフォームアクセスメモリと

、
ユニフォームアクセスインデックス付け処理によって生成された論理バケット識別子を、レコードデータエントリの各々のためのメモリの物理バケット位置にマップするための、変換テーブルと、

インデックスに書き込まれる集められたレコードデータエントリのためのキャッシュと

、
メモリの物理バケット位置にアクセスする手段であって、該物理バケット位置は、インデックス付け処理によって変換テーブルに供給された論理バケット識別子にマップされた

10

手段と、
レコードデータエントリの集まりを、キャッシュから、メモリの少なくとも1つの物理バケット位置のインデックスに順次書き込む手段と、

変換テーブルを更新して、少なくとも1つの物理バケット位置を論理バケット識別子に関連付ける手段と
を含む。

【0183】

2) 図面

[0147]インデックス付けの発明は、以下の図面と共に、詳細な説明を参照することによって、より十分に理解される。

20

【0184】

図9は、本発明の一実施形態により行われる、様々なインデックス付けオペレーションを例示する概略ブロック図である。

図10Aないし10Dは、本発明で使用されうるデータ構造の様々な実施形態を例示する図である。

【0185】

図11は、本発明の一実施形態によるルックアップオペレーションを例示する概略ブロック図である。

図12は、本発明の一実施形態による挿入オペレーションを例示する概略ブロック図である。

30

【0186】

図13は、本発明の一実施形態による削除オペレーションの概略ブロック図である。

図14は、本発明の一実施形態による更新オペレーションの概略ブロック図である。

図15A及び15Bは、本発明の一実施形態によるフリー消去ブロックを生成するためのランダム読み出し処理を例示する概略ブロック図である。

【0187】

図16A及び16Bは、スキャベンジング処理に従ってフリー消去ブロックを生成する別の方法を例示する概略ブロック図である。

図17は、本発明の実装を例示するための6層のビュー又はスタックを例示する概略ブロック図である。

40

【0188】

図18は、本発明の一実施形態で使用されるようなレコードエントリの概略図である。

図19A~19Eは、本発明の一実施形態によるカックウハッシングの実装を概略的に例示する図である。

【0189】

図20は、本発明の一実施形態による、各バケットが複数のレコードを保持する、複数のバケットの概略的例示の図である。

図21は、本発明の一実施形態によるバケットのコンテンツの概略図である。

【0190】

図22は、本発明の一実施形態による、複数のダイ、消去ブロック、ページ及びバケッ

50

トを有する物理フラッシュチップの一例を例示する概略ブロック図である。

図23A～図23Bは、本発明の一実施形態によるデバイス管理層のあるコンポーネントを例示する図である。

【0191】

3. 概観

[0148]本発明の1つ又は複数の実施形態によれば、特化されたメモリ技術及びアルゴリズムが使用されて、多数のレコード及びトランザクション要件を同時に有するインデックスが構築される。一実施形態は、置き換え(displacement)ハッシングインデックス付けアルゴリズム、例えばカックウハッシングを利用する。本発明は、フラッシュ、相変化及びソリッドステートディスク(SSD)メモリデバイス等、ノンユニフォームアクセスメモリ技術の使用を可能とする。

10

【0192】

[0149]本発明の様々な実施形態では、メモリデバイスが、メモリデバイスにとって効率的であるIO(入力/出力)パターンに遭遇しながら、インデックス付けアルゴリズムがアルゴリズムにとって自然(効率的)である方法で行うことを保証するために、新しいデータ構造及び方法が提供される。

【0193】

[0150]インデックス付けアルゴリズムによって見なされた論理バケットを、メモリデバイス上の物理バケットにマップする、あるデータ構造である間接指定テーブルが作成される。このマッピングは、ノンユニフォームアクセスメモリデバイスへの書込み性能が強化されるようにするものである。

20

【0194】

[0152]キャッシュの追い出し及びライトバックポリシーの一部として、バケットを集め、メモリデバイスへ順次書き出すために、別のデータ構造である連想キャッシュが使用される。

【0195】

[0153]インデックス付けアルゴリズムによって必要とされる(レコードの)バケットでキャッシュをポピュレートするために、方法が使用される。追加のバケットは、デマンド読出し中に又はスカベンジング処理によって、メモリデバイスからキャッシュへ読み出されてもよい。

30

【0196】

[0154]間接指定テーブルと併せた、キャッシュの使用は、メモリデバイスへの大きい順次書込みを可能にする。

[0155]フラッシュ技術は、インデックス付け問題のために必要とされた容量及びIOレート達成する基本機能を有するが、フラッシュアクセス特性はノンユニフォームである。このノンユニフォーム性はかなり十分なので、通常のインデックス付けアルゴリズムは、フラッシュメモリデバイスと共に、たとえ機能するとしても、うまく機能しない。

【0197】

[0156]本発明で使用されるノンユニフォームアクセスフラッシュメモリは、数百から数千ビットもの大きいブロックサイズで、即ち、バイトレベルのランダムアクセスではなく、読出し、書込み及び消去が行われなければならない、電氣的消去可能プログラマブル読出し専用メモリ(EEPROM)である。物理的には、フラッシュは、フローティングゲートトランジスタからなるメモリセルのアレイに情報を格納する、不揮発性メモリ形式である。2つのタイプのフラッシュメモリデバイス、NANDフラッシュ及びNORフラッシュがある。NANDフラッシュは、NORフラッシュより高密度及び大容量をより低コストで提供し、より高速な消去、順次書込み及び順次読出し速度を有する。本願において、且つ本発明において使用されるとき、「フラッシュ」メモリは、NANDフラッシュメモリを包含するものであり、NORメモリを包含しないものである。NANDは、各セルが1ビットの情報のみを格納するシングルレベルセル(SLC)デバイス、及び、1つのセルにつき2ビット以上を格納することができる、より新しいマルチレベルセル(MLC

40

50

）デバイスを共に含む。NANDフラッシュは高速なアクセス時間を提供するが、PC内でメインメモリとして使用された揮発性DRAMメモリほど高速ではない。フラッシュメモリデバイスは、フラッシュファイルシステムを含んでも含まなくてもよい。フラッシュファイルシステムは、典型的には、ウェアレベリング及び誤り訂正を行うために内蔵コントローラを有していない、内蔵フラッシュメモリと共に使用される。

【0198】

[0157]典型的なNANDフラッシュチップは、数GBのコンテンツを格納してもよい。コンピュータに取り付けられたメモリとは異なり、フラッシュチップ上のメモリは、あるサイズで、且つ、ある境界上でアクセスされなければならない。更に、メモリのセクションに書き込まれた後、それらのメモリ位置が再び書込み可能となる前に、消去オペレーションが行われなければならない。また、位置は徐々になくなるので、全ての位置が同様の数の書込みを得ることを保証することは、使用を更に複雑にする。読出し時間、書込み時間及び消去時間は、大幅に（数マイクロ秒から数ミリ秒まで）変わる可能性がある。よって、タイミング、ウェアレベリング及びアラインメント制約は、フラッシュの実際的な使用を困難にするのがせいぜいである。

10

【0199】

[0158]フラッシュメモリデバイスは、1つ又は複数のダイ（シリコンウエハ）を含んでもよい。各ダイは、大部分では、独立してアクセスされうる。

[0159]ダイは、数千もの消去ブロックからなる。消去ブロックは典型的には、128～512KBのサイズである。データがクリアされる必要があるとき、消去ブロック境界上でクリアされなければならない。

20

【0200】

[0160]NANDフラッシュのもう1つの制限は、データが順次書き込まれることのみが可能であることである。更に、書込みのためのセットアップ時間は長く、読出しのセットアップ時間のおよそ10倍である。

【0201】

[0161]データは、ページ粒度で読み出される。ページは、特定のフラッシュチップに応じて、1KBから4KBまでに及ぶ場合がある。各ページに関連付けられるものは、誤り訂正符号（ECC）チェックサムのために使用されうる、数バイトである。

【0202】

[0162]データは、ページ粒度で書き込まれる。書き込まれた後、ページは、その消去ブロック（ページを含む）が消去されるまで、再度書き込まれない場合がある。消去ブロックは、数十から100を超えるページを含む場合がある。

30

【0203】

[0163]上記の読出し及び書込みページ粒度に対する1つの例外は、サブページ書込み、又は、部分ページプログラミングである。技術に応じて、ページは、消去が必要とされる前に最大4回、部分的に書き込まれる場合がある。

【0204】

[0164]NANDフラッシュブロック内のページは順次、且つ、ブロック消去オペレーションの間に一度のみ書き込まれうるので、後続の書込みは、典型的には異なるフラッシュブロック内に位置する異なるページへの書込みを必要とする。ブロック消去の問題点は、フラッシュファイルシステムの機能である、書込み可能フラッシュブロックのプールを作成することによって、対処される。

40

【0205】

[0165]消去ブロックを消去することは、数ミリ秒を取りうるので、時間的に最も費用のかかるオペレーションである。（トラフィック的に）多用されるデバイスでは、消去ブロックが生成されうる速度（即ち、どのくらい速くフリー消去ブロックが使用可能にされるか）は、フラッシュ設計における制限要因であることが多い。

【0206】

[0166]多数のSSD（ソリッドステートディスク）は、フラッシュ技術を使用する。S

50

SSDにおけるファームウェアは、フラッシュ変換層（FTL）と呼ばれる層において、前述のアクセスの問題点に対処する。そうすることにおいて、しかし、ファームウェアは、どのようにSSDが使用されるか（例えば、大部分は読出し、大部分は書込み、読出し及び書込みのサイズ及びアラインメント）について仮定を行い、これらの仮定の結果として、SSDの性能特性は、インデックス付けアルゴリズムにとって準最適であることが多い。

【0207】

[0167]人が文献において、また実際に発見する多数のインデックス付けアルゴリズムは、ユニフォームメモリアクセスモデルに基づいており、即ち、全てのメモリが、読出しでも書込みでも時間的に等しくアクセス可能であり、アクセスサイズ又はアラインメントにおけるいかなる一次制約もない。

10

【0208】

[0168]インデックス付けソリューションを考察する場合、挿入、削除、ルックアップ及び修正等のオペレーションは典型的には、より多くの多様な時間の量を必要とし、ブロック、典型的には小さいブロック（4KBくらい）の読出し及び書込みは、より短い時間を必要とする。ブロックはランダムであるように見え、即ち、任意のブロックが読み出されてよく、任意の他のブロックが書き込まれてよい。幾つかのアルゴリズムでは、ランダム読出し - 修正 - 書込みI/Oプロファイルがあり、即ち、ランダムブロックが読み出され、次いで、わずかに修正されたデータで同じ位置に書き戻される。

【0209】

20

[0169]インデックス付けアルゴリズムの効率的な動作を必要とするこのランダムI/Oは、フラッシュが提供するように意図されるものではない。フラッシュは、ランダム読出しをうまく扱うことができるが、読出し - 修正 - 書込みがそうであるように、ランダム書込みは困難である。この理由は、すでに書き込まれたものを上書きすることはできず、まずそれを消去しなければならないからである。この状況を更に複雑にすることには、消去は時間を要し、大きい境界（典型的には64KB）上で起こらなければならない。

【0210】

[0170]消去ブロックが消去される時、そのブロック内のいかなる有効データも他のところへ移動される必要がある。アルゴリズムがフラッシュデバイス中でランダムな4KBブロックを書き込む場合、単純な実装は、その間ずっとブロックが消去され続ける結果となる。消去時間は遅いので、性能がかなり損なわれることになる。

30

【0211】

[0171]本発明によれば、フラッシュへの書込みを順次にすることを可能にするため、インデックス付けアルゴリズムが期待する論理ランダムアクセスをなお保ちながら、変換又は間接指定テーブルが作成される。このテーブルは、インデックス付けアルゴリズムによって必要とされるような（レコードの）論理バケットを、フラッシュデバイスの物理バケット（例えば、ページ）にマップする。

【0212】

[0172]インデックス付けアルゴリズムがバケット（例えば、フラッシュからのデータのページ）を読み込むとき、バケットコンテンツを修正（挿入、更新又は削除オペレーション）するために、バケットはキャッシュに移動される。フラッシュデバイス上の対応するバケットにはこのとき、有効ではない（フリー）としてマークが付けられうる。SSDの場合、これはTRIMコマンドの形式を取りうる。

40

【0213】

[0173]本発明の更なる実施形態によれば、フリー消去ブロックを生成するために、方法が提供される。いかなる所与の時間でも、消去ブロックは、有効及び無効データの組み合わせを有する場合がある。消去ブロックを解放するため、全ての有効データは、そのブロックから離れるように移動されなければならない。これを達成するために使用されうる、2つの機構がある。1つは、インデックス付けアルゴリズムによって生成されたランダム読出しを使用して、（インデックス付けアルゴリズムによって必要とされるよりも）より

50

多くを読み出し、消去ブロックを解放するようにすることである。インデックス付けアルゴリズムはランダム読み出しを生成する傾向があるので、経時的に、全ての消去ブロックは最終的に読み出され、空のページのために集積される。例えば、読み出しを含む消去ブロックが、幾つかのフリーページ及び幾つかの有効ページを有する場合、アルゴリズムは、消去ブロック全体を読み込み、且つ、全ての有効ページをキャッシュに入れることを選択してもよい。これは、後続の消去及び次いで書込みのために、その消去ブロックを解放する効果を有する。

【0214】

[0174]別法として、例えば、前述のランダム読み出し処理が十分に高速でない場合、消去ブロックを読み出し、別の消去ブロックと合体するために有効ページをキャッシュに入れるために、分離したスキベンジ処理（例えば、スレッド）が使用されうる。

10

【0215】

[0175]キャッシュが一杯になるとき、エントリが書き出されなければならない。キャッシュエントリのセットが集められ、部分ページの隣接するセット（部分ページ書込みがフラッシュデバイスによって可能にされる場合）、複数のページ及び1つ又は複数の消去ブロックのうちの1以上へと、順次書き込まれるようになる。キャッシュエントリがフラッシュデバイスに書き込まれるとき、間接指定テーブルは更新されるので、インデックス付けアルゴリズムはなお、エントリが固定された論理アドレスにあるものと見なす。

【0216】

4. インデックス付けオペレーション

20

[0176]本発明の様々な実施形態がここで、添付の図9～14を利用して説明され、本発明に従って行われた様々なインデックス付けオペレーションが例示される。図15～16は、ストレージ媒体（例えば、フラッシュメモリ）の効率的利用のためにフリー消去ブロックを生成する、2つの方法を例示する。これらの実施形態は例示的なものであり、限定するものではない。

【0217】

[0177]図9は、本発明の一実施形態による、バケット変換テーブル17'及びキャッシュ23'を利用する幾つかのインデックス付けオペレーションの概観である。図9の上部に、3つのインデックスオペレーション12'～14'が、ルックアップ関数15'及び変換関数16'への代替入力として示される。第1のインデックスオペレーション12'は、キーのための（レコードエントリ）から衛星データを返すための「キーをルックアップする」である。第2のインデックスオペレーション13'は、キーのためのレコードエントリを更新（修正）するための「キーのための衛星データを更新する」である。第3のインデックスオペレーション14'は、新しいレコードエントリを挿入するための「新しいキーを挿入する」である。もう1つのインデックスオペレーションである削除は、図9では図示されないが、図13に関して後述される。

30

【0218】

[0178]全ての3つのインデックスオペレーションは、最初にルックアップ関数15'を行い、キーのある関数 f （キー）が使用されて、インデックスが生成され、ここでは、ハッシュテーブルルックアップをサポート（例えば、加速）する論理バケット識別子である。バケット識別子（インデックス）は、変換関数16'への入力であり、論理バケット識別子のある関数 f （インデックス）は、フラッシュメモリ内の物理バケット位置を生成する。この変換関数はバケット変換テーブル17'によって実装され、バケット変換テーブル17'は、対象フラッシュメモリ位置（フラッシュ内の物理バケット位置）への（インデックス付けアルゴリズムによって提供されるような）論理バケット識別子のマップである。フラッシュメモリ26'に格納されたディクショナリ（インデックス）は、ルックアップキー（例えば、オブジェクト名）を衛星データ（例えば、ディスク上に格納されたオブジェクトへの位置ポインタ）にマップするレコードを備えてもよい。

40

【0219】

[0179]次に、3つのインデックス付けオペレーションのうちどれが行われているか（ル

50

ックアップ、更新又は挿入)に応じて、図9の下半分に示されたステップのうち1つ又は複数が行われる。

【0220】

[0180]ルックアップオペレーション18'では、変換関数によって識別されたバケットエントリは、キャッシュルックアサイドによる、フラッシュメモリ内の対象バケット22'からの読出し30'である(例えば、対象バケットがキャッシュに格納される場合、フラッシュメモリ26'からではなく、キャッシュ23'から読み出されてもよい)。

【0221】

[0181]更新オペレーション19'では、変換関数によって識別されたバケットエントリ(元のバケットエントリ)が、フラッシュメモリ(又はキャッシュ)の消去ブロック21a'内の対象バケット22'から読み出され30'、バケットが更新され、キャッシュへ移動され32'、且つ、後続の書込み24'では、複数のキャッシュバケットエントリが、フラッシュメモリ内の部分ページの隣接するセット、複数のページ、及び、消去ブロック(例えば、新しい消去ブロック21b')のうちの1以上に順次読み出される。この処理は、フラッシュ内の全ての移動されたバケットのステータスを、有効ではないデータ(例えば、フリー、又は、トリムオペレーションのために使用可能)に更新する33'。

10

【0222】

[0182]挿入オペレーション20'では、対象バケットがフラッシュから再度読み出され、修正されたバケットエントリが、フラッシュメモリ内の新しい位置への後続の順次書込み24'のために再度、キャッシュに移動される34'。

20

【0223】

[0183]図9は、隣接するフラッシュメモリバケットへの、キャッシュバケットエントリの集まりの順次書込み24'を行うことに先立って、複数のバケットエントリを集めるためのキャッシュ23'を概略的に示す。一実施形態では、スキベンジングオペレーション25'が、フリー消去ブロックを作成するために使用され、この処理は、スキベンジング処理中に(消去ブロックからの)任意の有効バケットをキャッシュに格納すること、及び、フラッシュ消去ブロックをフリーとして再割り振りすることを含む。

【0224】

[0184]図10に例示された新しいデータ構造の考察に続いて、図9で参照されたインデックス付けオペレーションが、図11~14のフロー図に関してより具体的に説明される。

30

【0225】

5. データ構造

[0186]図10は、本発明において有用なデータ構造の様々な実施形態を例示する。そのようなデータ構造は例示的なものであり、限定するものではない。

【0226】

[0187]図10Aは、(インデックス付けアルゴリズムによって生成された)論理バケットインデックスを、物理フラッシュバケットアドレスに変換するための、バケット変換テーブル(BTT)300'の一実施形態を例示する。BTTテーブルエントリは、3つのフィールド、即ち、有効301'、フラッシュ物理バケットアドレス302'、及び、拡張バケット状態303'を有して図示される。バケットアドレス粒度は、フラッシュデバイスの最小書込みサイズであり、即ち、(例えば、SLC NANDでの)部分ページ書込み、又は、(例えば、MLC NANDでの)ページ書込みである。BTTは、論理から物理バケットエントリへの1:1マッピングである。このテーブルは、より高いランダム性能(インデックス付けアルゴリズムによるランダム読出し及びランダム書込み)のために、フラッシュバケット割当ての再編成を可能とする。アルゴリズムのアクセラレーションを可能とするために、追加の状態情報がBTTへ、第3のフィールド内に追加されてもよい。

40

【0227】

[0188]図10Bは、バケット有効テーブル(BVT)305'の一実施形態を示す。こ

50

のテーブルは、トリミングのためのブロックへのバケットのスキベンジングを管理するために、フラッシュ内のどの物理バケットが有効であるかを追跡する。一例として、有効とラベル付けされたフィールド 306' は、小型のビット配列 (1ビット/バケット) であってもよい。BVTのサイズは、フラッシュバケットエントリの総数であり、そのサブセットのみが BTTによって使用中である。

【0228】

[0189] 図 10C は、バケット内に含まれた複数のレコード 310'、311'、312'、... を、逆 BTTポインタ 313' (バケット変換テーブル 17' へのセルフインデックス) と共に有する、フラッシュバケット 309' の一実施形態を例示する。よって、各バケットは、1つ又は複数のレコードのセット、及び、フラッシュバケット (例えば、ページ) が挿入、移動又は削除されるとき、BTTを更新するための逆ポインタを含む。バケットの各要素 (レコード又はポインタ) には、データ構造の個々の信頼性を向上させ、且つ、ストレージデバイスの有用寿命を著しく延ばすために、追加の ECCビット等、冗長コンテンツが追加されてもよい。例えば、任意選択のシーケンス番号フィールドが、停電発生中にデータ一貫性チェックを行うためにフラッシュバケット 309' に追加されてもよく、他の最適化フラグも同様に提供されてもよい。

10

【0229】

[0190] レコードサイズはバケットサイズに比べて小さいので、これは、個々のレコードベースで追加の誤り回復情報を実装する (任意選択の) 機会を提供する。この任意選択の機能は、訂正される可能性があるの下にあるストレージ技術の実効動作寿命を延ばす可能性のある、ビット誤り及び障害の数を増すことによって、ソリューションの全体の信頼性を向上させるようになる。

20

【0230】

[0191] 図 10D は、複数の消去ブロック 316' (1からM) を含む、SLC NANDフラッシュデバイス 315' の一例を示す。各消去ブロックは、複数のページ 317' (1からN) を含む。この例では、各ページは 4KBであり、各ページは複数のバケット 318' (1からB) を含み、各バケットは 1KBである。この例では、デバイスは部分ページ書込みをサポートする。

【0231】

[0192] バケットは、フラッシュデバイスの最小書込みサイズを表現する。典型的には、バケットは 1ページとなる。部分ページ書込みが可能にされる場合、1ページにつき4つのバケットをサポートする4部分ページ SLC NANDデバイス等、1フラッシュページにつき1つ又は複数のバケットが提供されてもよい。

30

【0232】

[0193] 1つの消去ブロックにつき、複数のフラッシュページが提供される。1つのフラッシュデバイスにつき複数の消去ブロックがあり、各ブロックは個々に消去される。

[0194] 典型的なフラッシュサブシステムは、複数のフラッシュデバイスからなる。NANDフラッシュデバイスは、消去オペレーション間に所与のブロック内で1ページ (又は部分ページ) につき1回、順次書き込まれ、複数のブロックが、同時に書込み及び読出しを行うために使用可能である。

40

【0233】

6. 処理フローチャート

[0195] 図 11 は、キーの存在を検証し、関連付けられた衛星データを返すためのルックアップオペレーション処理の一実施形態を例示する。ステップ 141' で、ルックアップキーがルックアップ関数に入力される。ステップ 242' で、ルックアップ関数 f (キー) は、ハッシュテーブルルックアップをサポート (例えば、加速) する論理バケット識別子を生成する。論理バケット識別子は変換関数に入力され、変換関数は、ステップ 343' で、バケット変換テーブル (BTT) 17' を介してフラッシュメモリ (物理バケット) 位置にマップされる。ステップ 444' で、バケットがキャッシュに格納されるまで、フラッシュメモリ内の対象バケットがフラッシュメモリから読み出され 45a'

50

、バケットがキャッシュに格納された場合、キャッシュ 23' から読出し可能である 45 b'。ステップ 6 46' で、キーのための衛星（レコード）データがインデックス付けアルゴリズムに返される。

【0234】

[0196] 図 12 は、挿入オペレーション処理の一実施形態を示す。最初のステップ 71' は、キーをルックアップ関数に入力する。ステップ 2 72' で、ルックアップ関数 f (キー) は、インデックス、ここでは論理バケット識別子を生成する。ステップ 3 73' で、バケット識別子が変換関数に入力され、変換関数は、バケット変換テーブル (BTT) 17' を利用して、バケット識別子を、挿入が起こるべきフラッシュメモリ物理バケット位置へマップする。ステップ 4 74' で、挿入処理は、対象バケット位置を変換関数から受信する。ステップ 5 75' で、挿入処理は、対象バケット 22' をフラッシュメモリの消去ブロック 21a' から 75a'、或いはキャッシュから 75b'、読み出す。ステップ 6 76' で、挿入処理は、レコードエントリを対象バケットに挿入し、修正済みバケットをキャッシュに書き込む。ステップ 7 77' で、挿入処理によって、複数のバケットエントリ（修正済み対象バケットを含む）がキャッシュ 73' から読み出される。ステップ 8 78' で、挿入処理は、キャッシュから読み出された修正済み対象バケット及び他のバケットを、フラッシュ 26' 内の新しい位置（消去ブロック 21b' 内のページ）へ書き込む。ステップ 9 79' で、挿入処理は、キャッシュからフラッシュへ移動された全てのバケットのための新しい位置により、バケット変換テーブル 17' を更新し 79a'、また、移動された全てのバケットについて、BVT 内のバケット有効エントリをも更新する 79b'。ステップ 10 80' で、挿入処理は、移動されたキャッシュエントリに、読出し専用（使用可能）とマークを付ける。ステップ 11 81' で、挿入処理は、元のフラッシュバケット（このとき、新しい消去ブロックに移動済み）に、フリーとしてマークを付ける。

【0235】

[0197] 図 13 は、削除オペレーション処理の一実施形態を例示する。最初のステップ 91' では、キーがルックアップ関数に提供される。ステップ 2 92' で、ルックアップ関数 f (キー) は、インデックス、ここでは論理バケット識別子を生成する。ステップ 3 93' で、バケット識別子が変換関数に提供され、変換関数は、バケット変換テーブル 17' を利用して、バケット識別子を物理フラッシュメモリバケット位置へマップする。ステップ 4 94' で、削除処理は、フラッシュメモリ位置を受信する。ステップ 5 95' で、対象バケットがフラッシュから読み出される。ステップ 6 96' で、処理は、バケット内の元のレコードエントリを削除し、修正済みバケットを（削除済みエントリと共に）キャッシュ 23' に書き込む。ステップ 7 97' で、バケットのグループ（集まり）がキャッシュから読み出される。ステップ 8 98' で、キャッシュ 23' から読み出された更新済み対象バケット及び他のバケットが、フラッシュ内のフリーページの隣接するセットに順次書き込まれる。ステップ 9 99' で、削除処理は、全ての移動されたバケットのためのフラッシュ内の新しい位置により、バケット変換テーブルを更新し 99a'、BVT 内のそれらのバケット有効ステータスを更新する 99b'。ステップ 10 100' で、削除処理は、キャッシュエントリに、読出し専用としてマークを付ける。ステップ 11 101' で、削除処理は、このとき、フラッシュ内の新しい位置に移動された元のフラッシュバケットに、フリーとしてマークを付ける。

【0236】

[0198] 図 14 は、フラッシュメモリに格納されたインデックス内のレコードを修正するための更新オペレーション処理の一実施形態を例示する。最初のステップ 51' で、キーが入力としてルックアップ関数に提供される。ステップ 2 52' で、ルックアップ関数 f (キー) は、インデックス、ここでは論理バケット識別子を生成する。バケット識別子は変換関数に入力される。ステップ 3 53' で、変換関数は、バケット変換テーブル 17' を利用して、バケット識別子を、更新が起こるべきフラッシュメモリ内の物理バケットへマップする。ステップ 5 55' で、対象バケットはフラッシュから 55a'、或い

10

20

30

40

50

はキャッシュから55b'、読み出される。ステップ656'で、エントリを更新した後、更新済みバケットがキャッシュ23'に書き込まれる。ステップ757'で、バケットのグループがキャッシュ23'から読み出され、ステップ858'で、キャッシュからフラッシュメモリ26'内の新しい位置へ順次書き込まれる。ステップ959'で、更新処理は、全ての移動されたバケットのための新しい位置により、バケット変換テーブル17'を更新し59a'、BVT内のそれらの有効ステータスを更新する59b'。ステップ1060'で、更新処理は、移動されたエントリに、キャッシュ23'内で読出し専用(及び、よって上書きされるために使用可能)としてマークを付ける。最後に、ステップ1161'で、更新処理は、このとき新しい位置に移動された元のフラッシュバケットに、フリー(使用可能)としてマークを付ける。

10

【0237】

[0199]図15Aは、フリー消去ブロックを生成するための処理の一実施形態を例示し、(ルックアップ、挿入又は修正等、アップストリームインデックス付けオペレーションによって生成された)デマンド読出しは、(対象バケットと)同じ消去ブロック内の追加のバケットを読み出す。図15Aで、この処理は、更新要求と共に例示される。ステップ1111'で、キーがルックアップ関数に提供される。ステップ2112'で、ルックアップ関数 $f(\text{キー})$ は、インデックス、ここでは論理バケット識別子を生成する。ステップ3113'で、バケット識別子が、フラッシュ内の物理対象バケット位置にマップされる。ステップ4114'で、更新及びスキベンジ処理は、対象フラッシュメモリ位置を受信する。ステップ5115'で、処理は、対象バケットと同じ消去ブロック内の全ての有効バケットを識別する。ステップ6116'で、更新処理は、対象バケット及び全ての識別された有効バケットを、対象バケットを含むフラッシュブロックから読み出す。ステップ7117'で、処理は、対象バケット内のレコードエントリを更新し、フラッシュブロックからの全ての有効バケットをキャッシュ23'に書き込む。ステップ8118'で、更新処理は、キャッシュからブロックのグループを読み出す。ステップ9119'で、更新処理は、キャッシュ23'から読み出された更新済み対象バケット及び他のバケットを、フラッシュ26'に書き込む。ステップ10120'で、更新処理は、移動された(キャッシュからフラッシュ内の新しい消去ブロック21b'に書き込まれた)全てのバケットのための新しい位置により、バケット変換テーブル17'を更新し120a'、BVT内のバケットエントリを更新する120b'。ステップ11121'で、更新処理は、このとき古いキャッシュエントリに、読出し専用としてマークを付ける。ステップ12122'で、更新処理は、元のフラッシュブロック(対象ブロック内の全てのバケット)に、フリーとしてマークを付ける。

20

30

【0238】

[0200]図15Bは、フリー消去ブロックを生成するためにのみ記載された、ランダム読出し処理の特定の実施形態を例示する。

[0201]この実施形態では、置き換えハッシングインデックス付けアルゴリズム125'は、論理バケット126'を生成する。インデックス付けアルゴリズムによって見なされる論理バケットサイズは、インデックス付けアルゴリズム及びフラッシュメモリが互換性を有するように、フラッシュ消去ブロックサイズに関連付けられる。これらのバケットは、インデックス読出し及び更新の結果として、ランダムに読み出されるようになる。

40

【0239】

[0202]バケット変換(間接指定)テーブル127'は、論理バケットインデックスを物理フラッシュデバイスバケット位置に変換する。この間接指定テーブルは、インデックス付けアルゴリズムが、読出し、書込み及び更新のために、ランダムに動作し、且つ、フラッシュデバイスレベルで行われる大きい順次書込みをなお有することを可能とする。好ましくは、間接指定テーブルは永続メモリに格納されるが、揮発性メモリに格納される場合、必要に応じて再構築されうる。

【0240】

[0203]間接指定テーブルの出力、即ち、物理デバイスバケット位置は、入力として、完

50

全連想バケットキャッシュ 128' に提供される。この実施形態では、空の消去ブロック f i f o 129' のコンテンツがハイウォーターマーク Q 未満である場合、消去ブロック全体（対象 4KB バケットを含む）が読み出される。

【0241】

[0204] 消去ブロックは論理バケットをホストし、典型的な構成は、1つの消去ブロックが 16個の 4KB 論理バケットを保持するものである。物理デバイスは、負荷、例えば、90%のために構成され、バケットの 90%が使用中であることを意味する。キャッシング及び犠牲化 (v i c t i m i z a t i o n) (追い出し) は、残りのバケットの 10%の大部分がフリー消去ブロックに集中させられるように、フラッシュメモリ内で論理バケットをパックする(集中させる)ために、使用される。

10

【0242】

[0205] キャッシュ犠牲化(追い出し処理)は、キャッシュ内に集められた 16個のバケットを取り、それら 16個のバケットをキャッシュからフリー消去ブロック 130' に書き出す。消去ブロックはランダム読出しオペレーションによってランダムに触れるるので、読出しオペレーションは、フリー消去ブロックを生成するために使用されうる。論理バケット識別子を生成するための暗号ハッシュ関数の使用は、読出しオペレーションのランダム性を増大させ、よって、フリー消去ブロックのランダム読出し生成を改善するようになる。

【0243】

[0206] 図 16A 及び 16B は、フリー消去ブロックを生成するための代替スキベンジ処理を例示する。このスキベンジ処理は、何れかのインデックス付けオペレーションの一部ではない。むしろ、より下位レベルのデバイス管理層の一部として実装される。この処理では、フラッシュ消去ブロック内の物理バケットのグループ(一部又は全部)がフラッシュから直接読み出され、バケット有効テーブル 27' が、消去ブロック内のどのバケットが有効であるかを決定するために使用される。

20

【0244】

[0207] 図 16A に例示されるように、ステップ 1 220' で、スキベンジ処理 25' は、完全な消去ブロック 21a' を読み出す。ステップ 2 222' で、スキベンジ処理は、バケット有効テーブル 27' を使用して、有効であるそれらの読出しの全てのバケットを識別する。ステップ 3 224' で、各有効バケットについて、論理バケット識別子がバケットから抽出される。ステップ 4 226' で、有効バケットがキャッシュ 23' に格納され、各々はその論理バケット識別子によってインデックス付けされる。

30

【0245】

[0208] 図 16B は、ステップ 1 で、スキベンジ処理 25' がバケット [94, 97] を、94 及び 97 を含めて読み出す一例を示す。ステップ 2 で、処理は、95 及び 96 でのバケットが有効であると決定する。有効バケットは、「1」によって指定されたバケット有効テーブルにおいて示され、非有効バケットは「0」による。ステップ 3 で、バケット 95 及び 96 のための論理バケット識別子、即ち、タグ 23 及び 49 がそれぞれバケットから抽出される。ステップ 4 で、この 2つのタグ並びにそれらの各バケット 95 及び 96 は、それらの各タグ 23、49 をインデックスとして使用して、キャッシュに挿入される。

40

【0246】

7. スタックレベルビュー及び実装

[0209] 本発明のもう 1つのより具体的な実施例が、ここで図 17 ~ 24 を参照して説明される。

【0247】

[0210] 図 17 は、物理フラッシュメモリデバイス 211' によって望まれるものとは大異なるビューである、インデックス付けアルゴリズム 203 によって望まれた I/O 使用プロファイルビューを、フラッシュアダプテーション層 207' が適合させる、本発明の

50

一実装を例示するための6層のビュー又はスタック200'を示す。最高レベル201'で、レコードのディクショナリ(インデックス)が提供され、それに対して、あるインデックス付けオペレーション204'(レコードをルックアップ、削除、挿入且つ修正する)が必要とされる。インデックス付けアルゴリズム層203'は、1つ又は複数のインデックス付けアルゴリズムによりディクショナリを実装し、例えば、カックウ置き換えハッシュアルゴリズムが一例である。インデックス付けアルゴリズムは、インデックスへのキーがどのようにインデックス永続層205'によって格納されるかのビューを有する。インデックス付けビューは、論理ビューであり、論理アドレス位置を特定する。このビューは更に、サイズ、アラインメント及びタイミングについて、インデックスへのユニフォームアクセスが存在するようになること、及び、インデックスが可変(安定)ストレージ上に格納されることを仮定する。

10

【0248】

[0211]インデックス永続層205'は、インデックスのレコードを格納する物理バケットを読み出し、それらの物理バケットに書き込むための論理バケットオペレーション206'を提示するようになる。これらの論理バケットオペレーション206'は、フラッシュアダプテーション層207'に提示され、フラッシュアダプテーション層207'は、前述のように、(インデックス付け処理の)論理バケットをフラッシュストレージデバイス上の物理バケット位置に変換する。フラッシュアダプテーション層は、よって、上記のインデックス付けアルゴリズムによって望まれたビュー及びI/O使用プロファイルを、下記の物理ストレージデバイス(フラッシュメモリ211')によって望まれた大変異なるビューに適合させる。ここで、物理バケットオペレーション208'は、ランダム読み出し及び集約(ブロック順次)書き込みを含み、これらのランダム読み出し及び集約(ブロック順次)書き込みは、バケットアクセスのノンユニフォームモデルを構成する。この実施例の物理バケットオペレーションは、トリムコマンドを更に含んでもよい。

20

【0249】

[0212]物理バケットオペレーションは、物理フラッシュデバイス上の資源を追跡且つ調節する、デバイス管理層209'によって実装される。これらの物理デバイスオペレーション210'はここでは、ランダム読み出し、大きい順次書き込み及びトリムコマンドを含む。

【0250】

[0213]物理デバイス層211'は、そのノンユニフォーム読み出し及び書き込み、並びに、サイズ、アラインメント及びタイミングについての不変性によって特徴付けられる。そのような物理デバイスの例は、生のフラッシュ、相変化、SSD、及びデバイス上に存在するフラッシュファイルシステムを有するフラッシュのうちの1以上を含む。

30

【0251】

[0214]本発明は、以下のもの等、デバイス管理層より下の追加の任意選択の強化を可能とする。

- ・バケットトリミング(細かいページトリミング)及びページ内のバケット追跡のモデルは、SSD又は同等のストレージデバイスのフラッシュファイルシステムに直接組み込まれる場合、よりよい消去ブロック管理を可能とする。

40

【0252】

- ・フラッシュページ上へのバケットのマッピングは、抽象化である。変化ごとにフラッシュに書き込まれたデータの量を最小化することによって、それらのデバイスの寿命を延ばすために、バケットは、SLC NANDのための部分ページへマップすることができる。バケットはまた、これがシステム全体の性能にとって有益であった場合、複数のフラッシュページ上にマップすることもできる。

【0253】

[0215]図18は、インデックスレコードの一例を示す。レコード140'は、フィンガープリント(キー)を格納するための最初の20バイトのフィールド141'を含め、合計で32バイトである。フィンガープリントは、好ましくはデータコンテンツの暗号ハッシュ

50

シュダイジェスト、例えば、SHA-1ハッシュアルゴリズムである。例示を容易にするため、フィンガープリントを、「A B 9 2 3 4 5 E 2 0 3 . . .」等、16進数で入力するのではなく、個々のフィンガープリントは、図19~22において、P、Q、R、S、T等、単一の大文字によって指定される。やはり例示のために簡潔にするため、これらの大文字はまた、レコード全体の代理としての役割も果たすようになる。レコードのフィールドはまた、2バイトの参照カウントフィールド142'、5バイトの物理ブロックアドレスフィールド143'、1バイトのフラグフィールド144'、及び、4バイトの雑フィールド145'をも含む。PBAフィールド143'は、指定されたフィンガープリント141'のために、ディスク上で格納されたデータの物理ブロックアドレスへのポインタを含む。参照カウントは、ディスク上で格納されたデータへの参照の数を追跡する。

10

【0254】

[0216]本発明の一実施形態によれば、インデックスレコードからのフィンガープリント141'は、前述(図9)のルックアップ関数 $f(\text{キー})$ への入力キーとして使用され、この実施例では、関数 $f(\text{キー})$ は、4つのハッシュ関数 H_0 、 H_1 、 H_2 及び H_3 のセットを備える。一般に、2つ以上のハッシュ関数の任意のセットを使用することができる。ハッシュ関数 H_x は、フィンガープリントを、0及び $N-1$ を含めて範囲 $[0, N-1]$ へマップし、但し、 N はハッシュテーブルのサイズである。この実施例では、フィンガープリント自体がハッシュであることを考えれば、BitFieldを抽出して、以下の4つのハッシュ値のファミリーを生成することができる。

【0255】

$$\begin{aligned} H_0(x) &= x < 0 : 31 > \text{mod } N \\ H_1(x) &= x < 0 3 2 : 63 > \text{mod } N \\ H_2(x) &= < 0 6 4 : 95 > \text{mod } N \\ H_3(x) &= x < 0 9 6 : 127 > \text{mod } N \end{aligned}$$

20

[0217]抽出されたBitField幅は、 $\log_2(N)$ 以上である。ばらばらのビットの任意の組み合わせを使用することができ、 $\log_2(N)$ 制約を受ける。図18に例示されるように、最初のフィールド141'内のフィンガープリントのみがハッシュされて、キーが形成される。レコード140'の残りのコンテンツ(フィールド142'~145')は、値又はペイロードを備える。

【0256】

[0218]図19は、カッコウハッシングとして知られる、置き換えハッシングインデックス付けアルゴリズムの一例を例示する。例示を容易にするため、2つの関数のみが使用される。図19Aは、フィンガープリントPがそれぞれ関数 $H_0(x)$ 及び $H_1(x)$ からハッシュ値2及び5を生成する一方、フィンガープリントQがこれらの同じ関数からハッシュ値1及び3を生成する、 2×3 グリッドを示す。カッコウハッシングアルゴリズムは、0~6のラベルが付けられた7個のロット(図19B)のうち1つにP及びQを入れるための、2つの代替ハッシュ値の中から選択するようになる。Pは、2つの位置のうち1つ、2又は5に入る可能性があり、Qは、2つの位置のうち1つ、1又は3に入る可能性がある。図19Cに示すように、このアルゴリズムは、Qを最下位の空ロット1に、Pをロット2に入れる。この実施例では、レコードコンテナは、1つのレコードを保持するロットとして見なされるが、本発明はそのように限定されず、インデックス付けアルゴリズムはまた、複数のレコードを保持するバケットをもコンテナと見なすことを理解されたい。ここでは、説明を簡素化するため、単一のレコードロットが使用される。

30

40

【0257】

[0219]このとき、もう1つのフィンガープリントRが提供され、Rは、同じハッシュ関数から1及び2のハッシュ値を生成する(図19Dのテーブルを参照)。ハッシングアルゴリズムは、Rを左の位置、即ちロット1に入れ、現在のエントリQを変位させるようになる(図19E)。Qはこのとき、 $H_1(Q)$ によって指定された他の任意選択の位置、即ち位置3に移動されるようになる。このアルゴリズムは、各レコードが空のロットに入るまで、レコードを変位させ続けるようになる。

50

【 0 2 5 8 】

[0219]この実施例では、「Rを挿入する」オペレーションを達成するため、インデックス付けアルゴリズムは、以下の読出し及び書込み要求を生成する。

- 1を読み出す(Qを得る)
- 2を読み出す(Pを得る)
- 1を書き込む(Rを書き込む)
- 3を読み出す(妥当性チェック)
- 3を書き込む(Q)

[0220]最初の2つの読出しは、Rがインデックス内にすでに存在しているのではないことを検証するために使用される。妥当性チェック(3を読み出す)は、スロット番号3が空であるかどうかを決定し、そうである場合、Qがスロット3に書込み可能となり、どのエントリもスロット3に再書込みされなかったとき、このアルゴリズムは終了される。スロット3が空でなかった場合、スロット3内の現在のエントリが別のスロットに移動されることが必要となる。我々がビットマップを有する場合、スロット3のコンテンツは既知であり、そうでない場合、我々は、そのステータスを決定するために、スロット3内のエントリを読み出す必要がある。各エントリは、そのエントリが有効であるかどうかを示す有効ビットを含む。有効とは、それが使用中である(且つ、その位置を現在占有しているものが変位されなければならない)ことを意味する。有効ではないとは、その位置が空であることを意味し、処理中のレコードがそこに書込み可能である。有効ビットのコンテンツもまた、あるメモリを犠牲にして、別のビットマップに格納されうる。

【 0 2 5 9 】

[0221]カックウハッシングアルゴリズムは、空のエントリに着地するまで、複数のエントリに渡って書き込み、以前のコンテンツを変位させ続けるという点において、再帰的である。実際には、この処理はめったに1回の変位を越えることはない。

【 0 2 6 0 】

[0222]インデックス付けアルゴリズムは、バケット及び個々のレコードオペレーションを共に有する。インデックス付けアルゴリズムは、1つのレコードを1つのコンテナ(スロット)に入れることとして(図19で)上述されるが、インデックス付けアルゴリズムによって、レコードがバケット、即ち、複数のレコードを含むバケットに集約されてもよいことは理解される。よって、上記の実施例は非限定的であり、概してレコードオペレーションを例示するものである。

【 0 2 6 1 】

[0223]前述のように、個々のレコードの読出し及び書込みは、フラッシュメモリにとって効率的ではないので、個々のレコードはバケットに集約される。図20は、各々が2つ以上のレコードを含む、4つのそのようなバケットを例示し、即ち、レコード位置0及び1を有するバケット B_0 、レコード位置2及び3を有する B_1 、レコード位置4及び5を有する B_2 、並びに、レコード位置6及びxを有する B_3 である。バケットサイズはフラッシュデバイスによって決定された最小書込みサイズ、即ち、完全ページ書込み又は部分ページ書込みの関数である(且つ、好ましくはそれに等しい)。典型的なバケットサイズは、4KBであってもよい。レコードの具体的な順序付けは、バケット内で必要とされない。ルックアップオペレーション中にバケット全体が有効レコードについて検索されるので、レコードはバケット内のいかなる点でも挿入可能である。変位時、カックウハッシングアルゴリズムによれば、バケット内のエントリはランダムに変位されうる。インデックス付けアルゴリズムは、よって、論理バケットを、ランダムな位置のように見えるものに、1つずつ書き込み、最終的にフラッシュアダプテーション層によって、フラッシュデバイスへのより大きい、物理的に隣接する(順次)書込みに集約される。

【 0 2 6 2 】

[0224]図21は、バケットエントリ160'の一例を例示する。4KBのバケットサイズは、下にあるデバイスの最小書込みサイズ、ここでは4KBのページに基づく。4KBのバケットは、バケットエントリ内のレコードの数を指定する、4バイトの最初のフィー

ルド161'を含む。4バイトのタグフィールド162'は、論理バケット識別子を指定する。この識別子(タグ)は論理アドレスであり、物理的なものではない。変換テーブルは、アルゴリズムバケットアドレス(ABA)をフラッシュバケットアドレスFBAにマップする。キャッシュは、(CPUの用語で)仮想キャッシュとして動作し、各キャッシュライン(エントリ)は、タグ、この場合はABAによって識別される。アルゴリズムがレコードを要求するとき、キャッシュを通過中にそれが知ることは、要求されたABAがキャッシュされることのみであり、ABAがマップされる場所(FBA)は、キャッシュの下端(例えば、図10Cにおける、BT Tへの逆ポインタ313'を参照)である。バケットは、複数のレコードR₀、R₁、R₂...を保持するためのフィールド163'を含み、各レコードは、32バイトのサイズである。この例では、4KBのバケットは、1バケットにつき(4096-4-4)/32個のレコード、即ち、約127個のレコードを保持するようになる。

10

【0263】

[0225]図22は、一実施形態における、バケット、ページ及び消去ブロックの相対サイズを例示する、フラッシュメモリデバイス164'の概略図である。物理フラッシュデバイスは、2GBのサイズであるチップ(パッケージ)165'である。チップ上に、2つのダイ(シリコンウエハ)166a'、167b'がある。各ダイ上に、2¹⁴個の消去ブロックがあってもよく、各消去ブロック167'は典型的には64KBである。ページ168'は、書込み可能である最小サイズ、ここでは4KBであり、スタック(図17を参照)内のより上位で使用されるような、同じく4KBである、バケット169'のサイズを決定する。

20

【0264】

[0226]図23は、物理フラッシュデバイス上の資源を追跡且つ調節するためのデバイス管理層(図17の209')の一実施形態による、選ばれたコンポーネントを例示する。図23Aは、(上部に)複数のページ(バケット)170'と、その後続く、どのページが有効であるか(1は有効、0は有効でない)を示す、ページアロケーションマップ171'を示す。この下は、将来トリムされることになるが、まだそうされていないページの、保留中トリムマップ172'である。ページアロケーション及び保留中トリムマップは、前述のような本発明の様々な実施形態では、バケットが有効データを保持するかどうかを決定する(図9に例示されたバケット有効テーブル27'を参照)するために使用されうる。

30

【0265】

[0227]図23Bは、消去ブロックインデックスによってインデックス付けされた、消去ブロック記述子テーブル175'の一例を例示する。各消去ブロック記述子エントリ176'は、消去済みの数177'、部分書込みの数178'、部分読出しの数179'、完全読出しの数180'、完全書込みの数181、及び、誤りの数182'を含む、複数のフィールドを含む。この情報は、本発明の様々な実施形態において前述されたように、フリー消去ブロックの生成において使用されうる。

【 図 5 】

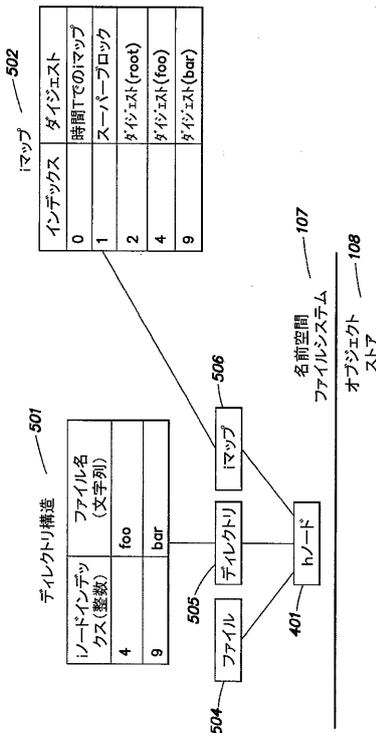


FIG. 5

【 図 6 】

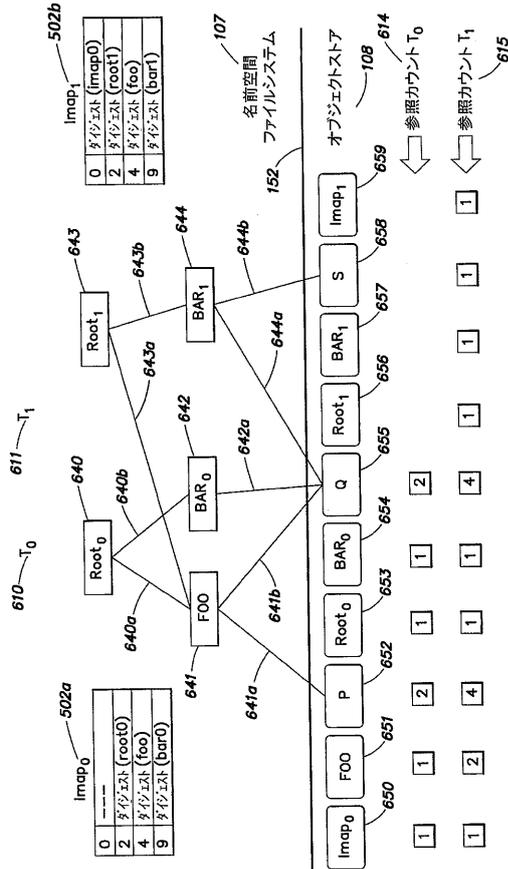


FIG. 6

【 図 7 】

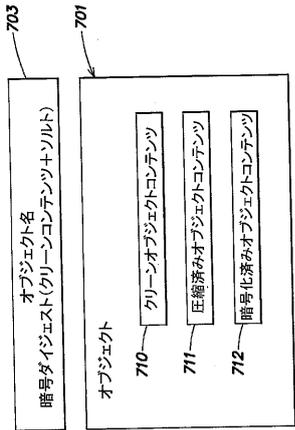


FIG. 7

格納されたオブジェクトは、
 クリア、
 圧縮済み、
 圧縮済み+暗号化済み
 暗号化済み、
 のうちの1つである

【 図 8 】

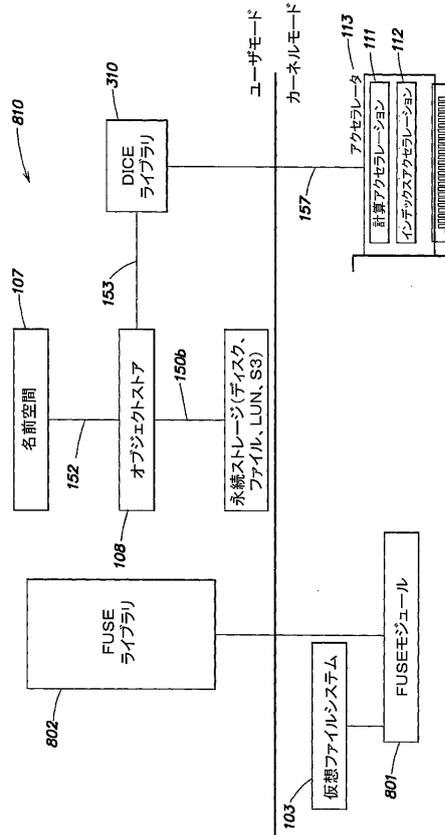


FIG. 8

【図18】

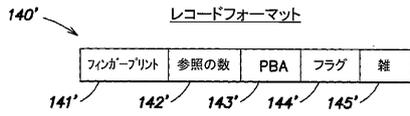


FIG. 18

【図19】

置き換えハッシング(カックウハッシング)

	H ₀ (X)	H ₁ (X)
P	2	5
Q	1	3

FIG. 19A

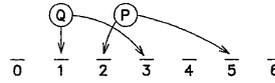


FIG. 19B

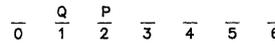


FIG. 19C

	H ₀ (R)	H ₁ (R)
R	1	2

FIG. 19D

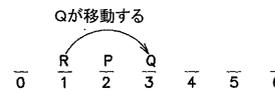


FIG. 19E

【図20】

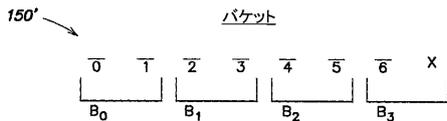


FIG. 20

【図22】

【図21】

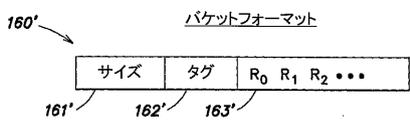


FIG. 21

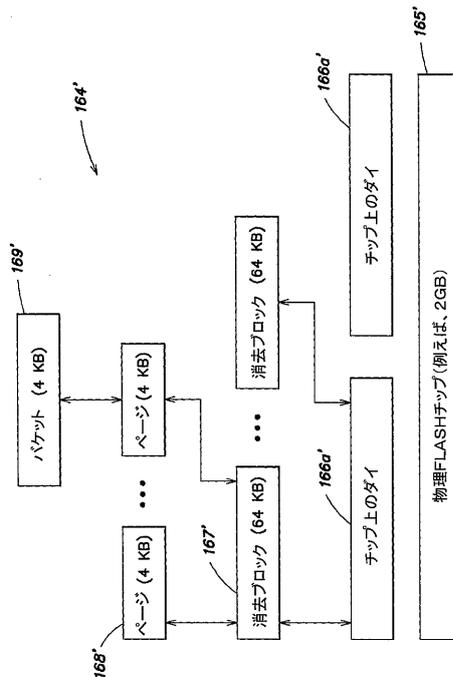


FIG. 22

【図 23】

デバイス管理

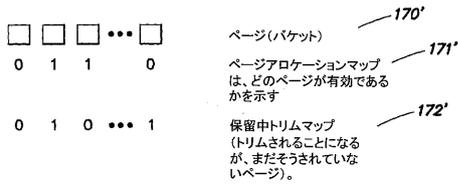


FIG. 23A

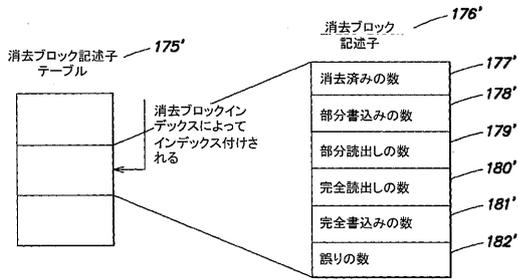


FIG. 23B

フロントページの続き

- (72)発明者 ビーバーソン,アーサー・ジェイ
アメリカ合衆国マサチューセッツ州01719,ボックスボロー,コブリー・ロード 37
- (72)発明者 ボウデン,ポール
アメリカ合衆国マサチューセッツ州01503,バーリン,ライマン・ロード 113

審査官 池田 聡史

- (56)参考文献 特開2007-305122(JP,A)
特開平06-195253(JP,A)
特開2005-018758(JP,A)
特開2004-118824(JP,A)
特開2004-157958(JP,A)
特開2008-257725(JP,A)
特開平10-083299(JP,A)
特開2008-251010(JP,A)
今泉貴史,UNIX流プログラミング:51,UNIX MAGAZINE,日本,株式会社アスキー,1995年 1月 1日,第10巻 第1号,pp.166~173
- (58)調査した分野(Int.Cl.,DB名)
G06F 12/00
JSTPlus(JDreamIII)