

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5302886号
(P5302886)

(45) 発行日 平成25年10月2日(2013.10.2)

(24) 登録日 平成25年6月28日(2013.6.28)

(51) Int.Cl. F I
G 0 6 F 12/00 (2006.01) G O 6 F 12/00 5 2 O P
G 0 6 F 3/06 (2006.01) G O 6 F 3/06 3 O 1 G
 G O 6 F 12/00 5 2 O E

請求項の数 10 (全 26 頁)

(21) 出願番号	特願2009-518192 (P2009-518192)	(73) 特許権者	303039534
(86) (22) 出願日	平成19年6月25日 (2007.6.25)		ネットアップ、インコーポレイテッド
(65) 公表番号	特表2009-543198 (P2009-543198A)		アメリカ合衆国 カリフォルニア 940
(43) 公表日	平成21年12月3日 (2009.12.3)		89, サニーヴェール, イースト ジ
(86) 国際出願番号	PCT/US2007/014664		ャバ ドライブ 495
(87) 国際公開番号	W02008/005211	(74) 代理人	100087642
(87) 国際公開日	平成20年1月10日 (2008.1.10)		弁理士 古谷 聡
審査請求日	平成21年6月16日 (2009.6.16)	(74) 代理人	100076680
(31) 優先権主張番号	11/477,886		弁理士 溝部 孝彦
(32) 優先日	平成18年6月29日 (2006.6.29)	(74) 代理人	100121061
(33) 優先権主張国	米国 (US)		弁理士 西山 清春
前置審査		(72) 発明者	ツェン, リン
			アメリカ合衆国カリフォルニア州9408
			9, サニーヴェール, イースト・ジャバ・
			ドライブ・495
			最終頁に続く

(54) 【発明の名称】 ブロック指紋を読み出し、ブロック指紋を使用してデータ重複を解消するシステム、及び方法

(57) 【特許請求の範囲】

【請求項1】

ストレージシステムにおいて重複データを識別するために、コンピュータで実施される方法であって、

複数のデータブロックにバイト単位でデータを記憶するステップと、

各チェックサム値が各データブロックについて計算される、複数のチェックサム値を計算するステップと、

前記データブロック内の所定のオフセットからデータをサンプリングするステップと、

各チェックサム値及び前記サンプリングされたデータを複数の指紋のうちの1つの指紋に記憶するステップであって、各指紋が、前記複数のデータブロックのうちの1つのデータブロックに関連し、前記データブロックについて計算されたチェックサム値、及び前記データブロックからのサンプリングされたデータを含む、各チェックサム値及び前記サンプリングされたデータを複数の指紋のうちの1つの指紋に記憶するステップと、

(1) 前記複数の指紋のうちの同一の指紋を識別し、(2) 前記同一の指紋に関連する前記識別された重複データブロックを比較し、前記識別された重複データブロックが実際に重複していることを確認することにより、前記複数のデータブロックのうちの重複データブロックを識別するステップであって、前記比較は、前記識別された重複データブロックのバイトごとの比較を実施することを含む、前記複数のデータブロックのうちの重複データブロックを識別するステップと

10

20

を含む、方法。

【請求項 2】

前記複数のデータブロックのうちの 1 つに記憶されたデータを変更するステップと、
前記所定のオフセットからサンプリングされたデータが、前記複数のデータブロックの
うちの 1 つの中の前記所定のオフセットに記憶された前記変更されたデータとは異なること
に回答して、前記複数の指紋のうちの 1 つを削除するステップと
をさらに含む、請求項 1 に記載のコンピュータで実施される方法。

【請求項 3】

前記複数のデータブロックのうちの 1 つに関連するメタデータを記憶するステップと、
前記メタデータを前記複数のデータブロックのうちの 1 つに関連する指紋レコードにコ
ピーするステップと、
前記複数のデータブロックのうちの 1 つに関連する指紋を前記指紋レコードに記憶する
ステップと
をさらに含む、請求項 1 または請求項 2 に記載のコンピュータで実施される方法。

10

【請求項 4】

前記メタデータは、生成番号コピーとして前記指紋レコードにコピーされたインデック
ス・ノード (i n o d e) の生成番号を含み、前記方法は、
前記 i n o d e を削除するステップと、
書き込みオペレーションに回答して、前記 i n o d e を再割り当てするステップと、
前記生成番号を変更するステップと、
前記生成番号コピーが、変更された生成番号とは異なることに回答して、前記指紋レ
コードを削除するステップと
をさらに含む、請求項 3 に記載のコンピュータで実施される方法。

20

【請求項 5】

前記メタデータは、CP カウントコピーとして前記指紋レコードにコピーされたコンシ
ステンシ・ポイント (C P) を含み、前記方法は、
書き込みオペレーションに回答して、前記 CP カウントを変更するステップと、
前記 CP カウントコピーが、変更された CP カウントとは異なることに回答して、前記
指紋レコードを削除するステップと
をさらに含む、請求項 3 または請求項 4 に記載のコンピュータで実施される方法。

30

【請求項 6】

前記メタデータは、データコンテナ内のデータブロックの数に関連するサイズ値を含み
、さらに、f b n コピーとして前記指紋レコードにコピーされたファイルブロック番号 (
f b n) を含み、前記方法は、
前記データコンテナ内のデータブロックの数の変化に回答して、前記サイズ値を変更
するステップと、
前記 f b n コピーが、変更されたサイズ値に関連するデータブロックの数よりも大き
なオフセットを示していることに回答して、前記指紋レコードを削除するステップと
をさらに含む、請求項 3 ~ 5 のいずれか一項に記載のコンピュータで実施される方法。

【請求項 7】

ストレージオペレーティングシステムを実行するプロセッサを含むストレージシステム
であって、前記ストレージオペレーティングシステムが、
ファイルシステム層と、
前記ファイルシステム層と協働し、書き込みオペレーションを処理するストレージモ
ジュール層であって、バイト単位でデータブロックを記憶し、前記データブロックについ
てチェックサム値を計算するようにさらに構成されたストレージモジュール層と、
前記データブロックの 1 以上のオフセットにおけるデータをサンプリングし、前記デ
ータブロックについての前記チェックサム値及び前記データブロックからの前記サンプリ
ングされたデータを含む指紋を生成し、前記指紋を指紋データベースに記憶するように構
成された指紋抽出層であって、前記指紋データベースが、複数の指紋を含み、各指紋が、

40

50

複数のデータブロックのうちの1つのデータブロックに関連する、指紋抽出層と、

(1) 前記指紋データベース中の同一の指紋を識別し、(2) 前記同一の指紋に関連する前記識別された重複データブロックを比較し、前記識別された重複データブロックが実際に重複していることを確認することにより、重複データブロックを識別するように構成された重複解消モジュールであって、前記比較は、前記識別された重複データブロックのバイトごとの比較を実施することを含む、重複解消モジュールと

を含む、ストレージシステム。

【請求項 8】

前記書き込みオペレーションにตอบสนองして、前記複数の指紋のうちの少なくとも1つの指紋を記憶するように構成された指紋キューをさらに含み、前記重複解消モジュールは、前記指紋キュー、及び前記指紋データベースの中から同一の指紋を識別することにより、重複データブロックを識別するようにさらに構成される、請求項 7 に記載のストレージシステム。

10

【請求項 9】

前記ストレージモジュール層は、RAID (Redundant Array of Inexpensive Disks) プロトコルを実施するように構成され、前記チェックサム値は、前記 RAID プロトコルに従って計算される、請求項 7 または請求項 8 に記載のストレージシステム。

【請求項 10】

前記チェックサム値は、32ビット RAID (Redundant Array of Inexpensive Disks) チェックサム値である、請求項 1 ~ 6 のいずれか一項に記載のコンピュータで実施される方法。

20

【発明の詳細な説明】

【技術分野】

【0001】

[発明の分野]

本発明は、ストレージシステムに関し、特に、ストレージシステムにおいて重複データを効率的に減らす技術に関する。

【背景技術】

【0002】

[発明の背景]

30

ストレージシステムは一般に1以上の記憶装置を備え、要望に応じてその中に情報を入れたり、そこから情報を取り出したりすることができる。ストレージシステムは、ストレージオペレーティングシステムを含み、ストレージオペレーティングシステムは、とりわけ、システムによって実施されるストレージサービスを支える種々のストレージオペレーションを実施することにより、システムを機能的に編成する。ストレージシステムは、種々のストレージアーキテクチャに従って実施することができ、限定はしないが、例えば、ネットワーク・アタッチド・ストレージ環境、ストレージ・エリア・ネットワーク、あるいは、クライアントコンピュータ、又はホストコンピュータに直接取り付けられたディスクアセンブリとして実施される場合がある。記憶装置は通常、ディスクアレイとして編成され、ストレージプロトコルに従って管理される幾つかのディスクドライブである。ここで、「ディスク」という用語は一般に、内蔵型の回転式磁気媒体記憶装置を意味する。この文脈におけるディスクという用語は、ハード・ディスク・ドライブ(HDD)やダイレクト・アクセス・ストレージ・デバイス(DASD)と同じような意味である。

40

【0003】

ディスクアレイへの情報の記憶は、好ましくは、ディスク空間の全体的論理構成を定義する、物理ディスクの1以上のストレージ「ボリューム」として実施される。ボリューム内のディスクは通常、1以上のグループに編成され、各グループが、RAID (Redundant Array of Independent or Inexpensive Disks) として運用される。大半の RAID 実施形態は、RAID グループ内の所与の数の物理ディスクにわたってデータ「ストライプ」を冗長書き込みし、ストライピングされたデータに関する冗長情報(パリティ)を適切

50

に記憶することにより、データ記憶の信頼性/完全性を向上させる。各RAIDグループの物理ディスクには、ストライピングされたデータを記憶するように構成されたディスク(すなわち、データディスク)と、データのパリティを記憶するように構成されたディスク(すなわち、パリティディスク)とがある。その後、ディスクが故障したときに、失われたデータの復旧を可能にするために、パリティは読みされる場合がある。「RAID」という用語、及びその種々の実施形態は周知であり、1998年6月、データ管理に関する国際会議(SIGMOD)の議事録において、D.A. Patterson、G.A. Gibson、及びR.H. Katzにより、「A Case for Redundant Arrays of Inexpensive Disks (RAID)」と題して開示されている。

【0004】

ストレージシステムのストレージオペレーティングシステムは、ファイルシステムのような上位モジュールを実施し、情報を入れるためのデータコンテナを論理編成する場合がある。例えば、情報は、ディレクトリ、ファイル、及びブロックの階層構造としてディスク上に記憶される場合がある。例えば、「ディスク上」の各ファイルは、ファイルの実際のデータのような情報を記憶するように構成された一組のデータ構造、すなわちディスクブロックとして実施される場合がある。これらのデータブロックは、ファイルシステムによって管理されるボリュームブロック番号(vbn)空間の中に編成される。ファイルシステムは更に、ファイル内の各データブロックに、対応する「ファイルオフセット」、又はファイルブロック番号(fbn)を割り当てる場合がある。ファイルシステムは通常、ファイルごとに連続したfbnを割り当てるのに対し、vbnは、より大きなボリュームアドレス空間に割り当てられる。ファイルシステムは、vbn空間内において幾つかのデータブロックを「論理ボリューム」として編成する。そして、もちろん必須ではないが、各論理ボリュームは、その論理ボリューム独自のファイルシステムに関連付けられる場合がある。ファイルシステムのサイズが、n+1ブロックである場合、ファイルシステムは通常、ゼロからnまでの連続した範囲のvbnから構成される。

【0005】

一つの既知のタイプのファイルシステムは、ディスク上でデータを上書きしないwrite-anywhereファイルシステムである。ディスクからストレージシステムのメモリ内へデータブロックを取り出し(読み出し)、新たなデータで「汚す」(すなわち、更新、又は変更)するとき、書き込み性能を最適化するために、以後、そのデータブロックは、ディスク上の新たな場所に記憶(書き込み)される。write-anywhereファイルシステムは、データがディスク上で実質的に連続的に配置されるような最適レイアウトを最初に仮定する場合がある。この最適なディスクレイアウトにより、ディスクに対する効率的なアクセスオペレーションが可能となり、特に、シーケンシャル読み出しオペレーションの場合に効率的なアクセスオペレーションが可能となる。ストレージシステム上で動作するように構成されたwrite-anywhereファイルシステムの一例は、カリフォルニア州サニーベイルにあるネットワーク アプライアンス、インコーポレイテッドから販売されているWrite Anywhere File Layout (WAFL)ファイルシステムである。

【0006】

大きなファイルシステムでは、個々のデータブロックの重複的発生を見付けることが一般的である。データブロックの重複が発生するのは、例えば、2以上のファイルその他のデータコンテナが、共通のデータを共有する場合や、あるいは、所与のファイル内の複数の場所に、所与のデータセットが発生する場合である。データブロックの重複の結果、ストレージシステムによって提供される複数の異なる場所に同一のデータが記憶されることにより、記憶空間は、十分に利用されないことになる。

【0007】

この問題を解決するために使用されている一つの技術として、「ファイル・フォールディング」と呼ばれるものがある。ファイル・フォールディングの基本原理は、アクティブファイルシステム上のファイルの新たなデータが、永久的イメージ中のそのファイルの古

10

20

30

40

50

いデータと同一のものである場合、新たなデータと古いデータとの間におけるディスクブロックの共有を許可するというものである。この技術は、カリフォルニア州、サニーベールにあるネットワーク アプライアンス、インコーポレイテッドから市販されている種々のストレージシステムにおいて既に実施されている。具体的には、それらのストレージシステムは、指定されたデータセットの空間保守的な、コンシステンシ・ポイント・イメージ（PCPI）を生成する機能を有する。PCPIとは、名前によるアクセスが可能な、空間保守的な、何らかの過去の時点におけるデータ（例えばストレージシステム）の読み出し専用イメージであり、何らかの過去の時点におけるそのデータの一貫性のあるイメージを提供する。具体的には、PCPIとは、記憶装置（例えばディスク）その他の永久的メモリに記憶されたアクティブファイルシステム、ファイル、又はデータベースのような記憶要素をある時点において表現した物であり、自分自身を他の時点において作成されたPCPIから区別するための名前その他の識別子を有する。PCPIの一例は、ネットワーク アプライアンス、インコーポレイテッドから市販されているストレージシステムにおいて実施されるようなSnapshotである。「PCPI」という用語と「Snapshot」という用語は、ネットワークアプライアンスの商標権を失墜させることなく、交換可能に使用される。

【0008】

PCPIを生成する技術については、1998年10月6日に発行された「Method For Maintaining Consistent States of a File System and for Creating User-Accessible Read-Only Copies of a File System」と題するDavid Hitz他による本件と同じ譲受人の米国特許第5,819,292号に詳しく記載されている。例えば、「スナップショット化」されたデータコンテナの中のブロックが変更された場合、ストレージシステムは、そのデータコンテナのもう一つの完全なコピー（変更後のもの）を作成する代わりに、単に、そのデータコンテナの変更されたブロックをアクティブファイルシステム上に生成する。未変更のブロックについては、ストレージシステムは単に、PCPI中の対応するブロックへのポインタをそのデータコンテナに与えるだけである。このようにして、PCPI中の未変更のブロックは、PCPIとアクティブファイルシステムとの間で共有されることになる。この技術については、2002年3月22日に出願された「File Folding Technique」と題する本件と同じ譲受人の米国特許出願第2003/0182317号に詳しく記載されている。

【0009】

ファイルフォールディングは、記憶空間の利用効率を向上させるためには役に立たない。しかしながら、SnapshotのようなPCPIに頼ることなく、アクティブファイルシステムにおけるデータ重複を減らすことが望ましい。また、重複ファイル、又は単一ファイル内の重複データによりアクティブファイルシステム上に発生することがある重複データブロックを識別し、減らすことが望まれている。より一般的には、ストレージシステムにおいて、のデータの位置とは無関係にデータ重複を減らすことが望まれている。

【0010】

データ重複の低減を実現する他の技術は、1999年11月23日に発行された「Partitioning a Block of Data into Blocks and for Storing and Communicating Such Subblocks」と題するRoss Williamsによる米国特許第5,990,810（以下、「810特許」）に記載されている。「810特許」に記載の方法は、まず、ローリングハッシュ関数を使用して、複数のサブブロックのデータを生成する。ローリングハッシュ関数は、固定サイズウィンドウのデータを使用し、2つのサブブロック間に位置する境界を生成する。一つのデータブロックを複数のサブブロックに分割した後、各サブブロックのハッシュ値が計算され、ハッシュ値のテーブルが形成される。次に、このハッシュテーブルを使用して、新たなサブブロックが、ハッシュテーブルに既に記憶されたハッシュ値を有するいずれかのサブブロックと同一であるか否かが判定される。この判定を実施するために、新たなサブブロックのハッシュ値が計算され、それが、ハッシュテーブルに含まれるハッシュ値と比較される。もし、新たなサブブロックのハッシュ値が、ハッシュテーブルに既に記憶

10

20

30

40

50

されていた場合、記憶されたハッシュ値を有するものとして識別されたサブブロックは、新たなサブブロックと同一であるとみなされる。その場合、新たなサブブロックは、既に記憶されているサブブロックへのポインタで置き換えられ、それによって、サブブロックに必要とされる記憶空間の量が低減される。' 8 1 0 特許に記載されている技術の大きな欠点は、計算上、膨大な数の集中的なハッシュ計算の処理能力が要求される点にあり、そのような計算は、そのような方法を実施するストレージシステムの全体的処理能力に悪影響を与えることがある。もう一つの大きな欠点は、データセットのサイズが増大するにつれて、ハッシュテーブルも大きくなるため、テラバイトやペタバイトのデータのような大きなデータセットには応用できないことがある点である。

【 0 0 1 1 】

重複データを無くすさらに別の技術は、2005年4月13日に出願された「Method and Apparatus for Identifying and Eliminating Duplicate Datablocks and Sharing Data Blocks in a Storage System」と題するLing Zheng他による本件と同じ譲受人の米国特許出願第11/105,895に記載されており、この文献の内容は、参照により、本明細書に援用される。この特許出願に記載されているシステムでは、全てのデータ重複解消オペレーションが、例えば4キロバイト(KB)サイズといった固定サイズのブロックに対して実施される。新たなブロックを記憶する場合、ハッシュ値は、4KBブロックの指紋として計算される。次に、この指紋は、以前に記憶された種々のブロックの指紋を含むハッシュテーブルと比較される。もし、新たなブロックの指紋が、以前に記憶されたブロックのものと同じであった場合、その新たなブロックは、以前に記憶されたブロックと同一である可能性が高い。その場合、新たなブロックは、以前に記憶されたブロックへのポインタで置き換えられ、それによって、ストレージリソースの消費が低減される。

【 0 0 1 2 】

しかしながら、このシステムの欠点は、新たな4KBブロックを記憶するときハッシュ計算が必要となる点にある。指紋やデータブロックのサイズによっては、そのような計算は、処理リソースをさらに必要とする場合がある。そして、それらのサイズは、さらに、ストレージシステムのサイズによって決まる場合がある。また、可能性は非常に低い、同一の複数の指紋が、重複データを意味しないこともありうる。例えば、重複解消プロセスにおいて、同一の指紋を有するデータブロックは、それらが実際に同一であることを確認するために、さらに比較される場合がある。従って、その後のデータ重複解消プロセスにおけるデータブロックの比較を容易にするためには、同一の指紋が重複データを意味するであろうという適度な確率を維持しながらも、より少ない処理リソースを使用して指紋を生成できることが望ましい。さらに、NetApp Data ONTAPオペレーティングシステムのような現代のストレージシステムの中には、システムに記憶されるデータブロックの完全性を確保するために、既にチェックサムオペレーションを実施しているものもある。既存のストレージ機能を利用して処理リソースを節約する重複解消プロセスが、望まれている。

【 発明の概要 】

【 課題を解決するための手段 】

【 0 0 1 3 】

[発明の概要]

本発明は、ストレージシステムにおける重複データブロックを効率的に識別できるようにブロック指紋を最適化することにより、従来技術の欠点を克服する。識別された重複データブロックは、次に、データ重複解消モジュールによって比較され、ストレージシステムの処理能力を最適化するために、重複データは削除される。

【 0 0 1 4 】

本発明の一実施形態によれば、ストレージシステムは、ファイルシステムのような高レベルモジュールを実施し、ストレージモジュールによって管理される複数のディスクのような記憶装置に記憶される情報を論理編成する。上位モジュールとストレージモジュールの間の境界において、指紋抽出層は、ストレージシステムにより記憶される各データブ

10

20

30

40

50

ックについて、例えば64ビットのような所定サイズの指紋を生成する。例えば、指紋の最初の所定数のビットは、ストレージモジュールによって実施されるチェックサム計算の結果から得られ、指紋の次の所定数のビットは、データブロック中の1以上の所定のオフセットに記憶されているデータから得られる。あるいは、ストレージモジュールと重複解消モジュールは、同じ64ビットチェックサム値を使用して、異なる目的に従事する場合がある。得られた指紋は、データブロックの中身を十分に反映し、望ましくない誤認識率を生じることなく、重複の識別を可能にする。また、指紋の生成のために、ハッシュ関数のような他の暗号化関数をさらに実施する必要はないため、プロセッサ使用率も低減される。

【0015】

各指紋は、指紋レコードに記憶され、指紋レコードはさらに、指紋データベースに記憶され、データ重複解消モジュールによってアクセスされる。データ重複解消モジュールは、指紋を定期的に比較し、重複する指紋を識別し、さらに重複データブロックを識別する場合がある。

【0016】

例示的实施形態によれば、各指紋レコードは、指紋が付けられたデータブロックに関する付加情報を有する。データ重複解消モジュールは、この付加情報をデータブロックに関連するメタデータと定期的に比較し、指紋データベースが、ストレージシステムの最新の状態を正確に反映するようにする。従って、「古くなった」指紋(例えば、消去されたデータブロックや変更されたデータブロックに関連する指紋)は、指紋データベースから削除される場合がある。あるいは、古くなった指紋を削除するプロセスは、指紋データベースが所定のサイズに達するといった、所定の条件に応じて実施され、又は、古くなった指紋を削除するプロセスは、指紋を比較する上記のプロセスの直前に実施される場合がある。

【図面の簡単な説明】

【0017】

【図1】本発明の一実施形態によるストレージシステム環境を示す略ブロック図である。

【図2】本発明の一実施形態によるストレージオペレーティングシステムを示す略ブロック図である。

【図3】本発明の一実施形態によるinodeを示す略ブロック図である。

【図4】本発明の一実施形態によるデータコンテナのバッファツリーを示す略ブロック図である。

【図5】本発明の一実施形態によるデータブロックへのポインタを有する間接ブロックを示す略ブロック図である。

【図6】本発明の一実施形態による変更されたデータブロックへのポインタを有する間接ブロックを示す略ブロック図である。

【図7】本発明の一実施形態による重複するデータブロックへのポインタを有する間接ブロックを示す略ブロック図である。

【図8】本発明の一実施形態による指紋を示す略ブロック図である。

【図9】本発明の一実施形態による指紋レコードを示す略ブロック図である。

【図10】本発明の一実施形態による、指紋レコードを生成する手順のステップの詳細を示すフロー図である。

【図11】本発明の一実施形態による、重複データブロックを識別する手順のステップの詳細を示すフロー図である。

【図12A】本発明の一実施形態による重複データブロックの識別を示す指紋キュー、及び指紋データベースを示す略ブロック図である。

【図12B】本発明の一実施形態による重複データブロックの識別を示す指紋キュー、及び指紋データベースを示す略ブロック図である。

【図12C】本発明の一実施形態による重複データブロックの識別を示す指紋キュー、及び指紋データベースを示す略ブロック図である。

10

20

30

40

50

【図12D】本発明の一実施形態による重複データブロックの識別を示す指紋キュー、及び指紋データベースを示す略ブロック図である。

【図13】本発明の一実施形態による、古くなった指紋レコードを削除する手順のステップの詳細を示すフロー図である。

【発明を実施するための形態】

【0018】

[例示的实施形態の詳細な説明]

A. ストレージシステム環境

図1は、本発明とともに有利に使用されるストレージシステム120を含む例示的ストレージシステム環境100を示す略ブロック図である。ストレージシステムは、ディスク
10 160のディスク130のような記憶装置上での情報の編成に関連するストレージ
サービスを提供するコンピュータである。ストレージシステム120は、システムバス1
25によって相互接続されたプロセッサ122、メモリ124、ネットワークアダプタ1
26、不揮発性ランダムアクセスメモリ(NVRAM)129、及びストレージアダプタ
128を含む。ストレージシステム120はさらに、ストレージオペレーティングシステ
ム200を含み、ストレージオペレーティングシステム200は、好ましくは、ファイル
システムのような上位モジュールを実施し、情報をディレクトリ、ファイル、及び仮想デ
ィスク(以下、「ブロック」と呼ばれる特殊なタイプのファイルの階層構造として論理
編成する。

【0019】

例示的実施形態として、メモリ124は、ソフトウェアプログラムコードを記憶するた
めに、プロセッサ122、及びアダプタによってアドレス指定可能な多数の記憶場所を有
する。メモリの一部は、本発明に関連する指紋キューのようなデータ構造を記憶するた
めに、「バッファキャッシュ」170としてさらに編成される場合がある。あるいは、NV
RAM129は、本発明に関連する指紋キュー、指紋データベース、その他のデータ構造
を記憶するために使用される場合がある。さらに、システムの永久的コンシステンシ
・ポイント・イメージ(PCPI)をサポートするシステムでは、NVRAM129を使用し
て、コンシステンシ・ポイント間においてファイルシステムに対して施された変更を記憶
する場合がある。そのような変更は、コンシステンシ・ポイントの進行中に不揮発性ログ
(NVL0G)131に記憶され、すなわち、ディスクにフラッシュ(書き込み)される
10 場合がある。さらに、プロセッサ122、及びアダプタは、ソフトウェアコードを実行し
、データ構造を操作するように構成された処理要素、及び/又は論理回路を含む場合が
ある。ストレージオペレーティングシステム200は、その一部が通常、メモリに常駐し
、処理要素によって実行され、とりわけ、ストレージシステムによって実施されるストレ
ージオペレーションを実行することにより、システム120を機能的に編成する。当業者
には明らかなように、本明細書に記載する本発明の技術に関連するプログラム命令の記憶
、及び実行には、他の処理手段や、種々のコンピュータ読取可能媒体を含む他の記憶手段
を使用してよい。

【0020】

ネットワークアダプタ126は、コンピュータネットワーク140を介してストレージ
システム120をクライアント110に接続するために必要な機械的、電氣的、及び信号
的回路を含み、コンピュータネットワーク140は、ポイント・ツーポイント接続、又は
ローカルエリアネットワークのような共有媒体を含む場合がある。例えば、コンピ
ュータネットワーク140は、イーサネット(R)ネットワーク、又はファイバチャネル(FC
)ネットワークとして実施される場合がある。クライアント110は、トランスミッシ
ョン・コントロール・プロトコル/インターネット・プロトコル(TCP/IP)のよう
な所定のプロトコルに従って個々のデータフレーム、又はデータパケット150をや
りとりすることにより、ネットワーク140を介してストレージシステムと通信するこ
とができる。

【0021】

10

20

30

40

50

クライアント110は、アプリケーション112を実行するように構成された汎用コンピュータであってよい。また、クライアント110は、クライアント/サーバモデルの情報配送に従ってストレージシステム120との間で情報をやりとりする場合がある。すなわち、ネットワーク140を介してパケット150をやりとりすることにより、クライアント110はストレージシステム120にサービスを要求する場合があり、システムはクライアント110から要求されたサービスの結果を返す場合がある。クライアント110は、ファイルやディレクトリの形をした情報をアクセスするときは、TCP/IPを介して、コモン・インターネット・ファイル・システム(CIFS)プロトコル、又はネットワーク・ファイル・システム(NFS)プロトコルのようなファイルベースのアクセスプロトコルを含むパケットを発行する場合がある。一方、クライアントは、ブロックの形をした情報をアクセスするときは、SCSI over TCP(iSCSI)プロトコルや、SCSI over FC(FCP)プロトコルのようなブロックベースのアクセスプロトコルを含むパケットを発行する場合がある。

【0022】

ストレージアダプタ128は、システム120上で実行されているストレージオペレーティングシステム200と協働し、ユーザ(すなわち、クライアント)から要求された情報をアクセスする。情報は、書き込み可能なストレージデバイス媒体の任意のタイプのアタッチド・アレイに記憶され、媒体には、例えば、ビデオテープ、光学的、DVD、磁気テープ、バブルメモリ、電氣的ランダムアクセスメモリ、又はMEMSデバイスなどがあり、あるいはデータやパリティ情報を記憶するように構成された任意の他の適当な媒体も含まれる。ただし、本明細書で例示するように、情報は、アレイ160のハードディスクドライブ(HDD)、及び/又はダイレクト・アクセス・ストレージ・デバイス(DASD)のようなディスク130に記憶されることが好ましい。ストレージアダプタ128は、従来の高性能FCシリアルリンクトポロジのようなI/O相互接続構成によってディスクに接続される入出力(I/O)インタフェース回路を含む。

【0023】

アレイ160への情報の記憶は、好ましくは、1以上のストレージ「ボリューム」として実施され、ボリュームは、一群の物理記憶ディスク130を含み、それらが協働して、ボリューム(複数の場合もあり)上のボリュームブロック番号(vbn)空間の全体的論理構成を規定する。当然ながら必須ではないが、各論理ボリュームは一般に、各自、自分のファイルシステムに関連している。論理ボリューム/ファイルシステム内のディスクは通常、1以上のグループに編成され、各グループは、RAID(Redundant Array of Independent or Inexpensive Disks)として運用され、RAIDプロトコルに従って管理される場合がある。RAID4レベル実施形態のような大半のRAID実施形態は、RAIDグループ内の所与の数の物理ディスクにわたってデータ「ストライプ」を冗長書き込みし、更に、そのストライピングされたデータに関連するパリティ情報を適切に記憶することにより、データ記憶の信頼性/完全性を向上させる。例示的RAIDプロトコルは更に、ディスクに書き込まれる各データブロックについて、データの完全性を確認するために32ビットチェックサム値を実施する。RAID実施形態の一例は、RAID4レベル実施形態であるが、当然ながら、本明細書に記載する本発明の原理によれば、他のタイプ、又は他のレベルのRAID実施形態を使用することも可能であるものと考えられる。同様に、論理ボリューム/ファイルシステムのディスクの編成には、他の実施形態、及び/又はプロトコルが使用される場合もある。

【0024】

B. ストレージオペレーティングシステム

ディスク130へのアクセスを容易にするために、ストレージシステム200は、write-anywhereファイルシステムを実施する。このファイルシステムは、仮想化モジュールと協働し、ディスク130により提供される記憶空間を仮想化する。例えば、ファイルシステムは、情報を名前付きディレクトリ、及びファイルの階層構造としてディスク上に論理編成する。「ディスク」上の各ファイルは、データのような情報を記憶す

10

20

30

40

50

るように構成された一組のディスクブロックとして実施される場合がある一方、ディレクトリは、特殊な形のファイルとして実施され、その中に他のファイルやディレクトリの名前、あるいはそれらへのリンクが記憶される場合がある。仮想化モジュールによれば、ファイルシステムは、情報をディスク上のブロックの階層構造としてさらに論理編成し、それらを名前付きの論理ユニット番号 (LUN) としてエクスポートすることが可能になる。

【0025】

例示的实施形態として、ストレージオペレーティングシステムは、カリフォルニア州サンニールにあるネットワーク アプライアンス、インコーポレイテッドから販売されている NetApp Data ONTAP オペレーティングシステムであることが好ましい。ただし、当然ながら、任意の適当なストレージオペレーティングシステムを本明細書に記載する本発明の原理に従って拡張して使用することも可能であるものと考えられる。したがって、「Data ONTAP」という用語を使用した場合でも、この用語は、本発明の教示に適合する任意のストレージオペレーティングシステムを指すものとして広い意味で捉えなければならない。

【0026】

図2は、本発明とともに有利に使用されるストレージオペレーティングシステム200を示す略ブロック図である。ストレージオペレーティングシステムは、統合ネットワークプロトコルスタックを形成するように編成された一連のソフトウェア層を含む。すなわち、より一般的に言えば、ストレージオペレーティングシステムは、ストレージシステムに記憶された情報をブロックアクセスプロトコル、及びファイルアクセスプロトコルを使用してアクセスするときのデータパスをクライアントに提供するマルチプロトコルエンジンを含む。プロトコルスタックは、IP層212、並びにその支援搬送手段であるTCP層214、及びユーザ・データグラム・プロトコル(UDP)層216のような種々のネットワークプロトコル層との間のインタフェースとして機能するネットワークドライバ(例えば、ギガビット・イーサネット(R)・ドライバ)のメディアアクセス層210を含む。ファイルシステムプロトコル層は、マルチプロトコルファイルアクセスを可能にし、その目的のために、ダイレクトアクセスファイルシステム(DAFS)プロトコル218、NFSプロトコル220、CIFSプロトコル222、及びハイパーテキスト・トランスファ・プロトコル(HTTP)プロトコル224をサポートする。仮想インタフェース(VI)層226は、VIAアーキテクチャを実施し、DAFSプロトコル218によって必要とされるようなRDMAのようなダイレクト・アクセス・トランスポート(DAT)機能を提供する。

【0027】

iSCSIドライバ層228は、TCP/IPネットワークプロトコル層を介したブロックプロトコルアクセスを提供する一方、FCドライバ層230は、ストレージシステムとの間で、ブロックアクセス要求、及び応答の送受信を行う。FCドライバ、及びiSCSIドライバは、ブロックに対するFC固有の、又はiSCSI固有のアクセス制御を提供し、したがって、ストレージシステム上のブロックをアクセスするときに、iSCSIとFCPのいずれか一方、あるいは両方へのLUNのエクスポートを管理する。さらに、ストレージオペレーティングシステムは、RAIDシステム240として実施され、I/Oオペレーションに従ってボリューム/ディスクに対する情報の記憶、又は読み出しを管理するストレージ(例えばRAID)プロトコルを実施するストレージモジュールと、SCSIプロトコルのようなディスクアクセスプロトコルを実施するディスクドライバシステム250とを含む。

【0028】

ディスクソフトウェア層を統合ネットワークプロトコルスタック層に橋渡しするのは、仮想化システムである。仮想化システムはファイルシステム280によって実施され、ファイルシステム280は、例えばvディスクモジュール290、及びSCSIターゲットモジュール270として実施される仮想化モジュールとの間で情報をやりとりする。vデ

10

20

30

40

50

ディスクモジュール290は、ファイルシステム280の上に層として形成され、ストレージシステムに対するユーザ（システム管理者）発行コマンドにตอบสนองして、ユーザインタフェース（UI）275のような管理インタフェースによるアクセスを有効化する。SCSIターゲットモジュール270は、FC、iSCSIドライバ228, 230と、ファイルシステム280との間に配置され、LUNがブロックとして表現される場合に、ブロック（LUN）空間とファイルシステム空間の間の変換層を仮想化システムに提供する。UI275は、種々の層、及びシステムへの管理者、又はユーザによるアクセスが可能となるような形で、ストレージオペレーティングシステム上に配置される。

【0029】

本発明の一実施形態によれば、ストレージオペレーティングシステム200は、重複解消モジュール284、及び指紋抽出層286をさらに含み、指紋抽出層286は、ファイルシステム280とRAIDシステム240の間に配置される。重複解消モジュール284は、所定の条件（例えば、ストレージシステムが所定の容量に達する等）にตอบสนองして定期的に、及び/又はUI275を通じて入力された管理者、又はユーザからのコマンドにตอบสนองして、データ重複解消プロセスを実施する場合がある。重複解消モジュール284、及び指紋抽出層286は、指紋キュー287、及び指紋データベース288との間で情報をやりとりする。指紋キュー287、及び指紋データベース288はそれぞれ、指紋レコードを記憶し、指紋レコードを使用して、重複データブロックを識別する。例えば、代替実施形態として、指紋キュー287、及び指紋データベース288は、ストレージシステムメモリ124、及び/又はNVRAM129に常駐する場合もあるが、それらは個別に、ディスク、及び/又は任意の他のコンピュータ読取可能媒体に記憶される場合もある。また、代替実施形態によれば、ディスクアレイ160への書き込みアロケーションオペレーションの実行中にデータ重複解消手順を実施するために、データ重複解消モジュール284は、書き込みアロケータプロセス282の中に統合される場合もある。「Method and Apparatus for Identifying and Eliminating Duplicate Data Blocks and Sharing Data Blocks in a Storage System」と題する上で援用した米国特許出願に記載されているように、そのようなデータ重複解消手順は、同一のデータを有する複数のデータブロックを指し示す代わりに単一データブロックを指し示すように、間接ブロック中のポインタを変更する場合がある。

【0030】

ファイルシステムは例えば、ディスクのような記憶装置に記憶された情報をアクセスする際に使用される論理ボリューム管理機能を備えたメッセージベースのシステムである。すなわち、ファイルシステム280は、ファイルシステムセマンティックを提供するだけでなく、通常ならばボリュームマネージャに関係する機能も提供する。そうした機能には、例えば、（1）ディスクのアグレゲーション、（2）ディスクの記憶帯域幅のアグレゲーション、及び（3）ミラーリング、及び/又はパリティ（RAID）のような信頼性保証がある。ファイルシステム280は、例えば4キロバイト（KB）ブロックを使用し、インデックス・ノード（「inode」）を使用してファイルやファイル属性（例えば、作成時刻、アクセス許可属性、サイズ、及びブロック位置）を識別するブロックベースのオンディスク・フォーマット表現を有するwrite-anywhereファイルシステム（例えば、カリフォルニア州サニーベイルにあるネットワーク・アプライアンス、インコーポレイテッドから販売されているWAFファイルシステムなど）を実施する。ファイルシステムは、ファイルを使用して、自分が有するファイルシステムのレイアウトを表すメタデータを記憶する。そのようなメタデータファイルとしては、とりわけ、inodeファイルがある。ディスクからinodeを読み出すためには、ファイルハンドル、すなわちinode番号を含む識別子が使用される。

【0031】

簡単に言えば、write-anywhereファイルシステムのinodeは全て、inodeファイルとして編成される。ファイルシステム（FS）infoブロックは、ファイルシステム内の情報のレイアウトを指定し、ファイルシステムの他の全てのinod

10

20

30

40

50

eを含むファイルのinodeを含む。各論理ボリューム（ファイルシステム）は、FS infoブロックを有し、FS infoブロックは、例えばRAIDグループ内の固定位置に記憶されることが好ましい。inodeファイルのinodeはinodeファイルのブロックを直接参照（指し示す）場合もあれば、inodeファイルの間接ブロックを参照し、さらに、その間接ブロックが、inodeファイルのブロックを指し示す場合もある。inodeファイルの直接ブロックの中にはそれぞれ、種々のinodeが埋め込まれ、それらのinodeがそれぞれ、間接ブロックを参照し、さらに、その間接ブロックが、ファイルのデータブロックを参照する。

【0032】

動作として、クライアント110からの要求は、パケット150としてコンピュータネットワーク140を介してストレージシステム120へ転送され、そこで、ネットワークアダプタ126により受信される。（層210、又は層230）のネットワークドライバは、そのパケットを処理し、必要に応じて、それをネットワークプロトコル層、及びファイルアクセス層へ渡し、更なる処理を施した後、write-anywhereファイルシステム280へと転送する。ここで、ファイルシステムは、要求されたファイルが、「コア」内、すなわちバッファキャッシュ170内に無ければ、それをディスク130からロードするオペレーションを生成する。情報がキャッシュ内に無ければ、ファイルシステム280は、inode番号を使用してinodeファイルを検索し、適当なエントリにアクセスし、論理vbnを読み出す。次にファイルシステムは、その論理vbnを含むメッセージ構造をRAIDシステム240へ渡す。そして、論理vbnは、ディスク識別子、及びディスクブロック番号(disk, dbn)にマッピングされ、ディスクドライバシステム250の適当なドライバへと送られる。ディスクドライバは、指定されたディスク130からdbnをアクセスし、ストレージシステムによる処理に備えて、要求されたデータブロック（複数の場合も有り）をバッファキャッシュ170に読み出す。要求が完了すると、ストレージシステム（及び、オペレーティングシステム）は、ネットワーク140を介してクライアント110に返答を返す。

【0033】

クライアント110からストレージシステム120に対して発行される書き込みオペレーションは、そのオペレーションが、後のディスク130へのコミットに備えて、NVRAM129のNVLOG131に一時的に記憶（「キューイング」）されることがある点を除けば、ストレージシステム100を貫通する同じ「パス」を通る。同様に、書き込みオペレーションは、その書き込みオペレーションによって変更された論理vbnに対応する指紋レコードを作成するために、指紋抽出層286、及び/又は重複解消モジュール284が、書き込みオペレーションの際にinodeファイル中の適当なエントリに関連するメタデータのようなデータを抽出する場合がある点を除けば、ストレージオペレーティングシステム200を貫通する同じソフトウェア「パス」を通る。

【0034】

なお、ストレージシステムにおいて受信されたクライアント要求に対し、データストレージアクセスや書き込みオペレーションを実施するために必要となる上記のストレージオペレーティングシステム層を貫通するソフトウェア「パス」は、代替としてハードウェアで実施してもよい。すなわち、本発明の代替実施形態では、ストレージアクセス要求データパスは、フィールド・プログラマブル・ゲート・アレイ(FPGA)、又は特定用途向け集積回路(ASIC)の中に埋め込まれた論理回路として実施される場合がある。この種のハードウェア実施形態によれば、クライアント110により発行された要求に回答してストレージシステム120によって提供されるストレージサービスの性能を向上させることができる。また、本発明の更に別の代替実施形態として、アダプタ126、128の処理要素は、パケット処理オペレーション、及びストレージアクセスオペレーションの負荷の一部、又は全部をプロセッサ122から取り除くように構成される場合があり、それによって、システムにより提供されるストレージサービスの処理能力を向上させる場合がある。当然ながら、本明細書に記載する種々のプロセス、アーキテクチャ、及び手順は、

10

20

30

40

50

単一のシステムとして実施しても、分散システムとして実施してもよく、また、ハードウェアで実施しても、ファームウェアで実施しても、ソフトウェアで実施してもよい。

【0035】

本明細書では、「ストレージオペレーティングシステム」という用語は通常、ストレージシステムにおいてストレージ機能を実施するためのコンピュータ実行可能コードを意味し、例えば、データアクセスを管理し、ファイルサーバの場合には、ファイルシステムセマンティックを実施することもあるコンピュータ実行可能コードを意味する。その意味で、Data ONTAPソフトウェアは、マイクロカーネルとして実施されるストレージオペレーティングシステムであって、かつWAF Lファイルシステムセマンティックを実施し、データアクセスを管理するためのファイルシステム280を含む、そのようなストレージオペレーティングの一例である。例えば、本発明によれば、ストレージオペレーティングシステム200は、指紋を比較し、データ重複解消を実施するための重複解消モジュール284を含む。しかしながら、代替実施形態では、指紋比較機能やデータ重複解消機能は、ストレージオペレーティングシステム200の他のモジュールで実施される場合もある。また、ストレージオペレーティングシステムは、UNIXやWindows XPのような汎用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プログラムとして実施してもよいし、本明細書に記載するようなストレージアプリケーションに合わせて構成された構成変更機能を備えた汎用オペレーティングシステムとして実施してもよい。

10

【0036】

さらに、当業者には分かるように、本明細書に記載する本発明の技術は、いかなるタイプの特殊目的のコンピュータ（例えばファイルサーバ、ファイラ、又はマルチプロトコルストレージアプライアンス）にも、汎用のコンピュータにも適用することができ、例えば、ストレージシステム120として実施され、又はストレージシステム120を含むスタンドアロンのコンピュータ、あるいはその一部にも適用することができる。本発明とともに有利に使用されるストレージシステムアプライアンスの一例は、2002年8月8日出願された「Protocol Storage Appliance that Provides Integrated Support for File and Block Access Protocols」と題するBrian Pawlowski他による米国特許出願公報第2004/0030668A1号に記載されている。また、本発明の教示は、種々のストレージシステムアーキテクチャに適合させることができ、限定はしないが、例えば、ネットワーク・アタッチド・ストレージ環境、ストレージ・エリア・ネットワーク、あるいは、クライアントコンピュータやホストコンピュータに直接取り付けされたディスクアセンブリにも適用することが可能である。したがって、「ストレージシステム」という用語は、ストレージ機能を実施するように構成され、他の装置又はシステムに関連する何らかのサブシステムだけでなく、そのような構成も含むものとして広い意味で捉えなければならない。

20

30

【0037】

C. オンディスク・ファイルシステム構造

例示的实施形態として、Write-anywhereファイルシステムでは、ファイル（又は、他のデータコンテナ）は、ディスク130への記憶に適したinodeデータ構造として表現される。図3は、inode300を示す略ブロック図である。inode300は、メタデータ部310、及びデータ部350を有することが好ましい。各inode300のメタデータ部310に記憶される情報はファイルを表し、例えば、ファイルのタイプ312（例えば、通常、ディレクトリ、仮想ディスクなど）、ファイルのサイズ314、ファイルのタイムスタンプ（アクセス、及び/又は変更）316、ファイルの所有者、例えばユーザID（UID318）、及びグループID（GID320）、コンシステンシ・ポイント（CP）フィールド322、並びに生成番号フィールド324を含む。

40

【0038】

CPカウントフィールド322は、特定バージョンのinodeが生成された時点におけるCPカウントを識別する。例えば、各CPは、単調増加CPカウンタによって生成さ

50

れた一意のCP番号によって識別される。代替実施形態では、CPカウンタは、各CPに一意の識別子を与えるような他の技術を使用して生成される場合がある。生成番号フィールド324は、特定inodeの生成を識別する。例えば、inodeが削除され、再使用されるたびに、そのinodeに関連する生成番号324は、インクリメントされる。write anywhereファイルシステムでは、inodeが変更されたときは常に、新たなコピーがディスクに書き込まれる。そのようなとき、すなわち、inode書き込みアロケーションの際に、それらのフィールド322、324は更新される。従って、生成番号は、inodeの生成/ア割り当てを反映し、CPカウンタは、inodeの変更を反映する。フィールド322、324は、特定のデータコンテナの指紋が生成された後、inodeによって表されるそのデータコンテナに変更があったか否かを迅速に判定するために使用される場合がある。

10

【0039】

ただし、各inodeのデータ部350の中身は、タイプフィールド312の中に規定されるファイル(inode)のタイプによって、異なる解釈をされる場合がある。例えば、ディレクトリinodeのデータ部350は、ファイルシステムによってコントロールされるメタデータを有する一方、通常inodeのデータ部は、ファイルシステムデータを有する場合がある。後者の場合、データ部350は、ファイルに関連するデータの表現物を含む場合がある。

【0040】

例示的实施形態によれば、通常オンディスクinodeのデータ部350は、ファイルシステムデータ、又はポインタを含む場合があり、後者は、ファイルシステムデータを記憶するために使用されるディスク上の4KBデータブロックを参照する。ディスク上のデータをアクセスするときのファイルシステムとRAIDシステム240の間の効率を上げるために、各ポインタは、論理vbnであることが好ましい。inodeのサイズには制限があるため(例えば128バイト)、64バイト以下のサイズのファイルシステムデータは、その全体が、そのinodeのデータ部の中に表現される。一方、ファイルシステムデータのサイズが64バイトよりも大きく、且つ64KB以下である場合、inode(例えば、第1レベルinode)のデータ部は、最大で16個までのポインタを有し、各ポインタが、ディスク上の4KBブロックのデータを参照する。

20

【0041】

また、データのサイズが64KBよりも大きく、且つ64メガバイト(MB)以下である場合、inode(例えば、第2レベルinode)のデータ部350にある各ポインタは、1024個のポインタを有する間接ブロック(例えば、第1レベルブロック)を参照し、各ポインタが、ディスク上の4KBデータブロックを参照する。64MBよりも大きなサイズのファイルシステムデータの場合、inode(例えば、第3レベルinode)のデータ部350にある各ポインタは、1024個のポインタを含む二重間接ブロック(例えば、第2レベルブロック)を参照し、さらに、各ポインタが、間接(例えば、第1レベル)ブロックを参照する。そして、その間接ブロックが、1024個のポインタを有し、各ポインタが、ディスク上の4KBのブロックを参照する。ファイルをアクセスする場合、ファイルの各ブロックは、ディスク130からバッファキャッシュ170へロードされる場合がある。

30

40

【0042】

オンディスクinode(又は、ブロック)が、ディスク130からバッファキャッシュ170へロードされる時、それに対応するコア内構造が、オンディスク構造に埋め込まれる。例えば、inode300の周りを囲む点線は、オンディスクinode構造のコア内表現を示している。コア内構造は、オンディスク構造、及びメモリ上の(ただし、ディスク上に無い)データの管理に必要な他の情報を記憶する一つのメモリブロックである。他の情報には、例えば、「ダーティ」ビット360がある。例えば、書き込みオペレーションによる命令に従って、inode(又はブロック)内のデータが更新/変更された後、変更されたデータは、ダーティビット360を使用して、「汚れた」とし

50

てマーキングされ、その後、その `inode` (ブロック) はディスクに「フラッシュ」(記憶) できるようになる。 `inode`、及び `inode` ファイルを含む、 `W A F L` ファイルシステムのコア内構造、及びオンディスクフォーマット構造については、1998年10月6日に発行された「Method for Maintaining Consistent States of a File System and for Creating User-Accessible Read-Only Copies of a File System」と題するDavid Hitz他による米国特許第5,819,292号に開示、及び記載されており、この文献は、参照により本明細書に援用される。

【0043】

図4は、本発明とともに有利に使用されるデータコンテナのバッファツリーを示す略ブロック図である。バッファツリーは、バッファキャッシュ170にロードされ、 `w r i t e - a n y w h e r e` ファイルシステム280によって管理されるデータコンテナ(例えばファイル400)の種々のブロックの内部表現である。埋め込み `inode` のようなルート(最上位) `inode` 300は、間接(例えば、第1レベル)ブロック404を参照する。間接ブロック(及び、 `inode`) は、ファイル400の実際のデータの記憶に使用されるデータブロック406を最終的に最終的に参照するポインタ405を含む。すなわち、ファイル400のデータは、種々のデータブロックに記憶され、それらのブロックの位置が、ファイルの間接ブロックに記憶される。各レベル1間接ブロック404は、1024個もの数のデータブロックへのポインタを有する可能性がある。ファイルシステムの「 `w r i t e a n y w h e r e` 」な性質によれば、これらのブロックは、ディスク130上のどこに配置されることもありうる。

【0044】

D. データ重複解消

図5は、本発明の一実施形態による、データブロック406への複数のポインタを含むレベル1間接ブロック404を示す略ブロック図である。例えば、各データブロック406は、4KBのデータを含む。特に、 `w r i t e a n y w h e r e` ファイルレイアウトによれば、レベル0データブロックは、ファイルシステム内のどこに配置されることもあり得る(すなわち、それらが、ディスク130上の物理的に連続したブロックに対応している必要はない)。データを削除、及び/又は重複解消するように構成されたストレージオペレーティングシステム200では、各データブロック406内のデータに、指紋を関連付ける場合がある。例えば、例示的なデータコンテナは、例えば `v b n 1`、 `v b n 2`、 `v b n 3`、及び `v b n 4` に記憶された一連のデータブロックを含む。各一意のデータブロックには、一意の指紋、例えば、A、B、C、及びDが関連付けられる。同様に、間接ブロック404の中にある一連のポインタ405、例えば、P1、P2、P3、及びP4は、データブロック `v b n 1`、 `v b n 2`、 `v b n 3`、及び `v b n 4` をそれぞれ参照する。

【0045】

図6は、書き込みオペレーションによってデータコンテナに変更が加えられた後の、図5のレベル1間接ブロック404を示す略ブロック図である。例えば、 `v b n 3`、及び `v b n 4` に以前に記憶された2つのデータブロックは、現在は、 `v b n 1` に記憶された第1のデータブロックのコピーを有するように変更され、その結果、それらに関連する指紋は同一になっている。変更が加えられたデータブロックは、今度は、指紋A、B、A、及びAにそれぞれ関連している。 `w r i t e a n y w h e r e` ファイルレイアウトによれば、変更の加えられたデータには、2つの新たな `v b n` (`v b n 5` と `v b n 6`) が割り当てられる。その結果、重複解消モジュール284は、それらの同一の指紋を使用して、 `v b n 5`、及び `v b n 6` に記憶されているデータブロック406が、 `v b n 1` にあるデータブロック406に記憶されているデータの複製を有しているか否かを判定する場合がある。同一データの3つのコピーを有することは、ファイルシステムリソースの浪費である。従って、重複解消モジュール284は、ポインタP3、及びP4を `v b n 1` を参照するように変更し、それによって、データコンテナの完全性を維持しつつも、 `v b n 5`、及び `v b n 6` の割り当てが解除され、記憶空間は節約される。

【0046】

図7は、例示的なデータ重複解消手順の実施後の、図6のレベル1間接ブロック404を示す略ブロック図である。例示の手順の後、各一意のデータブロックは、一意の指紋に関連し、データコンテナの複数の同一のファイルブロックが、一意のデータブロックに関連し、さらに、その一意のデータブロックが、一意の指紋に関連する。重複解消手順を実施する技術の詳細については、「Method and Apparatus for Identifying and Eliminating Duplicate Data Blocks and Sharing Data Blocks in a Storage System」と題する、上で援用した米国特許出願に記載されている。

【0047】

E. 指紋レコードの生成

本発明の一実施形態によれば、ストレージシステム120は、ファイルシステム280のような上位モジュールを実施し、RAIDシステム240のようなストレージモジュールによって管理されるディスク130のような複数の記憶装置に記憶される情報を論理編成する。上位モジュールとストレージモジュールの間の境界において、指紋抽出層286は、ストレージシステム120によって記憶される各データブロック406について、例えば64ビットのような所定サイズの指紋を生成する。例えば、指紋の最初の所定数のビット、例えば、最初の32ビットは、ストレージモジュールによって実施されるチェックサム計算の結果から得られ、指紋の次の所定数のビット、例えば、次の32ビットは、データブロック406中の1以上の所定のオフセットに記憶されたデータから得られる。得られた指紋は、データブロック406の中身を十分に反映し、望ましくない誤認識率を生じることなく、重複の識別を可能にする。また、指紋の生成のために、ハッシュ関数のような他の暗号化関数をさらに実施する必要はないため、プロセッサ使用率も低減される。

【0048】

さらに、指紋抽出層286は、特定のポリシーに基づいて、データブロックのサブセットに対して指紋を生成するように選択される場合もある。例えば、ユーザファイルに属するデータブロックに対して(ファイルシステムによって使用される内部的なメタデータファイルとは対照的に)指紋が生成される場合がある。あるいは、レベル0ブロック406に対して(間接ブロック404とは対称的に)、指紋が生成される場合がある。

【0049】

各指紋は指紋レコードに記憶され、指紋レコードは更に、指紋データベース288に記憶され、データ重複解消モジュール284によってアクセスされる。データ重複解消モジュール284は、指紋を定期的に比較し、重複する指紋を識別し、重複データブロック406を知らせる場合がある。次に、重複データブロック406が比較され、重複データは削除される。

【0050】

本発明の一実施形態によれば、指紋抽出層286は、ストレージシステム120によって記憶される各データブロック406に対して64ビットの指紋を生成する。図8は、例示的な指紋800の種々の要素を示す略ブロック図である。指紋800は例えば、RAIDチェックサムフィールド805、及びサンプルデータフィールド810を含み、代替実施形態では、さらに別のフィールド815を有する場合がある。なお、代替実施形態では、更に別のフィールド、及び/又は異なるフィールドが使用される場合もある。あるいは、RAIDチェックサムが、重複解消モジュールのニーズに十分に合致するチェックサム関数を含む場合、指紋800は、RAIDチェックサムだけしか有しない場合もある。

【0051】

通常書き込みオペレーションの一部として、ディスクに書き込まれた各データブロックに対するチェックサム計算を実施することにより、周知のRAIDプロトコルによってデータの完全性が確認される。例えば、データブロック406に対する各書き込みオペレーションについて、指紋抽出層286は、RAIDプロトコルに従って計算された32ビットのチェックサム値を識別し、それを指紋800のRAIDチェックサムフィールド805にコピーする。従って、指紋を生成するプロセスは、各データブロック406に関連する事前計算されたチェックサム値の利点が得られるように最適化される。あるいは、R

10

20

30

40

50

RAIDチェックサムフィールド805は、RAIDプロトコルに従って計算された32ビットチェックサムの一部だけを単独で、又は他の値と組み合わせる場合がある。同様に、ストレージプロトコルが、32ビットRAIDチェックサム値以外のチェックサム値を生成する場合、指紋800は、RAIDチェックサムフィールド805において、そのような他のチェックサム値を実施する場合がある。あるいは、RAIDチェックサムフィールド805は、データブロック406に記憶されたデータに関連する、ハッシュ関数のような他の暗号化関数の結果を含む場合がある。

【0052】

例えば、サンプルデータフィールド810は、指紋800に関連するデータブロック406内の1以上の所定のオフセットに記憶されたデータのコピーを有する。例えば、データブロック406が4KBのデータを有している場合、サンプルデータフィールド810は、32ビット(4バイト)の長さを有し、データブロック406に記憶されたデータの1024ビット目ごと(すなわち、128バイト目ごとにその最初のビット)のコピーを有する場合がある。あるいは、サンプルデータフィールド810は、データブロック406内の単一オフセットに記憶された32個の連続したビットのコピーを有する場合がある。得られる指紋800は、データブロック406の中身を十分に反映し、望ましくない誤認識率を生じることなく、重複の識別を可能にする。したがって、ここでも、指紋を生成するプロセスは、ストレージオペレーティングシステム200によって更に別のハッシュ計算、及び/又は暗号化関数を実施する必要なく、ストレージオペレーティングシステム200に既に提供されたデータの利点が見られるように最適化される。代替実施形態によれば、RAIDチェックサムフィールド805、及びサンプルデータフィールド810は、例えば、ストレージオペレーティングシステム200によって実施される更に別のハッシュ計算、及び/又は暗号化関数の結果といった、データブロック406に関連する他のデータを含むさらに別のフィールド815と組み合わせられる場合がある。

【0053】

各指紋は、データ重複解消モジュール284によって生成された指紋レコードに記憶される。図9は、本発明とともに有利に使用される指紋レコード900を示す略ブロック図である。指紋レコード900は、指紋800を記憶するフィールド、inode番号フィールド905、ファイルブロック番号(fbn)フィールド910、生成番号フィールド915、及びコンシステンシ・ポイント(CP)カウントフィールド920を含み、代替実施形態では、更に別のフィールド925を含む場合もある。inode番号フィールド905は、指紋が生成されたデータブロック406のinode300へのポインタを記憶する。fbnフィールド910は、データコンテナ内のブロックのオフセットを記憶する。例えば、データコンテナが複数のデータブロック406を含む場合、fbnフィールド910の値は、どのデータブロック406が指紋800に対応するものかを識別する。指紋レコード900の他のフィールド915、920、及び925は、書き込みオペレーションの際に、そのブロックのinode300のメタデータ部から収集される場合がある。例えば、inode300のCPカウントフィールド322、及び生成番号フィールド324に記憶される値は、指紋レコード900のCPカウントフィールド920、及び生成番号フィールド915からそれぞれコピーされる場合がある。従って、指紋レコード900は、指紋800に関連するデータを記憶し、更なる指紋処理機能をストレージオペレーティングシステムに与える。

【0054】

図10は、本発明の一実施形態による、指紋レコードを生成する手順1000のステップの詳細を示すフロー図である。手順1000は、ステップ1005から開始され、ステップ1010へ進み、そこで、ストレージオペレーティングシステム200のファイルシステム280は、書き込みコマンドをRAIDシステム240へ発行する。例えば、この書き込みコマンドは、ストレージシステム120に接続されたクライアント110によって生成される、新たなデータコンテナの新たなデータブロック406を書き込むためのコマンドである場合がある。あるいは、書き込みコマンドは、write anywhere

10

20

30

40

50

e ファイルレイアウトに従って新たなデータブロック 406 を割り当てることにより、既存のデータコンテナのデータを変更するコマンドである場合もある。

【0055】

手順 1000 はステップ 1020 へ進み、そこで RAID システムは、32 ビットチェックサム値を計算し、例えば、書き込みオペレーションをディスクドライバシステム 250 を通じて送信することにより、書き込みオペレーションを完了する。ステップ 1025 において、RAID システム 240 は、次に、確認メッセージをファイルシステム 280 へ返す。なお、書き込みオペレーションを実施し、チェックサム値のような値を返すために、他のストレージプロトコル、すなわち、他の RAID プロトコルを実施してもよい。ステップ 1030 において、指紋抽出層 286 は、RAID システム 240 からの確認メッセージを捕捉し、及び/又は監視する。ステップ 1035 において、指紋抽出層 286 は、確認メッセージから 32 ビットチェックサム値を抽出し、その値を指紋 800 の RAID チェックサムフィールド 805 にコピーする。また、指紋抽出層 286 は、データブロック 406 の種々のオフセットから 32 ビットのデータをサンプリングし、そのデータを指紋 800 のサンプルデータフィールド 810 にコピーする。従って、64 ビットの指紋 800 が生成される。さらに、指紋抽出層 286 は、データブロック 406 に関連するメタデータを抽出する。指紋レコード 900 が生成され（例えば指紋レコード 900 の種々のフィールドを記憶するための 1 以上のメモリエントリを割り当てることにより）、抽出されたメタデータは、指紋レコード 900 の適当なフィールドにコピーされる。代替実施形態によれば、指紋抽出層 286 は、書き込みコマンドに回答して、ファイルシステム 280 が提供するデータに基いてそれらのフィールドのデータを生成する場合がある。

10

20

【0056】

手順 1000 はステップ 1040 へと進み、そこで、指紋 800 は指紋レコード 900 に追加される。ステップ 1045 において、指紋レコード 900 は、指紋キュー 287 に追加され、そこで、新たに生成された指紋は、図 11 を参照して後で説明するような、重複解除モジュール 284 によって実施されるバッチ処理に備えて蓄積される。ステップ 1050 において、手順 1000 は終了する。

【0057】

なお、代替実施形態によれば、指紋レコード生成手順 1000 は、定期的実施される場合もあれば、書き込みオペレーション中に実施されるのではなく、例えば UI 275 から受信される管理者コマンドのような、所定の条件に回答して実施される場合もある。そのような実施形態では、ファイルシステム 280 は、ディスクアレイ 160、又はその一部をスキャンし、指紋レコードを生成する場合がある。

30

【0058】

例えば、データ重複解除モジュール 284 は、重複する指紋を識別し、重複するデータブロックを示すために、定期的指紋を比較する場合がある。なお、例示的实施形態によれば、重複する一対の指紋は、「誤認識」であることもある。「誤認識」とは、それらの指紋が、実際の重複データブロックには該当しないことを意味する。従って、重複解除を行う前に、識別されたデータブロックをバイトごとに比較し、それらが実際に重複であるか否かを確認する。例えば、各指紋 800 が、32 ビットの RAID チェックサム、及び 32 ビットのサンプルデータを有している場合、指紋レコード生成手順 1000 におけるストレージシステム処理リソースは、データ重複解除の際の誤認識の確率と引き換えに、節約される。ただし、誤認識の確率は非常に低く（40 億個のデータブロックにつき約 100 個の誤認識）、事前計算されたデータを使用して指紋 800 を生成することにより節約されるプロセッサリソースの量は、データ重複解除の際にバイトごとの比較を実施することにより消費されるプロセッサリソースの量を上回る。さらに、重複解除プロセスは、システムの活動が活発でない時、例えば深夜、又は週末に実施される場合がある。したがって、本明細書では、「同一」の、又は「重複」するデータブロックは、低い確率で実際には重複していないデータブロック（誤認識）も含む。

40

【0059】

50

一方、代替実施形態によれば、各指紋 800 において、RAID チェックサムフィールドとサンプルデータフィールド 805、810 に異なる値がコピーされる場合があり、及び/又は、更に別のフィールド 815 が実施される場合がある。例えば、指紋 800 は、誤認識の確率が確実にゼロになるようにするために暗号ハッシュ関数の結果を含む場合がある。従って、代替実施形態によれば、データ重複解消の際にバイト単位のデータブロック比較を実施する必要は無い場合もある。

【0060】

F. 重複データブロックの識別

図 11 は、本発明の一実施形態による重複データブロック 406 を識別する手順 1100 のステップの詳細を示すフロー図である。手順 1100 は、ステップ 1105 から開始され、手順 1110 へ進み、そこで、指紋キュー 287 に記憶されている指紋を指紋データベース 288 に記憶されている指紋と比較する。もし、同一の指紋があれば、手順 1100 はステップ 1115 へ進み、そこで、重複解消モジュール 284 は、例えば、同一の指紋に関連するデータブロック間でバイトごとの比較を実施し、その後、間接ブロック中のポインタを変更することにより、データ重複解消を実施する。次に、ステップ 1120 において、指紋キュー 287 に記憶されている指紋を指紋キュー 287 に記憶されている他の指紋と比較する。例えば、この比較は、2つの指紋間におけるビットごとの比較であり、重複解消モジュール 284 によって実施される。代替実施形態によれば、この比較は、例えば、バイトごとの比較であってもよく、及び/又は同一データを識別する他の方法であってもよい。同一の指紋があれば、手順 1100 はステップ 1125 へ進み、そこで、重複解消モジュール 284 は、図 5、図 6、及び図 7 を参照して上で例示したものと同様に、データ重複解消を実施する。次に、手順 1100 はステップ 1130 へ進み、そこで、指紋キュー 287 に記憶されている指紋レコードを指紋データベース 288 にコピーする。次に、ステップ 1135 において、指紋キュー 287 をフラッシュし、すなわち、指紋キュー 287 内の全ての指紋レコードを削除する。手順 1100 はステップ 1140 で終了する。

【0061】

なお、例示的实施形態によれば、上記手順 1100 は、アクティブファイルシステムに対して実施される場合がある。従って、ストレージオペレーティングシステム 200 には、2つ以上の指紋キューが存在する場合がある。手順 1100 のにおいて重複解消モジュール 284 がキューの一つをアクセスしている間も、書き込みオペレーションは継続され、新たに生成された指紋レコードは、第 2 の指紋キューに記憶されるか、又は、一時的記憶場所に記憶される場合がある。新たに生成された指紋レコードを指紋データベース 288 に記憶するのではなく、指紋キュー 287 に記憶することにより、新たに書き込まれたデータブロック、又は最近変更が加えられたデータブロックに対応する指紋を一つのグループにまとめて「バッチ」処理することができる。あるいは、新たに生成された指紋レコードは、直ぐに比較し、処理するために、指紋データベース 288 に直接記憶される場合もある。また、指紋データベース 288 の一端に指紋キューを有効に生成するために、指紋データベース 288 は、指紋生成の順番に従ってさらにソートされる場合があり、ソートされた指紋データベース 288 は、独立した指紋キュー 287 を必要とせずに、手順 1100 によってバッチ処理される場合がある。さらに別の実施形態によれば、新たに生成された指紋レコードは、指紋データベース 288 に直接記憶され、指紋値 800 によってソートされる場合があり、重複解消モジュール 284 は、指紋データベース 288 の全体、又は種々の部分を定期的にスキャンすることにより、重複を調べる場合がある。

【0062】

例示的实施形態によれば、指紋 800 間の比較を容易にするために、指紋キュー 287、及び指紋データベース 288 の指紋レコード 900 は、指紋レコードが有している指紋 800 の値に基いてソートされる場合がある。例えば、図 12A は、ソートされた指紋キュー 287、及びソートされた指紋データベース 288 を示している。例えば、指紋 800 は文字 A ~ G で表され、指紋キュー 287 内の種々の場所は Q1 ~ Q6 で表され、指紋

データベース 288 内の種々の場所は D1 ~ D4 で表されている。

【0063】

上記手順 1100 のステップ 1110、及び 1115 では、指紋キュー 287 と指紋データベース 288 の両方にある同一の指紋が識別され、それに対応するデータブロックの重複が解消される。例えば、図 12B は、識別された重複する指紋（「C」、及び「D」）を示している。記憶場所 Q3、及び D1 に記憶されている指紋レコードに関連するデータブロックを比較し、それらが、実際の重複データブロックであるか否かが、判定される。同様に、記憶場所 Q4、及び D2 に記憶されている指紋レコードに関連するデータブロックが、比較される場合がある。このような比較を実施し、実際の重複データブロックの重複を解消する技術は、「Method and Apparatus for Identifying and Eliminating Duplicate Data Blocks and Sharing Data Blocks in a Storage System」と題する上記の米国特許出願に開示されている。

10

【0064】

上記手順 1100 のステップ 1120、及び 1125 では、指紋キュー 287 内の同一の指紋が識別され、それらに対応するデータブロックの重複が解消される。例えば、図 12C は、指紋キュー 287 内で識別された重複する指紋（「G」）を示している。記憶場所 Q5、及び Q6 に記憶されている指紋レコードに関連するデータブロックを比較し、それらが実際の重複データブロックであるか否かが、判定される場合がある。

【0065】

実際の重複データブロック、及びそれに対応する指紋レコードは削除され、上記手順 1100 のステップ 1130 において、指紋キュー 287 は指紋データベース 288 に追加される。結果得られるマージ後の指紋データベース 288 は、指紋値 800 に従ってソートされる。例示的实施形態によれば、誤認識の確率がある場合、誤認識された指紋を含む指紋レコードは、保持される場合がある。例えば、図 12D は、手順 1100 の完了後の、マージされ、記憶された指紋データベース 288 を示している。例えば、指紋「C」、及び「G」は実際の重複データブロックに相当するが、指紋「D」は誤認識であった。従って、指紋「C」と「G」の重複は解消され、それらに対応するデータブロックの重複は解消された。一方、指紋「D」の重複は保持され、指紋データベース 288 にマージされた。もし、指紋「D」を有する新たなデータブロックが生成された場合、そのデータブロックは、指紋「D」を有する既存のデータブロックのいずれか、又は両方と比較することが望ましい。

20

30

【0066】

G . 古くなった指紋の削除

例示的实施形態によれば、各指紋レコード 900 は、指紋が生成されたデータブロック 406 に関する更に別の情報を有する場合がある。データ重複解消モジュール 284 は、その別の情報を、データブロック 406 に関連するメタデータと定期的に比較し、指紋データベース 288 が、ストレージシステムの最新の状態を確実に反映するようにする。そのため、「古くなった」指紋（例えば、削除されたデータブロック、及び/又は変更されたデータブロックに関連する指紋）は、指紋データベース 288 から削除される場合がある。あるいは、古くなった指紋を削除するプロセスは、例えば、指紋データベース 288 が所定のサイズに達するといった所定の条件にตอบสนองして実施され、あるいは、古くなった指紋を削除するプロセスは、指紋を比較する上記プロセスの直前に実施される場合がある。あるいは、指紋は、指紋が古くなったとき、すなわち、対応するデータが削除され、又は変更されたときに削除される場合がある。

40

【0067】

一実施形態によれば、RAID システム 240 に対して書き込みオペレーションが実施されるたびに、新たな指紋レコード 900 が生成される。従って、同じデータブロックに対して複数の書き込みオペレーションが実施された場合、複数の指紋レコードが生成されることもある。ただし、そのデータブロックを正確に反映するのは、最新の指紋レコードだけである。同様に、データブロックは、対応する指紋レコードが陳腐化したときに、削

50

除され、及び/又は再割り当てされる場合がある。従って、処理能力を向上させ、及び/又は指紋データベース288によって消費される空間を減らすためには、「古くなった」指紋レコード(例えば、削除されたデータブロック、及び/又は変更されたデータブロックに関連する指紋レコード)を、指紋データベース288から削除する必要がある。

【0068】

図13は、本発明の一実施形態による、古くなった指紋レコードを削除する手順1300のステップを詳細に示すフロー図である。例えば、手順1300は、重複解消モジュール284によって実行されるが、手順1300は、ストレージオペレーティングシステム200との間で情報をやりとりする機能を有するものであれば、他のいかなるモジュールによって実行されてもよい。手順1300は、ステップ1305から開始され、ステップ1310へと進み、そこで、重複解消モジュール284は、指紋レコード900に記憶されているinode番号905を読み出す(ロードする)。次に、重複解消モジュール284は、そのinode番号905を使用して、データブロック406のinode300をロードする。複数の指紋レコードのinodeをロードしつつ、シーケンシャルディスクアクセスを可能にするため、手順1300の前に、又は手順1300において、指紋データベース288は、inode番号905、及びfbn910に従ってソートされる場合がある。ステップ1315において、重複解消モジュール284は、inode300からメタデータを取り出し(読み出し)、指紋レコード900に記憶されているメタデータと照合する。

【0069】

例えば、ステップ1320では、生成番号324、915を比較し、指紋レコード900の生成後に、inode300が削除され、再利用された(すなわち、異なっている)か否かが判定される。そうであれば、指紋レコード900は古くなっているため、手順1300はステップ1340へ進み、そこで、指紋レコード900は削除される場合がある。そうでなければ、手順1300はステップ1325へ進み、そこで、CPカウント322、920を比較し、指紋レコード900の生成後に、inode300が変更されたか否か(その結果、複数のCPカウントが発生する)が判定される。そうであれば、指紋レコード900は古くなっており、このinode300が複数回変更されていれば、inode300に対応する他の指紋レコードも、同様に古くなっている場合がある。従って、ステップ1330では、このinode300に対して最新の指紋レコード900だけを保持し、他の指紋レコード(低いCPカウントを有する)は全て削除する場合がある。複数の指紋レコードのCPカウント間の比較を容易にするために、指紋データベース288は、まず、inode番号905によってソートされ、次に、CPカウント920によってソートされる場合がある(複数の指紋レコードが同一のinode番号を共有している場合)。

【0070】

さらに、データコンテナは、切り詰められる場合がある(例えば、データコンテナ中の所与のオフセット以降のデータを全て削除することにより)。データコンテナの切り詰めは、関連inode300に記憶されるメタデータに影響を与えないことがある。従って、ステップ1335において、手順1300は、指紋800のサンプルデータ810を使用して、データブロック406との間の整合性(すなわち、キューオフセットにおいてデータが相違するか否か)を確認する。サンプルデータ810(データブロック406内の1以上の所与のオフセットから最初にコピーされたもの)が、データブロック406のその所与のオフセットに現在記憶されているデータと相違する場合、手順1300はステップ1340へと進み、そこで、指紋レコード900を削除する場合がある。先行するデータブロックにあるデータコンテナも、切り詰められる場合がある(例えば、データコンテナに関連するinode300のサイズフィールドに記憶される値を減らすことにより)。従って、ステップ1350において、手順1300は、指紋レコード900のfbn910を関連inode300のサイズ314と比較する。サイズ314がfbn910よりも小さかった場合、手順1300はステップ1340へと進み、そこで指紋レコード9

10

20

30

40

50

00は削除される場合がある。そして、手順1300はステップ1345において終了する。

【0071】

例えば、古くなった指紋の削除手順1300は、定期的実施され、及び/又は所定の条件(例えば、UI245から受信される管理者コマンド等)に応じて実施される。また、手順1300は、指紋データベース288が所定の閾値サイズを超えたときに、それに応じて実施される場合もある。例えば、複数のデータブロックが繰り返し書き込まれ、及び/又は変更される際に、指紋データベース288における指紋レコードの数は、ストレージシステムにおけるデータブロックの最大数を上回ることがある。なぜなら、指紋レコードの多くが、古くなっているからである。手順1300は、例えば、指紋レコードの数がデータブロックの最大数の120%を超えたときに実施される場合がある。

10

【0072】

代替実施形態によれば、古くなった指紋の削除手順1300は、重複指紋識別・データ重複解消手順1100の一部として実施される場合もある。例えば、手順1300は、手順1100のステップ1110の直前に、重複解消モジュール284によって実施される場合がある。その結果、古くなった指紋、及びそれに関連するデータブロックを比較する際に消費されるストレージシステム処理リソースの量を減らせる場合がある。

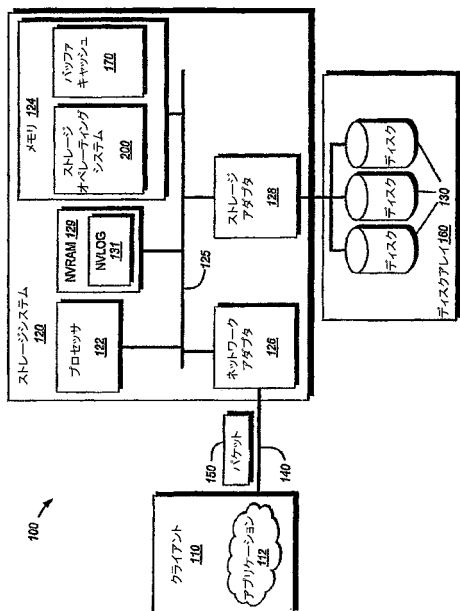
【0073】

上記の説明は、本発明の特定の幾つかの実施形態に関するものである。しかしながら、当業者には明らかなように、それらの実施形態の利点の一部、又は全部を維持しながらも、記載した実施形態に対し、他にも変形、又は変更を施すことが可能である。当然ながら、本明細書に記載する手順、プロセス、及び方法は、他の順序で実施される場合もある。例えば、手順1100の指紋キューと指紋データベースの比較ステップ1120は、指紋キューと指紋キューの比較ステップ1110の後に実施される場合や、それと同時に実施される場合もある。同様に、指紋生成、重複識別、及び古くなった指紋の削除等のプロセス/手順は、同時に実施される場合もあれば、個別に(すなわち、任意の順序で)実施される場合もあれば、プロセス/手順の所定の順番で実施される場合もある。また、本発明の教示は、ソフトウェア(コンピュータ上で実行されるプログラム命令を含むコンピュータ読取可能媒体を含む)でも、ハードウェアでも、ファームウェアでも実施することができ、あるいはそれらの組み合わせによって実施される場合もある。本明細書の説明は、ファイルシステムを例として書かれているが、本発明は、ファイルシステム以外のストレージ、例えば、LUN、及び/又はブロックベースのストレージにも使用することが可能である。従って、本明細書の説明は、単なる例として捉えるべきものであり、本発明の範囲を何ら制限するものではない。従って、添付の特許請求の範囲の目的は、そうした変形や変更も、本発明の思想、及び範囲に含めることにある。

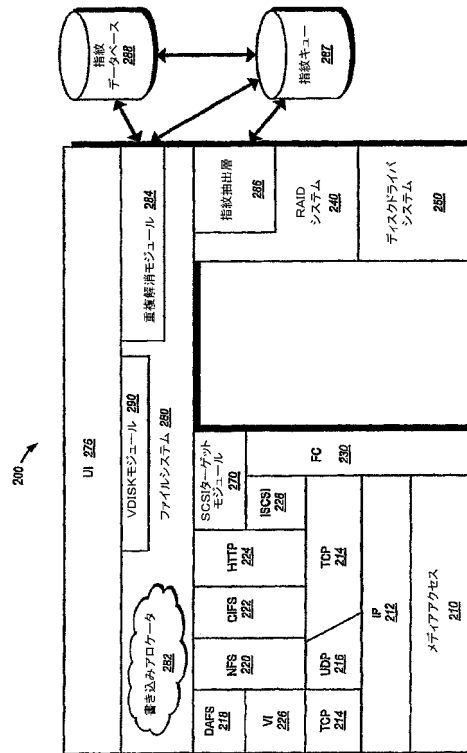
20

30

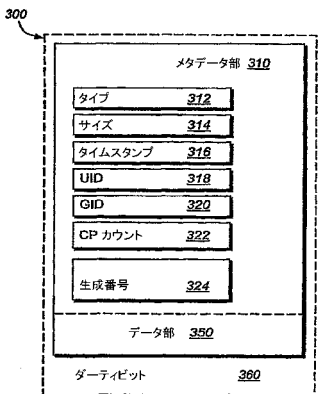
【図1】



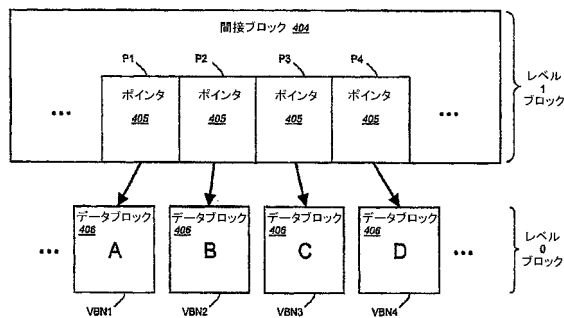
【図2】



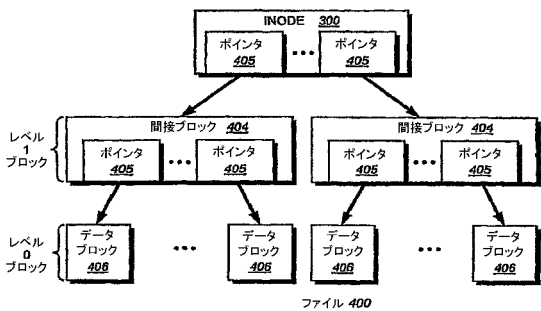
【図3】



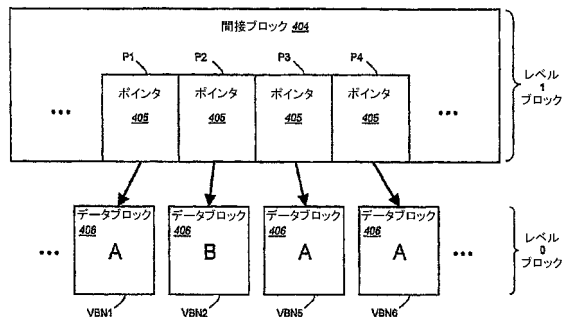
【図5】



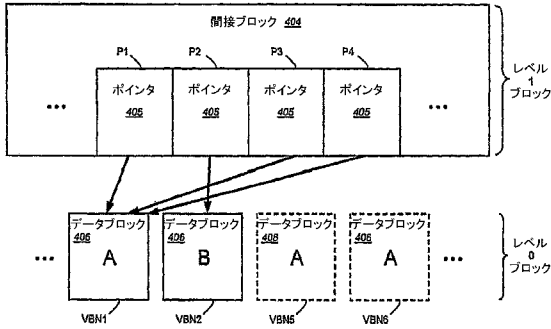
【図4】



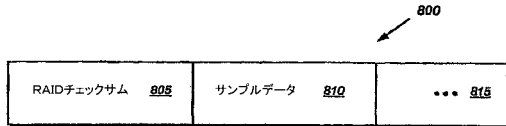
【図6】



【図7】



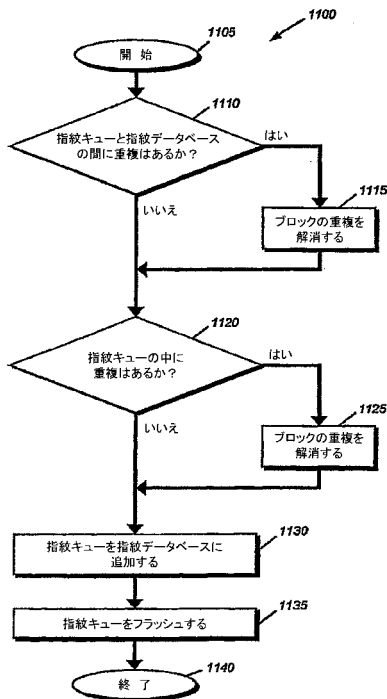
【図8】



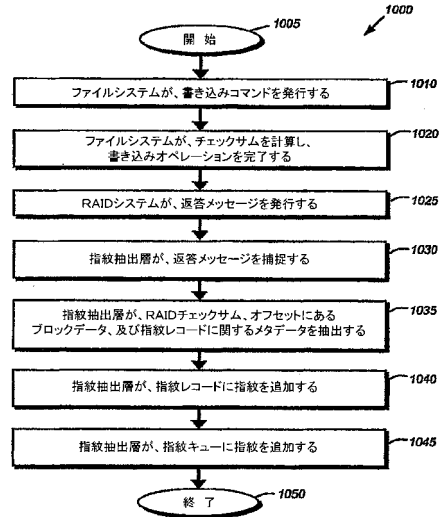
【図9】

指紋	800
INODE番号	905
FBN	910
生成番号	915
CP カウント	920
⋮	925

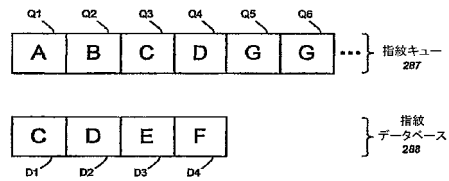
【図11】



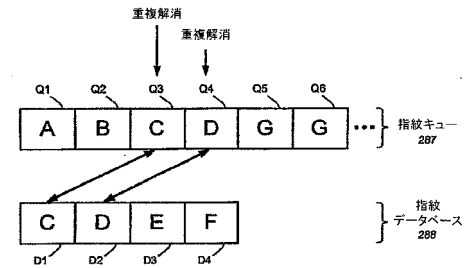
【図10】



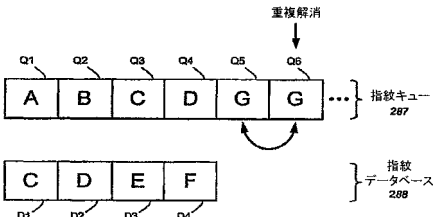
【図12A】



【図12B】



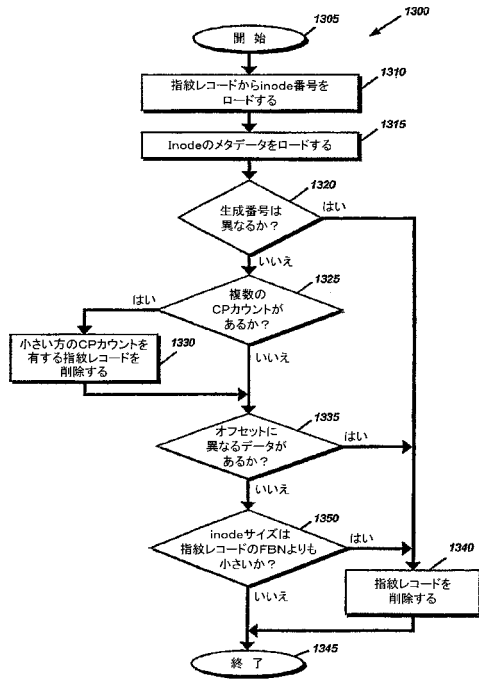
【図12C】



【図12D】



【図13】



フロントページの続き

(72)発明者 フォン, インフン

アメリカ合衆国カリフォルニア州 9 4 0 8 9 , サニーヴェール, イースト・ジャバ・ドライブ・4
9 5

審査官 桜井 茂行

(56)参考文献 特表 2 0 0 3 - 5 2 4 2 4 3 (J P , A)

国際公開第 0 1 / 0 6 1 5 6 3 (W O , A 1)

Bo Hong et al, "Duplicate Data Elimination in a SAN File System", Proceedings of the 2
1st IEEE / 12th NASA Goddard conference on mass storage system and technologies, 米国
, IEEE, 2 0 0 4年 4月, P. 103, [Online] [http://www.ssrc.ucsc.edu/Papers/hong-mss04.
pdf](http://www.ssrc.ucsc.edu/Papers/hong-mss04.pdf) [2011/11/2取得]

(58)調査した分野(Int.Cl., DB名)

G 0 6 F 1 2 / 0 0

G 0 6 F 3 / 0 6