

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6947670号
(P6947670)

(45) 発行日 令和3年10月13日 (2021. 10. 13)

(24) 登録日 令和3年9月21日 (2021. 9. 21)

(51) Int. Cl.	F I
G 0 6 F 3/06 (2006.01)	G 0 6 F 3/06 3 0 5 Z
G 0 6 F 3/08 (2006.01)	G 0 6 F 3/06 3 0 5 C
	G 0 6 F 3/06 3 0 4 F
	G 0 6 F 3/08 H
	G 0 6 F 3/06 3 0 4 N

請求項の数 20 (全 19 頁)

(21) 出願番号	特願2018-52379 (P2018-52379)	(73) 特許権者	390019839
(22) 出願日	平成30年3月20日 (2018. 3. 20)		三星電子株式会社
(65) 公開番号	特開2018-156656 (P2018-156656A)		S a m s u n g E l e c t r o n i c s
(43) 公開日	平成30年10月4日 (2018. 10. 4)		C o . , L t d .
審査請求日	令和3年2月17日 (2021. 2. 17)		大韓民国京畿道水原市靈通区三星路129
(31) 優先権主張番号	62/474, 039		129, S a m s u n g - r o , Y e o n
(32) 優先日	平成29年3月20日 (2017. 3. 20)		g t o n g - g u , S u w o n - s i , G
(33) 優先権主張国・地域又は機関	米国 (US)		y e o n g g i - d o , R e p u b l i c
(31) 優先権主張番号	62/561, 625	(74) 代理人	110000051
(32) 優先日	平成29年9月21日 (2017. 9. 21)		特許業務法人共生国際特許事務所
(33) 優先権主張国・地域又は機関	米国 (US)	(72) 発明者	キ ヤ ン ソク
			アメリカ合衆国, 94303, カリフ
			ォルニア州, パロ アルト, アルテア
			ウォーク 873
			最終頁に続く

(54) 【発明の名称】 仮想装置階層を利用した複数のメモリ装置を含む仮想装置に対する客体の格納及び読み取り方法とこれを用いたストレージ装置

(57) 【特許請求の範囲】

【請求項 1】

各々が許容可能な最小サイズの値を有し、ステートレスのデータ保護を利用する仮想装置として構成された複数からなるメモリ装置と、

第1客体及び第2客体のそれぞれのサイズに応じて、前記第1客体及び前記第2客体の各々に第1データ保護及び第2データ保護の中の1つを適用することにより、各々がそれぞれのサイズを有する前記第1客体及び前記第2客体を格納するように、前記仮想装置を管理するように構成された仮想装置階層と、を備え、

前記仮想装置階層は、前記第1データ保護を大型客体に適用し、

前記第2データ保護を小型客体に適用し、

前記第1データ保護又は前記第2データ保護のいずれかを中型客体に適用するように構成されることを特徴とするストレージ装置。

【請求項 2】

前記メモリ装置は、1つ以上のデータ装置及び1つ以上のパリティ装置として構成されることを特徴とする請求項1に記載のストレージ装置。

【請求項 3】

前記第1データ保護は、消去コーディングを含み、

前記第2データ保護は、複製を含むことを特徴とする請求項2に記載のストレージ装置。

【請求項 4】

前記消去コーディングは、前記第 1 客体及び前記第 2 客体の中の該当客体が大型客体として分類される場合、データ保護のために利用されることを特徴とする請求項 3 に記載のストレージ装置。

【請求項 5】

前記第 1 客体及び前記第 2 客体の中の該当客体は、以下の不等式を満足する場合、前記大型客体として分類されることを特徴とする請求項 4 に記載のストレージ装置。

$$((P + 1) \times O > (S + P) \times m \text{ AND } O \geq S \times m)$$

(ここで、O は客体のサイズ、P はパリティ装置の個数、S はデータ装置の個数、m は前記複数からなるメモリ装置のそれぞれの許容可能な最小サイズの値の中の最大値を示す。)

10

【請求項 6】

前記複製は、前記第 1 客体及び前記第 2 客体の中の該当客体が小型客体として分類される場合、データ保護のために利用されることを特徴とする請求項 3 に記載のストレージ装置。

【請求項 7】

前記第 1 客体及び前記第 2 客体の中の該当客体は、以下の不等式を満足する場合、前記小型客体として分類されることを特徴とする請求項 6 に記載のストレージ装置。

$$((P + 1) \times O = < (S + P) \times m)$$

(ここで、O は客体のサイズ、P はパリティ装置の個数、S はデータ装置の個数、m は前記複数からなるメモリ装置のそれぞれの許容可能な最小サイズの値の中の最大値を示す。)

20

【請求項 8】

前記消去コーディング又は前記複製のいずれかは、前記第 1 客体及び前記第 2 客体の中の該当客体が大型客体でも小型客体でもなく中型客体として分類される場合、前記複数からなるメモリ装置で利用可能な 1 つ以上の空間、前記仮想装置に格納された前記第 1 客体又は前記第 2 客体に対するアクセス時間、及び前記第 1 客体又は前記第 2 客体がアクセスされる頻度に基づいて、データ保護のために利用されることを特徴とする請求項 3 に記載のストレージ装置。

【請求項 9】

前記第 1 客体及び前記第 2 客体の中の該当客体は、以下の不等式を満足する場合、前記中型客体として分類されることを特徴とする請求項 8 に記載のストレージ装置。

$$((P + 1) \times O > (S + P) \times m) > S \times m > O)$$

(ここで、O は客体のサイズ、P はパリティ装置の個数、S はデータ装置の個数、m は前記複数からなるメモリ装置のそれぞれの許容可能な最小サイズの値の中の最大値を示す。)

30

【請求項 10】

前記パリティ装置は、前記メモリ装置の固定されたサブセットであることを特徴とする請求項 2 に記載のストレージ装置。

【請求項 11】

前記パリティ装置は、前記メモリ装置の可変されるサブセットを含み、

40

前記複数からなるメモリ装置の各々は、データ装置及びパリティ装置の両方として動作するように構成されることを特徴とする請求項 2 に記載のストレージ装置。

【請求項 12】

前記メモリ装置は、ソリッドステートドライブを含むことを特徴とする請求項 1 に記載のストレージ装置。

【請求項 13】

仮想装置階層を利用して、各々が許容可能な最小サイズの値を有する複数からなるメモリ装置を含む仮想装置に複数からなる客体を格納する方法であって、

前記方法は、前記仮想装置階層によって、格納される客体を受信する段階と、

前記仮想装置階層によって、O、P、S、及び m に基づき前記客体が大型であるか又は

50

小型であるかを判定する段階と、を有し、

ここで、 O は客体のサイズ、 P はパリティ装置の個数、 S はデータ装置の個数、 m は前記客体を大型又は小型に分類するための前記複数からなるメモリ装置のそれぞれの許容可能な最小サイズの値の中の最大値を示し、

前記客体が大型として分類される場合、消去コーディングのためのチャンクのサイズ及び前記客体のデータチャンクのパディング量を決定する段階と、

各々が前記消去コーディングのためのチャンクのサイズを有する複数からなるデータチャンク内に前記客体を分類する段階と、

前記消去コーディングを用いて P 個のパリティチャンクを計算する段階と、

前記データチャンク及び前記パリティチャンクを格納するメモリ装置を決定する段階と

10

、
前記メモリ装置に前記データチャンク及び前記パリティチャンクを書き込む段階と、を含み、

前記客体が小型として分類される場合、データ及び複製版のためのメモリ装置を決定する段階と、

前記メモリ装置に前記データ及び前記複製版を書き込む段階と、を含むことを特徴とする方法。

【請求項 14】

前記客体は、大型でも小型でもない場合、中型として分類され、

複製又は前記消去コーディングは、前記複数からなるメモリ装置で利用可能な1つ以上の空間、前記仮想装置に格納された前記客体に対するアクセス時間、及び前記客体がアクセスされる頻度に基づいて適用されることを特徴とする請求項13に記載の方法。

20

【請求項 15】

前記客体の中の第1客体及び第2客体のそれぞれに対応するパリティチャンクは、前記メモリ装置の固定されたサブセットに格納されることを特徴とする請求項13に記載の方法。

【請求項 16】

前記客体の中の第1客体及び第2客体のそれぞれに対応するパリティチャンクは、前記メモリ装置の可変されるサブセットに格納されることを特徴とする請求項13に記載の方法。

30

【請求項 17】

前記客体の中の第1客体及び第2客体のそれぞれに対応するデータ及び複製版は、前記メモリ装置の中の異なる装置に格納されることを特徴とする請求項13に記載の方法。

【請求項 18】

前記複数からなるデータチャンクの中の少なくとも1つは、ゼロでパディングされることを特徴とする請求項13に記載の方法。

【請求項 19】

仮想装置階層によって、各々が許容可能な最小サイズの値を有する複数からなるメモリ装置を含む仮想装置から、キーと共に、 O 、 P 、 S 、及び m に基づき小型、中型、又は大型の客体を読み取る方法であって、

40

ここで、 O は客体のサイズ、 P はパリティ装置の個数、 S はデータ装置の個数、 m は前記複数からなるメモリ装置のそれぞれの許容可能な最小サイズの値の中の最大値を示し、

前記方法は、前記仮想装置階層によって、データ読み取り要請に対応するキーを受信する段階と、

前記仮想装置階層によって、前記メモリ装置の全てに読み取り要請を送信する段階と、

前記仮想装置階層によって、前記メモリ装置から応答を受信する段階と、を有し、

前記客体が大型又は中型である場合、前記仮想装置階層によって、データチャンク及びパリティチャンクを受信し、消去コーディングを利用して前記客体を復元し、

前記客体が小型である場合、前記受信された応答は、前記客体又は前記客体の複製版であることを特徴とする方法。

50

【請求項 20】

前記キーは、前記複数からなるメモリ装置の中から開始装置又は第 1 装置を決定するためのハッシュ（キー）を含むことを特徴とする請求項 19 に記載の方法。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、キーバリューストレージシステムに関し、より詳細には、仮想装置階層を利用した複数のメモリ装置を含む仮想装置に対する客体の格納及び読み取り方法とこれを用いたストレージ装置に関する。

10

【背景技術】**【0002】**

従来のソリッドステートドライブ（Solid State Drive：SSD）は、通常ブロックインターフェースのみを使用し、複数配列の独立ディスク（Redundant Array of Independent Disk：RAID）、消去コーディング（erasure coding）又は複製（replication）を通じてデータ信頼性を提供する。客体（object）フォーマットの大きさが多様になり非構造化されると共に、客体とブロックレベルのインターフェース間の効果的なデータ変換に対する要求がある。また、空間効率性及び速いアクセス時間（access time）の特性を維持しつつも、データ信頼性を保証することが好ましい。

20

【発明の概要】**【発明が解決しようとする課題】****【0003】**

本発明は、上記従来技術に鑑みてなされたものであって、本発明の目的は、仮想装置階層を利用した複数のメモリ装置を含む仮想装置に対する客体の格納及び読み取り方法とこれを用いたストレージ装置を提供することにある。

【課題を解決するための手段】**【0004】**

本発明は、ブロック装置とは異なるキーバリューストレージシステム（例：キーバリュースソリッドステートドライブ）に関するものである。

30

【0005】

一部の実施形態は、キーバリュースソリッドステートドライブに対してデータ信頼性を具現する方法に関するものである。空間オーバーヘッド（space overhead）に基づいて、複製及び消去コーディングのハイブリッドがキーバリュースソリッドステートドライブ群に適用され、これは客体に対するステートレス（stateless）の可変長消去コードを具現する。

【0006】

一部の実施形態は、以下の特徴のうちの 1 つ以上の特徴を有する。1）信頼性は、固定ブロック毎ではなく、可変客体毎に提供される。2）単一のディスク群に対する客体の目標信頼性を具現するために、複製及び消去コーディングを混合して適用する。3）客体に対する好ましい技法の決定において、空間効率性は第 1 メトリック（primary metric）であり、性能は第 2 メトリックである。4）メカニズムは、複数配列の独立ディスク（RAID）に類似するステートレスである。5）複製又は消去コーディングのための更なる情報が格納される必要はない。6）客体のサイズに拘らず、アップデートのための読み取り - 修正 - 書き込み（read-modify-write）動作の必要がない。

40

【0007】

一部の実施形態は、キーバリュースソリッドステートドライブ群の信頼性を具現する方法を提供する。また、本発明は、ブロック内のデータの一部がアップデートされた場合、ブ

50

ロック装置に発生する読み取り - 修正 - 書き込み動作を避けることができる。本発明によると、信頼性は、ブロック（例：固定ブロック）毎ではなく、客体（例：可変客体）毎に提供されるためである。

【0008】

上記目的を達成するためになされた本発明の一態様によるストレージ装置は、ステートレス（stateless）のデータ保護を利用する仮想装置として構成された複数からなるメモリ装置と、複数からなる客体のそれぞれのサイズに応じて、前記客体の一部に第1データ保護を適用し、前記客体の他の一部に第2データ保護を適用することにより、前記客体を格納するように前記仮想装置を管理する仮想装置階層と、を有する。

【0009】

前記メモリ装置は、1つ以上のデータ装置及び1つ以上のパリティ（parity）装置として構成され得る。

第1データ保護は、消去コーディング（erasure coding）を含み、前記第2データ保護は、複製（replication）を含み得る。

前記消去コーディングは、前記客体のうちの該当する客体が大型客体として分類される場合、データ保護のために利用され得る。

前記客体のうちの該当客体が不等式 $((P + 1) \times O > (S + P) \times m \text{ AND } O \geq S \times m)$ を満足する場合、前記大型客体として分類され得る。ここで、Oは客体のサイズ、Pはパリティ装置の個数、Sはデータ装置の個数、mは許容可能な最小サイズの値を示す。

前記複製は、前記客体のうちの該当客体が小型客体として分類される場合、データ保護のために利用され得る。

前記客体のうちの該当客体が不等式 $((P + 1) \times O = < (S + P) \times m)$ を満足する場合、小型客体として分類され得る。ここで、Oは客体のサイズ、Pはパリティ装置の個数、Sはデータ装置の個数、mは許容可能な最小サイズの値を示す。

前記消去コーディング又は前記複製のいずれかは、前記客体のうちの該当客体が大型客体としても小型客体としても分類されない場合、性能メトリック（performance metrics）及びデータ使用特性に基づいて、データ保護のために利用され得る。

前記複数の客体のうちの該当客体が不等式 $((P + 1) \times O > (S + P) \times m) > S \times m > O$ を満足する場合、中型客体として分類され得る。ここで、Oは客体のサイズ、Pはパリティ装置の個数、Sはデータ装置の個数、mは許容可能な最小サイズの値を示す。

前記パリティ装置は、1つ以上の大型客体を格納する場合、固定され得る。

前記パリティ装置は、1つ以上の大型客体を格納する場合、循環（rotate）し得る。

前記メモリ装置は、ソリッドステートドライブ（Solid State Drive）を含み得る。

【0010】

上記目的を達成するためになされた本発明の一態様による仮想装置階層を利用して複数からなるメモリ装置を含む仮想装置に複数からなる客体を格納する方法は、前記仮想装置階層によって、前記客体のうちの該当客体が大型であるか又は小型であるかを判定する段階と、前記客体のうちの該当客体が大型として分類される場合、消去コーディングのためのチャンク（chunk）のサイズ及び前記客体のうちの該当客体の1つ以上のデータチャンクのパディング（padding）量を決定する段階と、前記消去コーディングを用いてP個のパリティチャンクを計算する段階と、前記データチャンク及び前記パリティチャンクを格納するメモリ装置を決定する段階と、前記メモリ装置に前記データチャンク及び前記パリティチャンクを書き込む段階と、前記客体のうちの該当客体が小型として分類される場合、複製を通じて生成されたデータ及び複製版（replicas）のためのメモリ装置を決定する段階と、前記メモリ装置に前記データ及び前記複製版を書き込む段階

10

20

30

40

50

と、を有する。

【0011】

前記客体のうちの該当客体は、大型でも小型でもない場合、中型として分類され、前記複製又は前記消去コーディングは、性能メトリック及びデータ使用特性に基づいて適用され得る。

前記客体のうちの少なくとも2つの客体に対応するパリティチャンクは、前記メモリ装置の固定されたサブセット (s u b s e t) に格納され得る。

前記客体のうちの異なる客体に対応するパリティチャンクは、前記メモリ装置の固定されたサブセットに格納され得ない。

前記客体のうちの少なくとも2つの客体に対応するデータ及び複製版は、前記メモリ装置のうちの異なる装置に格納され得る。

10

前記1つ以上のデータチャンクのうちの少なくとも1つは、ゼロでパディングされ得る。

【0012】

上記目的を達成するためになされた本発明の一態様による仮想装置階層を利用して複数からなるメモリ装置を含む仮想装置からキー (k e y) を含む客体を読み取る方法は、前記仮想装置階層によって、前記メモリ装置の全てに読み取り要請を送信する段階と、前記仮想装置階層によって、前記メモリ装置から応答を受信する段階と、を有し、前記客体が大型である場合、前記仮想装置階層によって、データチャンク及びパリティチャンクを受信し、消去コーディングを利用して前記客体を復元し、前記客体が小型である場合、前記データチャンクは、前記客体又は前記客体の複製版である。

20

【0013】

前記キーは、前記複数のメモリ装置のうちから開始装置 (s t a r t d e v i c e) 又は第1装置 (p r i m a r y d e v i c e) を決定するためのハッシュ (即ち、キーのハッシュ) を含み得る。

【発明の効果】

【0014】

本発明によれば、空間オーバーヘッドに基づいて、大型客体の場合に消去コーディング、小型客体の場合に複製を用いたステートレスのハイブリッドが利用される。また、中型客体は、アクセスパターン等に基づいて消去コーディングと複製との間の切り換えが行われる。また、チャンクのサイズは客体毎に変わる。これにより、ブロック内のデータの一部がアップデートされた場合にブロック装置に発生する他の客体との空間共有による読み取り - 修正 - 書き込み動作を避けることができ、空間効率性及び速いアクセス時間の特性を維持しながらデータ信頼性を保証することができる。

30

【図面の簡単な説明】

【0015】

【図1】本発明の一実施形態によるキーバリュースリッドステートドライブの概略的な構成図である。

【図2】本発明の一実施形態による装置グループを含む仮想装置及び仮想装置における客体のストレージを示す概念図である。

40

【図3】本発明の一実施形態による仮想装置に客体を書き込む方法のフローチャートである。

【図4】本発明の一実施形態による共有パリティ方法を用いた図2の仮想装置における大型客体のストレージを示す概念図である。

【図5】本発明の一実施形態による専用パリティ方法を用いた図2の仮想装置における大型客体のストレージを示す概念図である。

【図6】本発明の一実施形態による図2の仮想装置における小型客体のストレージを示す概念図である。

【図7】本発明の一実施形態による仮想装置から客体を読み取る方法のフローチャートである。

50

【図 8】本発明の一実施形態による仮想装置から客体を読み取る方法のフローチャートである。

【発明を実施するための形態】

【0016】

以下、本発明を実施するための形態の具体例を、図面を参照しながら詳細に説明する。本発明は、多様な様相の実施形態を有することができ、本明細書で説明する実施形態に限定して解釈すべきではない。本明細書で説明する実施形態は、本発明が徹底且つ完全になるように例示として提供し、これを通じて本発明が属する技術分野の通常の知識を有する者に、本発明の様相及び特徴を十分に伝える。従って、本発明が属する技術分野の通常の技術者が本発明の様相及び特徴を完全に理解するに当たって、不要なプロセス、構成、及び技法について説明を省略する。特に言及しない限り、図面及び詳細な説明全般に亘って同一の参照符号は、同一の構成要素を示し、同一の構成要素に関する重複説明はしない。図面に示す構成要素、階層、及び領域は、明確性のために誇張することがある。

【0017】

本明細書で本発明の特定の実施形態を示して説明しても、特許請求の範囲によって定義される本発明の思想や範囲、及びこれらの均等範囲から逸脱することなく、説明する実施形態について、ある変更及び変形が行われ得るということは、本発明が属する技術分野の通常の技術者にとって当然である。例えば、通常の技術者が理解するように、多様な図面に示す実施形態の特徴は、本発明の思想及び範囲から逸脱することなく、結合される。

【0018】

本明細書において、「第 1」、「第 2」、「第 3」等の用語は、多様な要素、構成、領域、階層、及び／又はセクションを説明するために用いられるが、このような用語によって、要素、構成、領域、階層、及び／又はセクションが制限されるべきではない。このような用語は、何れかの要素、構成、領域、階層、又はセクションを他の要素、構成、領域、階層、又はセクションと区分するために用いられる。従って、以下に記載する第 1 要素、構成、領域、階層、又はセクションは、本発明の思想及び範囲を逸脱することなく、第 2 要素、構成、領域、階層、又はセクションと指称することができる。

【0019】

要素や階層が他の要素や階層「上に」、「に連結された」、又は「に結合された」と言及する場合、これは直接、「他の要素や階層上に」、「他の要素や階層に連結された」、又は「他の要素や階層に結合された」ことが可能である、或いは 1 つ以上の媒介要素や媒介階層が存在してもよいことが理解される。また、1 つの要素や階層が 2 つの要素や階層の「間」にあると言及する場合、2 つの要素や階層の間に 1 つの要素や階層があるか、又は 1 つ以上の媒介要素や媒介階層があってもよいということも理解される。

【0020】

本明細書で用いる用語は、特定の実施形態を説明するための目的のみであって、本発明を制限しようとするものではない。本明細書で用いる単数形態の用語は、文脈上明らかに特に指示がない限り、複数形態の用語も含む。本明細書において、「含む」という用語を用いる場合、このような用語は、言及する特徴、数字、段階、動作、要素、及び／又は構成の存在を明示するが、1 つ以上の他の特徴、数字、段階、動作、要素、構成、及び／又はこれらの集合の存在や付加を排除するものではない。本明細書で用いる「及び／又は」、「及び／若しくは」という用語は、挙げられる 1 つ以上の関連項目の任意の何れか 1 つ及び全ての組み合わせを含む。要素が挙げられた後に続く「少なくとも 1 つ」のような表現は、挙げられた要素の全目録を修飾するものであり、目録の個別の要素を修飾するものではない。

【0021】

本明細書で用いる「実質的に」、「約」、及びこれと類似する用語は、近似 (approximation) を示す用語として用いられるものであり、程度 (degree) を示す用語として用いられるものではなく、本発明が属する技術分野の通常の技術者に認識される測定値又は計算値に内在した偏差を説明するためのものである。また、本発明の実

10

20

30

40

50

施形態に関する説明で用いる「し得る及び／又はできる」という表現は、「本発明の１つ以上の実施形態」を言及するものである。本明細書で用いられる「用いる」及び「用いられる」という用語は、それぞれ「利用する」及び「利用される」という用語と同じ意味と見なされる。更に、「例示的な」という用語は、例示又は一例を指称する。

【 0 0 2 2 】

本明細書で説明する本発明の実施形態によると、電子／電気装置及び／又は任意の他の関連装置や構成（例：ホスト、ソリッドステートドライブ、メモリ装置、及び仮想装置階層）は、任意の適切なハードウェア、ファームウェア（例：特定用途向け集積回路）、ソフトウェア、又ははソフトウェア、ファームウェア、及びハードウェアの適切な組み合わせを利用して具現される。例えば、キーバリュ（Key Value：KV）のソリッドステートドライブ（Solid State Drive：SSD、以下「SSD」と称する）、ホスト、SSD、メモリ装置、及び仮想装置階層などの装置の多様な構成要素は、１つの集積回路（Integrated Circuit：IC）チップ又は個別のICチップ上に形成される。また、このような装置の多様な構成要素は、フレキシブル印刷回路フィルム（flexible printed circuit film）、テープキャリアパッケージ（Tape Carrier Package：TCP）、印刷回路基板（Printed Circuit Board：PCB）上に具現されるか又は１つの基板（substrate）上に形成される。更に、このような装置の多様な構成要素は、１つ以上のコンピューティング装置の１つ以上のプロセッサで実行され、コンピュータプログラムの命令（語）を実行し、本明細書に記載する多様な機能を行うための他のシステム構成要素と相互作用するプロセス又はスレッド（thread）である。コンピュータプログラムの命令（語）はメモリに格納され、メモリは、例えばランダムアクセスメモリ（Random Access Memory：RAM）又はフラッシュメモリ（例：NANDフラッシュメモリ）装置などの標準メモリ装置を用いるコンピューティング装置で具現される。コンピュータプログラムの命令（語）は、例えばCD-ROM、フラッシュドライブ等のような他の非一過性のコンピュータ読み取り可能な記録媒体に格納される。また、本発明が属する技術分野の通常の技術者は、本発明の実施形態の思想や範囲を逸脱することなく、多様なコンピューティング装置の機能が単一のコンピューティング装置に結合若しくは統合されるか、又は特定のコンピューティング装置の機能が１つ以上の他のコンピューティング装置に分散されるということを認識する。

【 0 0 2 3 】

本明細書で用いる技術用語及び科学用語を含む全ての用語は、特に定義しない限り、本発明が属する技術分野の通常の知識を有する者が一般的に理解するものと同様の意味を有する。また、通常用いられる辞典に定義されているような用語は、関連技術及び／又は本明細書の文脈上の意味と一致すると解釈されるべきであり、本明細書で明らかに定義しない限り、理想的又は過度に形式的な意味として解釈されるべきではない。

【 0 0 2 4 】

図１は、本発明の一実施形態によるキーバリュSSD 10の概略的な構成図である。本実施形態によるストレージシステム（又はストレージ装置）は、図１に示すように、１つ以上のキーバリュSSDを含む。但し、本発明が図１に限定されるわけではない。

【 0 0 2 5 】

本実施形態によると、キーバリュSSD 10のキーバリュAPI 15は、従来のブロックマッピングを必要としないユーザーキーバリュ装置ドライバー 20で動作する。

【 0 0 2 6 】

本実施形態によると、キーバリュSSD 10を含むストレージシステムは、所望の信頼性（例：目標信頼性）を得るために、客体のそれぞれのサイズに応じて、第１データ保護（例：消去コーディング）を客体の一部に適用し、第２データ保護（例：複製）を客体の他の一部に適用するハイブリッドステートレスのデータ保護方法を利用する。このような方法により、信頼性を損なわずに空間効率的なソリューションが提供される。一部の実施形態において、キーバリュSSD 10自体がハイブリッドステートレスのデータ保護

方法を行うこともできるが、ストレージシステムによりハイブリッドステートのデータ保護方法が行われる場合には、例えば複数のドライブ (SSD) に亘る管理 (例: 仮想装置階層の動作) がより容易になる。

【0027】

本実施形態によると、客体は、空間効率性のために分類 (classification) され、サイズに基づいて分類され、それぞれのサイズのクラス毎に相違するバックアップの方法が用いられる。

【0028】

客体に対する消去コーディングの空間オーバーヘッドが客体に対する複製の空間オーバーヘッドよりも小さい場合、該当客体は、大型客体 (large object) と見なされる。その場合、消去コーディングがより好ましい。消去コーディングは、空間フットプリント (footprint) が少ないためである。言い換えると、客体が不等式 $((P+1) \times O > (S+P) \times m \text{ AND } O \geq S \times m)$ 、ここで、 O は客体のサイズである) を満足する場合、大型客体と見なされる。上記の不等式と以下の不等式において、 O は客体のサイズ (即ち、客体サイズ)、 P はパリティ装置カウント (即ち、仮想装置におけるパリティ装置の個数)、 S はデータ装置カウント (即ち、仮想装置におけるデータ装置の個数) を意味し、 m は許容可能な最小サイズの値 (即ち、個々の装置における全ての最小値のサイズのうちの最大値) を意味する。例えば、本実施形態による「許容可能な最小サイズの値」は、任意の装置の最小値のサイズの要件に違反しないシステムにおいて、任意の装置に格納される値のサイズを指称する。それぞれの装置毎に、該当装置が支援する最小の客体のサイズを有する。本実施形態によると、客体は、全ての装置に対して同一のサイズに分割されるため、分割のサイズは、装置が支援する任意の最小のサイズよりも大きくなければならない。 m よりも小さいサイズの客体を格納しようとする場合、少なくとも1つの装置は、客体を格納できなくなる。

【0029】

言い換えると、以下の条件の両方満足する場合、客体は大型客体と見なされる。条件1) 客体のサイズ O にパリティ装置の個数よりも1つ大きい値 $(P+1)$ をかけた値が、許容可能な最小サイズの値 m にデータ装置の個数 S とパリティ装置の個数 P との合計 $S+P$ をかけた値よりも大きい。そして、条件2) 客体のサイズ O がデータ装置の個数 S に許容可能な最小サイズの値 m をかけた値よりも大きい。

【0030】

本実施形態によると、大型客体は消去コーディングされる。即ち、客体は、 S 個のチャンク (chunk) (即ち、データチャンク又は S 個のポーション (portion)) に分割され、パリティチャンク (即ち、パリティポーション) は、 S 個のチャンクを用いて計算される。 S 個及び P 個のチャンクのそれぞれは、本明細書で記載するように、対応する装置に格納される。

【0031】

客体に対する複製の空間オーバーヘッドが客体に対する消去コーディングの空間オーバーヘッドよりも小さい場合、該当客体は小型客体 (small object) と見なされる。その場合、複製が好ましい。複製は、よりよい読み取り性能を提供し、相対的に複雑な消去コーディングよりもアップデートをうまく処理できるためである。これは、アプリケーションのメタデータが小さくなる傾向があるという観測に鑑みても合理的である。言い換えると、客体が不等式 $((P+1) \times O \leq (S+P) \times m)$ を満足する場合、小型客体と見なされ、複製される。

【0032】

即ち、客体のサイズ O にパリティ装置の個数よりも1つ大きい値 $(P+1)$ をかけた値が、データ装置の個数 S とパリティ装置の個数 P との合計 $S+P$ に許容可能な最小サイズの値 m をかけた値よりも小さい場合、客体は小型客体と見なされる。

【0033】

客体が小型又は大型として分類される一部のグレー領域 (gray area) が存在

10

20

30

40

50

する。例えば、客体が不等式 $((P+1) \times O > (S+P) \times m > S \times m > O)$ を満足する場合、客体は中型客体 (medium object) と見なされ、性能メトリック (performance metrics) (例：空間対アクセス時間) 及び/又はデータ使用特性 (例：アップデートの頻度) に基づいて、複製や消去コーディングが用いられる。

【0034】

言い換えると、客体のサイズ O にパリティ装置の個数よりも1つ大きい値 $(P+1)$ をかけた値が、データ装置の個数 S とパリティ装置の個数 P との合計に許容可能な最小サイズの値 m をかけた値よりも大きく、データ装置の個数 S とパリティ装置の個数 P との合計に許容可能な最小サイズの値 m をかけた値がデータ装置の個数 S に許容可能な最小サイズの値 m をかけた値よりも大きく、データ装置の個数 S に許容可能な最小サイズの値 m をかけた値が客体のサイズ O よりも大きい場合、客体は中型客体と見なされる。

【0035】

例えば、性能がより重要であり、客体が頻繁にアップデートされる場合、複製がよりよい選択になる。このような場合、中型客体は小型客体として分類される。例えば、不等式 $((P+1) \times O = (S+P) \times m) \text{ OR } ((P+1) \times O > (S+P) \times m) \text{ AND } S \times m > O$ 、即ち $((P+1) \times O \leq (S+P) \times m \text{ OR } O < S \times m)$ を満足する場合、客体は、本実施形態において小型として分類される。

【0036】

他の例として、空間効率性がより重要な場合、消去コーディングが用いられる。このような場合、中型客体は大型客体として分類される。例えば、不等式 $((P+1) \times O > (S+P) \times m \text{ AND } O \geq S \times m) \text{ OR } ((P+1) \times O > (S+P) \times m) > S \times m > O = ((P+1) \times O > (S+P) \times m)$ 、即ち $((P+1) \times O > (S+P) \times m)$ を満足する場合、客体は本実施形態において大型として分類される。

【0037】

図2は、本発明の一実施形態による装置グループを含む仮想装置及び仮想装置における客体のストレージを示す概念図であり、装置グループ $SSD1$ 、 $SSD2$ 、 $SSD3$ 、 $SSD4$ 、 $SSD5$ 、 $SSD6$ を含む仮想装置200及び仮想装置200内の客体 (大型客体202及び小型客体204) のストレージを示す。装置のうち、 $SSD1$ 、 $SSD2$ 、 $SSD3$ 、及び $SSD4$ はデータ装置として構成され、 $SSD5$ 及び $SSD6$ はパリティ装置として構成される。図2には、例示の目的として、4つのデータ装置 ($SSD1$ 、 $SSD2$ 、 $SSD3$ 、 $SSD4$) 及び2つのパリティ装置 ($SSD5$ 、 $SSD6$) のみを図示しているが、仮想装置200内のデータ及びパリティ装置の個数は、これに限定されない。更に、他の SSD がデータ装置及びパリティ装置として構成される。

【0038】

例えば、仮想装置200は、 S 個のデータ装置及び P 個のパリティ装置の全体を含み、パリティ装置は、固定 (fixed) 又は循環する (図2を参照すると、例として S 値に4、 P 値に2が用いられる)。例えば、パリティ装置が循環する場合、異なる大型客体の全てのパリティチャンク (又はパリティポーション) がパリティ装置のうちの同一の装置に格納されずに、データ装置の一部は、1つ以上の大型客体のためのパリティ装置として用いられる。言い換えると、パリティ装置が固定された場合は、客体に対応する「 P 」個のパリティチャンクがメモリ装置の同一の「 P 」セットに格納されるのに対し、パリティ装置が循環する場合は、客体に対応する「 P 」個のパリティチャンクが必ずしもメモリ装置の同一のセットに格納されない。また、装置は、平面 (flat) 又は階層的 (hierarchy) 構成で組織される。多数の装置に亘って拡散 (spread) 又は複製された客体に対する開始装置は、キーのハッシュ値 (hash value) によって決定される。

【0039】

また、データ装置は、必要及び／又はユーザーの設計上の選択に応じて、パリティ装置として再構成され、その逆の場合も同様である。例えば、仮想装置 200 の装置の個数は、目標信頼性に基づいて設定可能である。消去コーディングにおいて、P 個の欠陥に耐えるための装置の総数は、データ装置の個数 S とパリティ装置の個数 P との合計になる。複製において、P 個の欠陥に耐える装置の総数は、P + 1 になる。装置の容量は、互いに同一又は類似する。

【0040】

本実施形態によると、仮想装置 200 内の装置セット、又は仮想装置 200 に対応する装置セットは、信頼性管理の単位となるグループを構成する。装置 SSD 1、SSD 2、SSD 3、SSD 4、SSD 5、SSD 6 のグループは、単一のサーバーやラック (rack) 内に存在するか、又はサーバーやラック全域に亘って存在し、階層アーキテクチャ又は平面アーキテクチャを有するように構造化される。

10

【0041】

装置グループを含む仮想装置 200 は、仮想装置階層 210 と称される階層によって管理されるため、装置グループは、単一の仮想装置として提供される。仮想装置階層 210 は、ステートレスである。仮想装置階層 210 は、実行の際 (runtime) に、客体の個数、使用可能な容量、及び／又は同類のものなどの装置の最小限のメタデータ情報をキャッシュして維持する。特に、本実施形態による仮想装置階層 210 は、キー情報を維持する必要がない (例：キーに対するマッピングがない)。仮想装置 200 の容量は、全ての装置容量のうちの最小容量 (例：図 2 の場合、SSD 1、SSD 2、SSD 3、SSD 4、SSD 5、及び SSD 6 の容量のうちの最小容量) にグループ内装置の個数をかけた値によって決定される。

20

【0042】

仮想装置階層 210 は、各装置が処理可能な最小値のサイズ及び最大値のサイズを知る。仮想装置階層 210 は、仮想装置 200 の最小値のサイズ及び最大値のサイズを決定する。例えば、本実施形態によると、個々の装置の全ての最小値のサイズ m_i のうちの最大値は、仮想装置 200 の最小値のサイズ m として定義され、個々の装置の全ての最大値のサイズ M_i のうちの最小値は、仮想装置 200 の最大値のサイズ M として定義される。他の実施形態において、仮想装置の最大値のサイズ M は、個々の装置の全ての最大値のサイズ M_i のうちの最小値にデータ装置の個数 S をかけた値によって定義される。

30

【0043】

一部の実施形態による仮想装置 200 は、本発明の技術分野に属する通常の技術者に知られている任意の適切な消去コーディングアルゴリズムを利用し、使用可能な MDS (Maximum Distance Separable) アルゴリズム (例：リードソロモン符号 (Reed-Solomon code)) を用いる。図 2 に示すように、大型客体 202 には、「2」のパリティ値を有する消去コード (即ち、消去コーディングアルゴリズム) が適用され、パリティ 1 及びパリティ 2 が用いられる。

【0044】

本実施形態によると、客体 (例：大型客体 202) は、S 個のチャンクに分割され符号化される (データ及びパリティ装置 (即ち、S + P 個の装置) に同一のサイズで分散される)。例えば、消去コーディングにより符号化された大型客体 202 は、データ 1、データ 2、データ 3、データ 4、パリティ 1、及びパリティ 2 に分割される。客体が占める実際のストレージ空間は、本発明の実施形態において帯域 (band) と指称する。帯域は、消去コーディングの場合、S + P 個の装置に亘って存在するが、複製の場合には、P + 1 個の装置に亘って存在する。例えば、複製は、小型客体 204 に適用される。帯域は、客体を完全に含む (即ち、客体全体が帯域内に格納される)。一部の実施形態において、帯域は、客体の S 個のチャンクが格納される S 個の装置に亘って存在する。

40

【0045】

客体のサイズが装置の割当単位又は整列単位と合わない場合、帯域内の客体に対して割り当てられた余分の空間がパディングされる (例：「0」でパディングされる)。例えば

50

、図2は、大型客体202のデータ4が「0」でパディングされ、データ4の全てのデータビットを格納するために必要でない余分の空間を占めることを示している。帯域のサイズは、実施形態によって変わる。

【0046】

図3は、本発明の一実施形態による仮想装置（例：図2、図4～6の仮想装置200）に客体を書き込む方法のフローチャートである。図4は、本発明の一実施形態による共有パリティの方法を用いた仮想装置200における大型客体（242、244）のストレージを示す概念図である。図5は、本発明の一実施形態による専用パリティの方法を用いた仮想装置200における大型客体（242、244）のストレージを示す概念図である。図6は、本発明の一実施形態による仮想装置200の小型客体（客体1（262）、客体2（264））のストレージを示す概念図である。

10

【0047】

図3に示すように、段階300で、仮想装置階層（例：図2、4～6の仮想装置階層210）は、キーを含むサイズ0の客体を仮想装置（例：図2、4～6の仮想装置200）に書き込みというインストラクション（instructions）又は命令（command）を（例えば、ホスト装置から）受信する。他の実施形態において、書き込みのインストラクション又は命令は、ホストによって提供された書き込みのインストラクションに応答して、仮想装置階層によって生成される。

【0048】

段階302で、仮想装置階層は、上述のような不等式を用いて客体が大型か否かを判定する。例えば、 $((P+1) \times O > (S+P) \times m \text{ AND } O \geq S \times m)$ である場合、客体は大型と見なされる。ここで、Oは客体のサイズ、Pはパリティ装置の個数、Sはデータ装置の個数、mは許容可能な最小サイズの値（即ち、個々の装置における全ての最小値のサイズのうちの最大値）を意味する。

20

【0049】

段階302で客体が大型客体として分類された場合、段階312に示すように、仮想装置階層は、消去コーディングのためのデータチャンクのサイズを決定し、1つ以上のデータチャンクに対するパディングの量（例：「0」でパディング）を決定する。その後、客体は、パディングとの整列（alignment with padding）を考慮して同一サイズのS個のチャンクに分割され、段階314に示すように、本発明の技術分野に属する通常の技術者に知られている適切な消去コーディングアルゴリズムを利用して、S個のチャンクからP個のコードチャンク（即ち、P個のパリティチャンク）が生成（例：計算）される。

30

【0050】

段階316で、仮想装置階層は、分散ポリシー（distribution policy）に基づいて、データチャンク及びパリティチャンクを格納するための装置（即ち、S個の装置及びP個の装置）を決定する。例えば、分散ポリシーは、キーのハッシュ値により客体に対する開始装置を決定すること、並びにノ又は固定された装置及びノ若しくはスポット（spot）にデータ及びノ若しくはパリティチャンクを格納することを含む。段階318で、データチャンク及びパリティチャンクは、対応する装置に書き込まれる。例えば、S+P個のチャンクは、S+P個の装置（例：図2のSSD1、SSD2、SSD3、SSD4、SSD5、SSD6）へのストレージのために分散される。例えば、図4に示す循環するパリティ装置において、キーのハッシュにより決定された装置でデータの書き込みが始まり、それぞれのブロック、即ちチャンク（及びパリティブロック、即ちパリティチャンク）は、第1装置の最初のデータから始まって順に書き込まれる。例えば、図5に示すように、固定されたパリティ装置において、全てのデータ及びパリティブロック（即ち、チャンク）は、予め指定された装置に格納される。その場合、開始装置もまた、予め指定される。小型客体の場合、開始装置及び複製装置もまた、キーをハッシュすることにより決定される。

40

【0051】

50

図4に示すように、パリティ装置は、共有（即ち、循環）される。即ち、格納された大型客体に応じて、データチャンクを格納するためのデータ装置又はパリティチャンクを格納するためのパリティ装置の両方に単一の装置が用いられる。例えば、客体242は、データ1、データ2、データ3、（「0」でパディングされた）データ4、パリティ1、及びパリティ2に分割され、客体244もまた、データ1、データ2、データ3、（「0」でパディングされた）データ4、パリティ1、及びパリティ2に分割される。図4から分かるように、大型客体242のデータ1、データ2、データ3、及びデータ4は、それぞれ仮想装置200のSSD1、SSD2、SSD3、及びSSD4に格納されるのに対し、大型客体244のデータ1、データ2、データ3、及びデータ4は、それぞれ仮想装置200のSSD6、SSD1、SSD2、及びSSD3に格納される。

10

【0052】

また、客体242のパリティ1及びパリティ2がそれぞれ仮想装置200のSSD5及びSSD6に格納されるのに対し、客体244のパリティ1及びパリティ2は、それぞれ仮想装置200のSSD4及びSSD5に格納される。従って、合わせてS個のデータ装置及びP個のパリティ装置が存在しても、パリティ装置が循環するため、専用のパリティ装置はない。

【0053】

図4に示す例とは異なり、図5は、専用パリティ装置、即ち仮想装置200のSSD5及びSSD6を利用する具現例を示す。例えば、大型客体242及び大型客体244の両方のデータ1、データ2、データ3、（「0」でパディングされた）データ4、パリティ1、及びパリティ2は、それぞれ仮想装置200のSSD1、SSD2、SSD3、SSD4、SSD5、及びSSD6に格納される。

20

【0054】

循環するパリティの実施形態及び小型客体のために、客体に対する開始装置は、キーのハッシュ値によって決定される。例えば、図4の共有パリティ装置において、開始装置は、 $\text{ハッシュ(キー)} \% (S + P)$ によって決定される。その後、後続のデータ及びパリティチャンク（即ち、 $S + P$ 個のチャンク）が順に $(\text{ハッシュ(キー)} + 1) \% (S + P)$ 、 $(\text{ハッシュ(キー)} + 2) \% (S + P)$ 、...、 $(\text{ハッシュ(キー)} + S + P - 1) \% (S + P)$ に書き込まれる。専用パリティ装置が存在する場合、 $(S + P)$ の代わりにS個の装置が用いられる。

30

【0055】

データ及びパリティチャンクが対応する装置に書き込まれた後、段階320で大型客体の書き込みプロセスが完了する。

【0056】

段階302で、客体が大型と判定されない場合、プロセスは段階304に進み、客体が小型か否か（即ち、 $((P + 1) \times O = (S + P) \times m)$ を満足するか）に関して判定する。客体が小型と判定された場合、仮想装置階層は、複製を開始し、段階308で、分散ポリシーに基づき、どの装置を利用してデータ及び複製版を格納するかを決定する。例えば、分散ポリシーは、キーのハッシュ値によって客体に対する開始装置を決定すること、並びに/又は固定された装置及び/若しくはスポットにデータ及び/若しくは複製版を格納することを含む。その後、段階310で、データ及び複製版が対応する装置に書き込まれる。

40

【0057】

本実施形態によると、パディングとの整列を考慮して、 $P + 1$ 個の複製版（1つのデータコピー及びP個のパリティコピーを含む）が客体に対して生成され、複製版は、 $P + 1$ 個の装置に分散される。例えば、図6に示すように、客体1（262）は3回複製され（1つのデータ及び2つの複製版を含む）、複製版はそれぞれ仮想装置200のSSD1、SSD2、及びSSD3に格納される。同様に、客体2（264）も3回複製され（1つのデータ及び2つの複製版を含む）、複製版はそれぞれ仮想装置200のSSD3、SSD4、及びSSD5に格納される。図6に示す例において、仮想装置200は、合わせて

50

S 個のデータ装置及び P 個のパリティ装置を含む。また、客体 1 (2 6 2) 及び客体 2 (2 6 4) の両方が小型客体であるため、図 6 に示す例において、消去コーディングは用いられない。

【 0 0 5 8 】

キーのハッシュ値を用いて、S + P 個の装置のうちの第 1 装置が選択される。P 個の複製版は、ストレージ組織、性能、及びノ又は同類のもの等に基づいて決定的に選択される。例えば、データは第 1 装置に格納され、複製版は、専用パリティ装置の使用に拘らず、(ハッシュ (キー) + 1) % (S + P) 、 (ハッシュ (キー) + 2) % (S + P) 、 ... 、 (ハッシュ (キー) + P) % (S + P) 、又は異なるノード、ラックに格納される。

【 0 0 5 9 】

再び図 3 を参照すると、段階 3 0 4 で客体が小型ではないと判定された場合、即ち客体が大型でもなく (段階 3 0 2 参照) 小型でもない (段階 3 0 4 参照) 場合、客体は中型客体 (即ち、 $(P + 1) \times O > (S + P) \times m > S \times m > O$) と判定され、プロセスは段階 3 0 6 に進み、中型客体が小型客体として処理されるか否かを判定する。段階 3 0 6 で中型客体が小型客体として処理される場合、プロセスは段階 3 0 8 に進んで小型客体のストレージプロセスを開始し、段階 3 0 6 で中型客体が大型客体として処理される場合、プロセスは段階 3 1 2 に進んで大型客体のストレージプロセスを開始する。

【 0 0 6 0 】

図 7 及び図 8 は、本発明の一実施形態による仮想装置 (例：図 2、図 4 ~ 6 の仮想装置 2 0 0) から客体を読み取る方法のフローチャートである。仮想装置階層 (例：図 2、図 4 ~ 6 の仮想装置階層 2 1 0) は、キー及び値のサイズなどの客体のメタデータを維持しないため、読み取る客体が小型か大型かについて知らない。従って、仮想装置階層は、客体のユーザーキーを用いて、全ての物理装置 (即ち、S + P 個の装置) に読み取り要請を送信することで、読み取りプロセス 7 0 0 を開始し、段階 7 0 2 に示すように、サブ - 読み取り要請を全ての物理装置に送信する。段階 7 0 4 で、仮想装置階層は、装置から応答を受信する。ユーザー (例：ホスト) が要請する客体が大型である場合、段階 7 0 6 でエラーが存在しないと判定されると、S + P 個の全ての装置は、ユーザーキーを含む要請に対してそれぞれの応答を返信する。

【 0 0 6 1 】

例えば、読み取る客体が大型客体であり、エラーが無い場合、全ての装置 (即ち、S + P 個の装置) は応答する。しかし、N 個の装置にエラーがある場合、S + P - N 個の装置のみが応答する。仮想装置階層が同一サイズの任意の S 個のチャンクを受信する限り (即ち、データチャンク S の総数と同一のデータチャンク S 及びパリティチャンク P の任意の組み合わせ) 、ユーザー客体を復元 (r e b u i l d) することができる。言い換えると、装置のパリティ個数 (即ち、P と同一の装置の個数) を超える失敗が発生しない限り、大型客体の場合、データが復元される。

【 0 0 6 2 】

受信したチャンクの総数が S よりも小さいか又はチャンクのサイズが同一でない場合、エラーが存在する。全ての装置が N O N _ E X I S T エラーを返信するか又は復旧不能なエラーが発生した場合、存在しない客体の読み取りになる。

【 0 0 6 3 】

最初は、仮想装置階層が客体のタイプを知らないため、客体のタイプを「N O N E」として初期化する。段階 7 0 8 で判定されるように客体が大型である場合、段階 7 1 8 で客体タイプが判定される。段階 7 1 8 で客体タイプが「N O N E」である場合、段階 7 2 0 で客体タイプは大型に設定される。段階 7 1 8 で客体タイプが「N O N E」と判定されない場合、仮想装置階層は、段階 7 3 2 で客体タイプが大型であるか否かを判定する。段階 7 3 2 で客体タイプが大型でない場合、段階 7 3 4 に示すように、エラーと決定される。段階 7 2 0 で客体タイプが大型に設定された後、又は段階 7 3 2 で客体タイプが大型と判定された場合、仮想装置階層は、段階 7 2 2 で全てのデータチャンクを有するか否かを判定する。

10

20

30

40

50

【 0 0 6 4 】

段階 7 2 2 で仮想装置階層が全てのデータチャンクが受信されていると判定した場合、段階 7 3 0 に示すように、読み取りプロセスが完了する。一方、段階 7 2 2 で全てのデータチャンクが受信されていないと判定した場合、仮想装置階層は、段階 7 2 4 で全ての装置から応答が受信されているか否かを判定する。段階 7 2 4 で全ての装置が応答した場合、仮想装置階層は、段階 7 2 6 でデータチャンクを少なくとも S 個（全てのデータチャンク及びパリティチャンクをカウント）有するか否かを判定する。S 個よりも小さい数のチャンクが受信された場合、仮想装置階層は、段階 7 3 4 に示すように、エラーと決定する。少なくとも S 個のチャンク（受信された全てのデータチャンク及びパリティチャンクをカウント）が正確に受信された場合、仮想装置階層は、段階 7 2 8 で消去コーディングアルゴリズムを用いて S 個のチャンクで客体を復元し、段階 7 3 0 で読み取りプロセスを完了する。例えば、1 つ以上の装置が予想外にオフラインであった場合、1 つ以上の装置が応答しないこともある。従って、一部の実施形態において、全ての装置が応答しなくても、少なくとも S 個のチャンクが受信される限り、仮想装置階層は、段階 7 2 8 に示すように、客体の復元を進める。

10

【 0 0 6 5 】

仮想装置階層が段階 7 0 8 で客体が大型でないと判定した場合、プロセスは段階 7 1 0 に進み、客体タイプが「NONE」か否かを判定する。段階 7 1 0 で客体タイプが「NONE」である場合、段階 7 1 2 で客体タイプは小型に設定される。段階 7 1 0 で客体タイプが「NONE」でない場合、段階 7 1 6 で客体タイプが小型か否かに関して判定される。その際、客体タイプが小型でない場合、段階 7 3 4 に示すように、エラーと決定する。段階 7 1 2 で客体タイプが小型に設定されるか、又は段階 7 1 6 で仮想装置階層が客体タイプを小型と判定した場合、仮想装置階層は、段階 7 1 4 で、受信されたチャンクが有効か否かを判定する。段階 7 1 4 で、受信されたチャンクが有効である場合、段階 7 3 0 に示すように、読み取りプロセスが完了する。

20

【 0 0 6 6 】

ユーザー（例：ホスト）が要請した客体が小型である場合、複製版（即ち、第 1 コピー（primary copy）及び複製版（replicas）のうちの 1 つ）を有する P + 1 個の装置は、エラーが無いかなかを返信し、残りの装置は客体が存在しないことを知らせるエラーを返信する。段階 7 1 4 で仮想装置階層が任意の有効なチャンクを受信する限り、装置は客体を有する。全ての装置がNOT__EXISTエラーを返信する場合、このような客体は存在しない（又はエラーが存在する）。全ての装置が返信するのではなく、返信した全ての装置がNOT__EXISTを報告する場合、段階 7 3 4 に示すように、復旧不能なエラーが発生する。

30

【 0 0 6 7 】

仮想装置階層が段階 7 1 4 でチャンクが有効でないと判定した場合、段階 7 2 4 で全ての装置から応答が受信されているか否かに関して判定する。段階 7 2 4 で全ての装置から応答が受信されていない場合、仮想装置階層は、段階 7 0 4 で全ての装置から応答を得るように進め、図 7 に示すように、段階 7 0 6 でエラーが存在するか否か等を判定するプロセスを引き続き行う。

40

【 0 0 6 8 】

本実施形態によると、仮想装置階層は、各装置に全ての客体キーを列挙するように要請し、読み取りに失敗した場合に復元（reconstruction）するために、全てのキーに対する概念的に全体的な順序を有する。仮想装置階層は、キーを順に 1 つずつ確認する。

【 0 0 6 9 】

客体が大型である場合、仮想装置階層は、ハッシュ（キー）を用いることにより、キーに対する開始装置を決定し、固定されたパリティ装置が用いられない場合は、開始装置の情報に基づいて、どのチャンクが生成されるべきか（データチャンク又はコードチャンク）を決定する。パリティ装置が用いられる場合には、どのチャンクが復元されなければな

50

らないかが明確である。新たな装置のチャンクは、大型客体の読み取りの場合と同様に、有効なチャンクで復元される。

【 0 0 7 0 】

客体が小型である場合、仮想装置階層は、ハッシュ（キー）を用いることにより、キーに対する第 1 装置を決定し、第 1 装置情報に基づいて、どの装置が複製版を有するかを決定する。新たな装置が複製版を有さなければならない場合、客体は新たな装置に書き込まれる。これは、装置の全ての客体を巡回し（ v i s i t ）、失敗した装置が復元されるまで繰り返される。

【 0 0 7 1 】

このように、本発明の 1 つ以上の実施形態によると、空間オーバーヘッドに基づいて消去コーディング及び複製のステートレスのハイブリッドが利用される。また、中型客体は、例えばアクセス（ a c c e s s ）パターン等に基づいて、消去コーディングと複製との間の切り換えが行われる。また、チャンクのサイズは、客体毎に変わる。また、他の客体との空間共有による読み取り - 修正 - 書き込み動作はこれ以上必要ではない。

【 0 0 7 2 】

以上、本発明の実施形態について図面を参照しながら詳細に説明したが、本発明は、上述の実施形態に限定されるものではなく、本発明の技術的範囲から逸脱しない範囲内で多様に変更実施することが可能である。

【 符号の説明 】

【 0 0 7 3 】

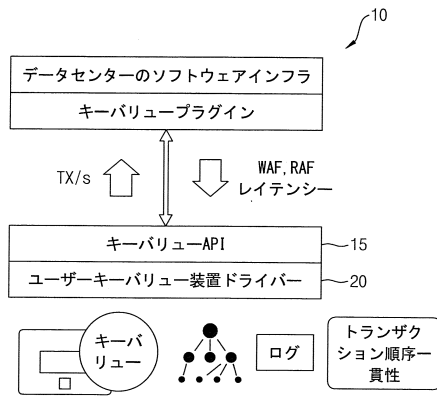
- 1 0 キーバリュース S D
- 1 5 キーバリュース A P I
- 2 0 ユーザーキーバリュース装置ドライバー
- 2 0 0 仮想装置
- 2 0 2、2 4 2、2 4 4 大型客体
- 2 0 4 小型客体
- 2 1 0 仮想装置階層
- 2 6 2、2 6 4 客体 1、客体 2（小型客体）

10

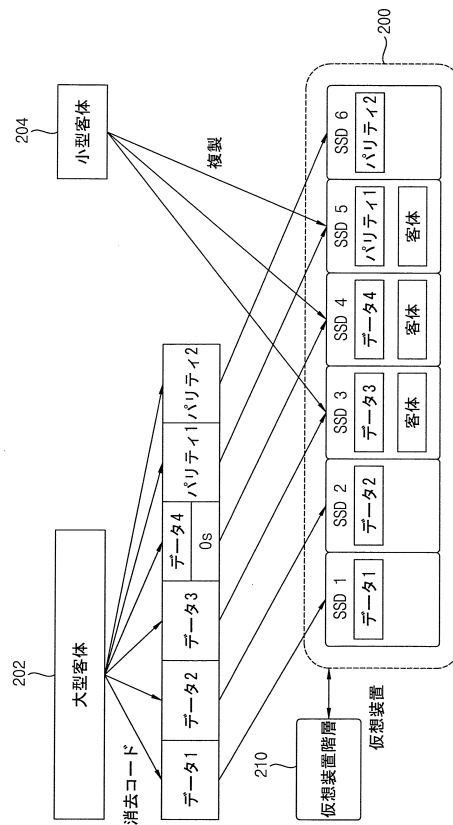
20

30

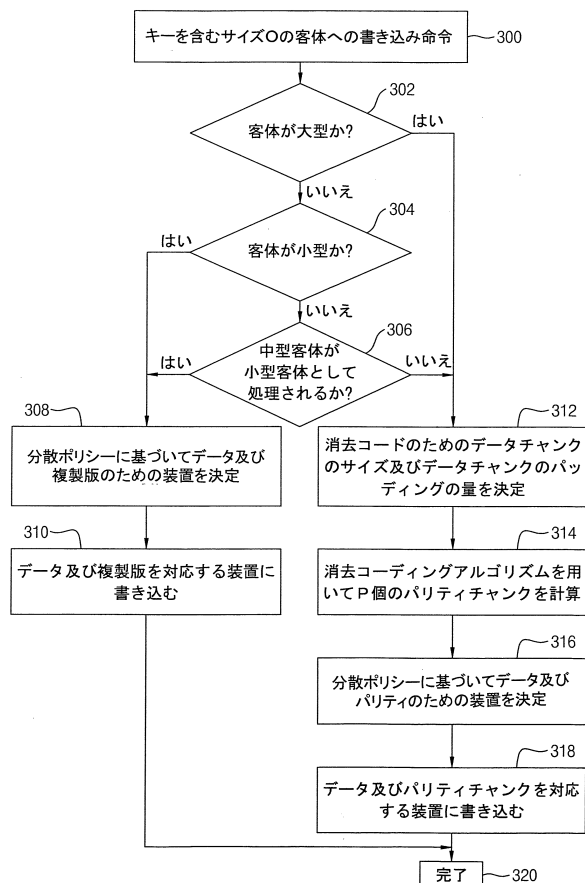
【図 1】



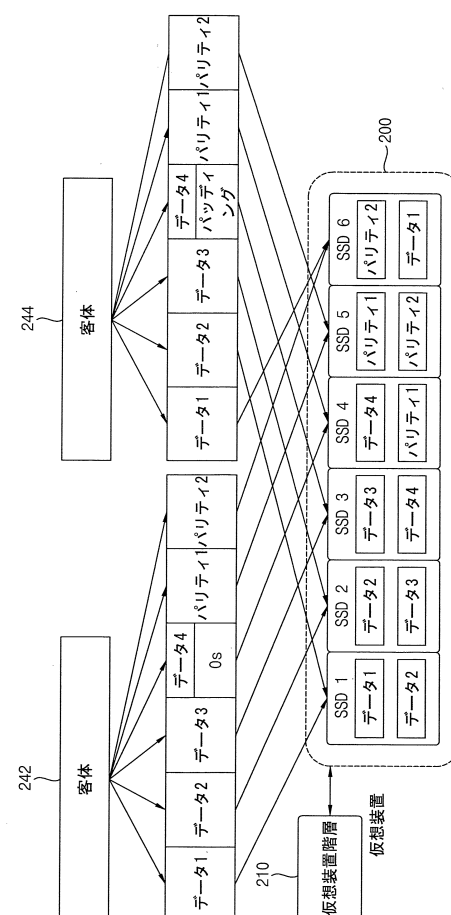
【図 2】



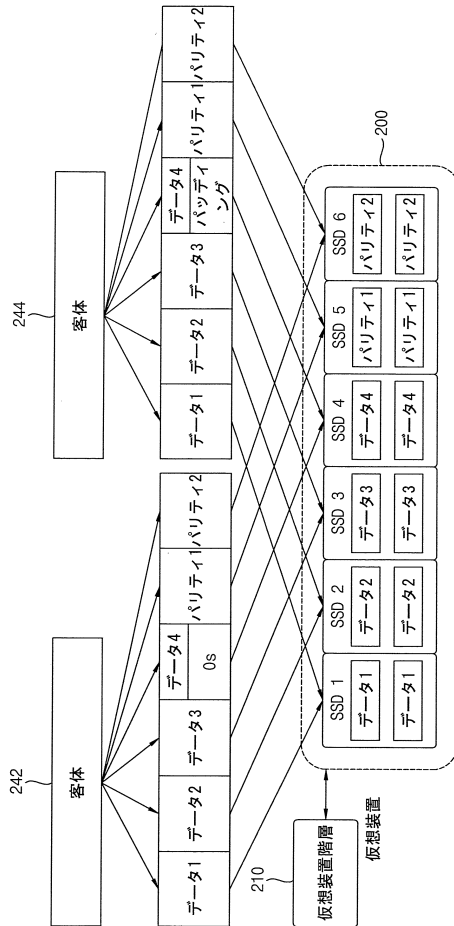
【図 3】



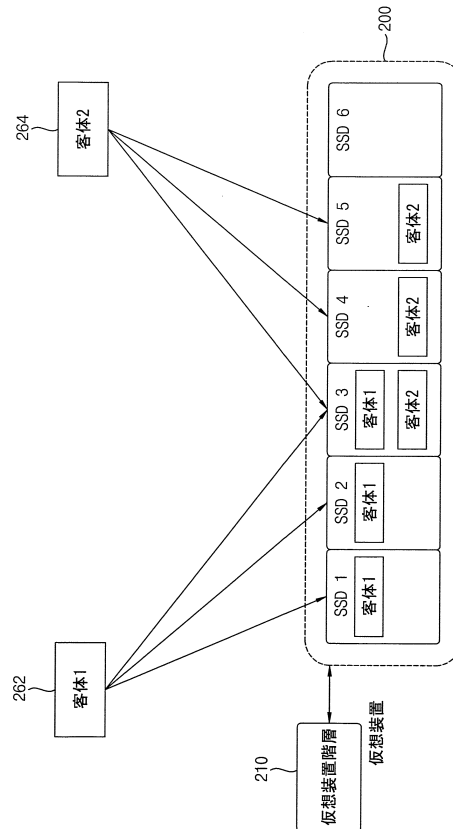
【図 4】



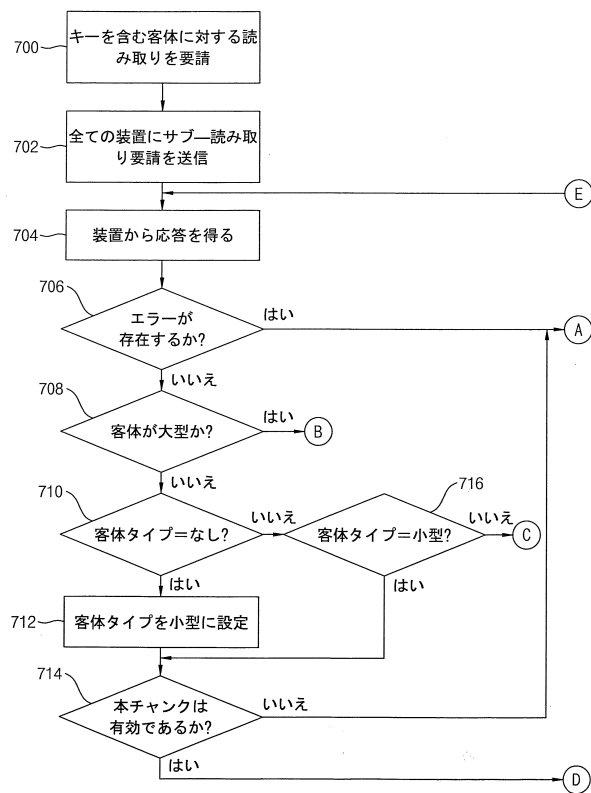
【図 5】



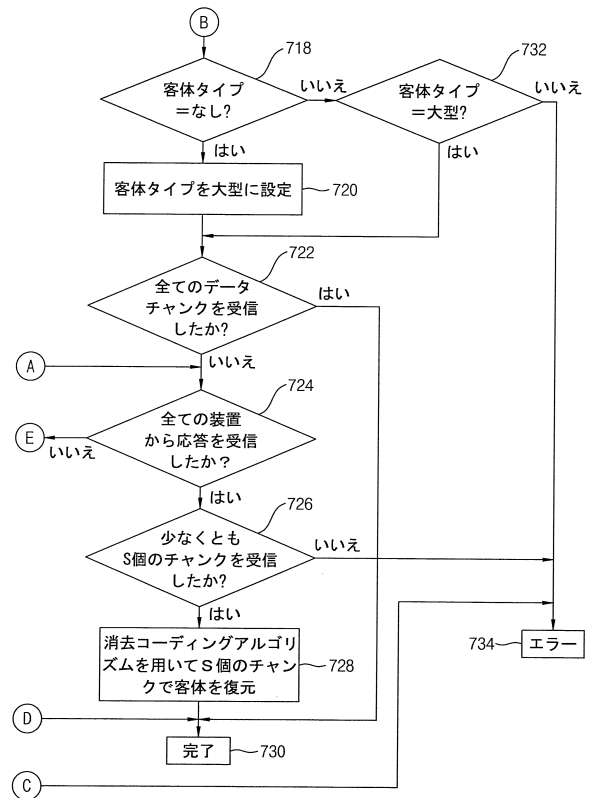
【図 6】



【図 7】



【図 8】



フロントページの続き

- (31)優先権主張番号 62/562,219
(32)優先日 平成29年9月22日(2017.9.22)
(33)優先権主張国・地域又は機関
米国(US)
(31)優先権主張番号 15/876,028
(32)優先日 平成30年1月19日(2018.1.19)
(33)優先権主張国・地域又は機関
米国(US)

早期審査対象出願

審査官 松平 英

- (56)参考文献 特表2015-519674(JP,A)
特表2015-520588(JP,A)
特表2016-500183(JP,A)
国際公開第2016/137402(WO,A1)
米国特許第08504535(US,B1)
井上陽治,第10回 オブジェクトストレージの宿命、容量効率を高める方法とは?,online
 ,日本,ITmedia,2016年02月10日,[検索日 2021年 5月10日]、インターネット<<https://www.itmedia.co.jp/enterprise/articles/1602/10/news021.html>>
浅沼,HDFS Erasure Codingの紹介とYahoo!JAPANにおける運用
事例,online,日本,Yahoo!JAPAN,2017年03月07日,[検索日 2021年 5月
10日]、インターネット<https://techblog.yahoo.co.jp/infrastructure/hdfs_erasure_coding/>

(58)調査した分野(Int.Cl.,DB名)

G06F 3/06-3/08
12/08-12/128
13/10-13/14
13/38-13/42
16/00-16/958
17/30