

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6255496号
(P6255496)

(45) 発行日 平成29年12月27日 (2017.12.27)

(24) 登録日 平成29年12月8日 (2017.12.8)

(51) Int. Cl.	F I				
G06F 17/30	(2006.01)	G06F	17/30	110C	
G06F 12/00	(2006.01)	G06F	17/30	412	
G06F 11/14	(2006.01)	G06F	12/00	531D	
		G06F	12/00	545A	
		G06F	11/14	664	

請求項の数 15 (全 42 頁)

(21) 出願番号	特願2016-540658 (P2016-540658)	(73) 特許権者	507303550
(86) (22) 出願日	平成26年12月18日 (2014.12.18)		アマゾン・テクノロジーズ・インコーポレ ーテッド
(65) 公表番号	特表2017-504885 (P2017-504885A)		アメリカ合衆国・98108-1226・ ワシントン州・シアトル・パイオーボク ス・81226
(43) 公表日	平成29年2月9日 (2017.2.9)	(74) 代理人	100098394
(86) 国際出願番号	PCT/US2014/071159		弁理士 山川 茂樹
(87) 国際公開番号	W02015/095521	(74) 代理人	100064621
(87) 国際公開日	平成27年6月25日 (2015.6.25)		弁理士 山川 政樹
審査請求日	平成28年6月16日 (2016.6.16)	(72) 発明者	ドンラン, ブライアン・ジュームズ
(31) 優先権主張番号	14/133,522		アメリカ合衆国・98109-5210・ ワシントン州・シアトル・テリー アヴェ ニュー ノース・410
(32) 優先日	平成25年12月18日 (2013.12.18)		
(33) 優先権主張国	米国 (US)		
(31) 優先権主張番号	14/133,575		
(32) 優先日	平成25年12月18日 (2013.12.18)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 ボリュームコーホート内の小ボリュームの照合調整

(57) 【特許請求の範囲】

【請求項1】

1つ以上のデバイス上に組み込まれた格納サービスによって、コーホート内の複数の格納ノード上に組み込まれた複数の小ボリュームの選択されたサブセットに、データオブジェクトから生成されたデータ要素のセットを格納することと、

前記小ボリュームの各々2つにおいて、他の前記小ボリュームのための共通のオブジェクトリストを生成することとあって、前記小ボリュームのうち1つでの前記共通のオブジェクトリストは、前記他の小ボリュームにも格納されているはずである、前記小ボリューム上のデータオブジェクトを示す、前記生成することと、

前記2つの小ボリュームの各々において、前記小ボリュームにおける前記共通のオブジェクトリストについてのハッシュ値を生成することと、

前記2つの小ボリューム上で生成された前記ハッシュ値が一致しないことを判定することと、

前記判定にตอบสนองして、前記小ボリュームの両方に格納されているはずであるが格納されていない1つ以上のデータオブジェクトを識別することと、

を含む、方法。

【請求項2】

前記判定すること及び前記識別することが、前記2つの小ボリュームの各々で行われ、前記方法が、前記格納サービスの照合調整プロセスに、前記識別された1つ以上のデータオブジェクトを報告することをさらに含む、請求項1に記載の方法。

10

20

【請求項 3】

前記方法が、前記格納サービスの照合調整プロセスに、前記生成されたハッシュ値を提供することをさらに含み、前記照合調整プロセスが、前記判定すること及び前記識別することを行う、請求項 1 に記載の方法。

【請求項 4】

前記 2 つの小ボリュームを照合調整して、前記 2 つの小ボリュームが各々前記識別された 1 つ以上のデータオブジェクトから生成されたデータ要素を格納するようにする、前記格納サービスの照合調整プロセスをさらに含む、請求項 1 に記載の方法。

【請求項 5】

オブジェクト冗長化データ要素の前記セットが、オブジェクト冗長化手法によって前記データオブジェクトから生成され、前記オブジェクト冗長化手法が、レプリケーション手法またはイレージャ符号化手法のいずれか一方である、請求項 1 に記載の方法。

10

【請求項 6】

小ボリュームにおいて共通のオブジェクトリストを前記生成することが、前記小ボリューム上に前記データ要素とともに格納されたメタデータによって、前記共通のオブジェクトリストを生成することを含み、前記小ボリューム上の所与のデータ要素についての前記メタデータが、データ要素の前記生成されたセットにおける他のデータ要素が格納される 1 つ以上の他の小ボリュームを示す、請求項 1 に記載の方法。

【請求項 7】

小ボリュームにおいて共通のオブジェクトリストを前記生成することが、前記格納ノード上の前記データオブジェクトのオブジェクト識別子によって、前記共通のオブジェクトリストを生成することを含み、オブジェクト識別子の変換式が、前記それぞれのデータオブジェクトについて生成された他のデータ要素が格納される前記コーホート内の 1 つ以上の位置を示す、請求項 1 に記載の方法。

20

【請求項 8】

小ボリュームにおいて共通のオブジェクトリストについてのハッシュ値を前記生成することが、前記共通のオブジェクトリスト内のオブジェクト識別子からハッシュ木を生成することを含み、前記ハッシュ値が、前記ハッシュ木のルートハッシュである、請求項 1 に記載の方法。

【請求項 9】

1 つ以上のコンピューティングデバイス上に組み込まれた格納サービスによって、オブジェクト冗長化格納システムに格納されるデータオブジェクトを受信することであって、前記オブジェクト冗長化格納システムが、N 個の格納ノードにおよぶコーホートを含む、前記受信することを含み、

30

受信されたデータオブジェクトごとに、

オブジェクト冗長化手法によって、前記データオブジェクトから N 未満である M 個のオブジェクト冗長化データ要素を生成し、

前記コーホート内の前記 N 個の格納ノードの中から M 個の格納ノードを選択し、

前記選択された M 個の格納ノードに前記 M 個のデータ要素を格納し、前記 M 個のデータ要素の 1 つは、前記 M 個の格納ノードの各々に格納され、

40

前記 M 個のデータ要素のうち少なくとも 1 つとともに、前記 M 個のデータ要素の他の 1 つ以上が格納される前記コーホート内の 1 つ以上の位置を示すメタデータを格納し、

前記データ要素の個々に関して、前記メタデータと前記データ要素とが、対応する格納ノードの格納ブロック内に互いとともに格納される、方法。

【請求項 10】

前記データ要素の個々に関して、前記データ要素とともに格納される前記メタデータは、前記データ要素が格納される前記コーホート内の場所を示さない、請求項 9 に記載の方法。

【請求項 11】

前記オブジェクト冗長化手法が、レプリケーション手法またはイレージャコード化手法

50

の一方である、請求項 9 に記載の方法。

【請求項 1 2】

前記選択手法が、前記 M 個の格納ノードの少なくとも 1 つをランダムに選択する手法か、または前記それぞれのデータオブジェクトの識別情報によって、前記 M 個の格納ノードの少なくとも 1 つを選択する手法かの方である、請求項 9 に記載の方法。

【請求項 1 3】

前記格納システムに格納されたデータオブジェクトについての要求を受信することと、
前記コーホート内の前記 N 個の格納ノードのサブセットを選択することと、
前記要求されたデータオブジェクトに対応するデータ要素のための格納ノードの前記選択されたサブセットについて問い合わせることと、
格納ノードの前記問い合わせを受けたサブセットから、前記要求されたデータオブジェクトに対応するデータ要素の少なくとも 1 つを受信することと、
をさらに含む、請求項 9 に記載の方法。

10

【請求項 1 4】

複数の格納ノードと、
1 つ以上のプロセッサと、
オブジェクト冗長化格納システムを組み込むために、前記 1 つ以上のプロセッサの少なくとも 1 つによって実行可能であるプログラム命令を格納するメモリであって、前記オブジェクト冗長化格納システムが、
前記複数の格納ノードの N 個におよぶコーホートを確立し、
前記格納システムに格納されるデータオブジェクトを受信するように構成され、
データオブジェクトごとに、
オブジェクト冗長化手法によって、前記データオブジェクトから N 未満の M 個のオブジェクト冗長化データ要素を生成し、

20

選択手法によって、前記コーホート内の前記 N 個の格納ノードの中から M 個の格納ノードを選択し、
前記選択された M 個の格納ノードに前記 M 個のデータ要素を格納することであって、前記 M 個のデータ要素の 1 つが、前記 M 個の格納ノードの各々に格納される、メモリと、
を含み、

30

前記生成すること、前記選択すること、及び前記格納することが、前記受信されたデータオブジェクトについて生成された前記データ要素を、前記コーホート内の前記 N 個の格納ノードにわたって分配し、前記 N 個の格納ノードの所与の 2 つが、異なるデータ要素のセットを含み、

前記オブジェクト冗長化格納システムがさらに、各データオブジェクトに関して、
前記データオブジェクトについて生成された前記 M 個のデータ要素のうちの少なくとも 1 つに関して、前記 M 個のデータ要素のうちの他の 1 または複数が格納される前記コーホート内の 1 または複数の位置を示すメタデータを生成し、

前記 M 個のデータ要素のうちの前記少なくとも 1 つとともに前記メタデータを格納する、ように構成され、

前記データ要素の個々に関して、前記メタデータと前記データ要素とが、対応する格納ノードの格納ブロック内に互いとともに格納される、システム。

40

【請求項 1 5】

前記オブジェクト冗長化手法が、イレージャ符号化手法であり、前記オブジェクト冗長化データ要素が、前記データオブジェクトから生成されたシャードであり、前記冗長な符号化手法によって所与のデータオブジェクトについて生成された、前記 M 個のシャードの少なくとも R 個のサブセットが、前記それぞれのデータオブジェクトを再生するために必要とされる、請求項 1 4 に記載のシステム。

【発明の詳細な説明】

【技術分野】

【0 0 0 1】

50

通常のデータ格納アプリケーションまたはサービスは、要求を受信して、1つ以上のクライアントの代わりにデータオブジェクトを格納し、当該データオブジェクトを、1つ以上の格納ノードに格納することができる。いくつかのデータ格納サービスは、オブジェクト冗長化格納システムと呼ばれることがあり、冗長手法または方式を用いてデータオブジェクトを格納し、格納されたデータのためのより高いレベルの存続性を提供し得る。たとえば、データ格納サービスは、データオブジェクトを複製して、それを2つ以上の異なる格納ノードまたは位置にわたって格納し、データオブジェクトが、任意の所与の格納ノードまたはノードの組み合わせの故障を切り抜ける可能性を増大させることができる。いくつかのオブジェクト冗長化格納システムでは、各々のレプリカは、オブジェクトデータの正確な写しに相当する必要はない。たとえば、いくつかのオブジェクト冗長化格納システムでは、データオブジェクトは、冗長な符号化手法（たとえば、イレージャ符号化）に従って数多くの部分、すなわち「シャード（shard）」に分割されることがあり、シャードの各々が、異なる格納ノードに格納されることがある。

10

【0002】

データオブジェクトが、多数のノード間で簡単に複製されるシステムでは、データオブジェクトを抽出するために、1つのレプリカのみを抽出する必要がある。しかしながら、冗長な符号化手法、たとえばイレージャ符号化を用いる場合、データオブジェクトは、一般的には、生成されたシャードのうち2つ以上であるが全てより少ないものから再生される。たとえば、データオブジェクトから20個のシャードを生成するイレージャ符号化手法を用いると、データオブジェクトを再生するために、少なくとも10個のシャードが必要とされる場合がある。

20

【図面の簡単な説明】

【0003】

【図1】データオブジェクトが格納位置の集合体群に格納されるオブジェクト冗長化格納システムを図示し、ここでは、所与のデータオブジェクトから生成された1つのデータ要素が、群中の各々の位置に格納される。

【図2A - 2B】少なくともいくつかの実施形態による、ボリュームコーホートを組み込んでいるオブジェクト冗長化格納システムを図示する。

【図3A - 3B】少なくともいくつかの実施形態によるコーホートの例を図示する。

【図4A - 4C】実施形態による、コーホート内でのデータ要素の、メタデータとのタグ付けを図示する。

30

【図5】少なくともいくつかの実施形態による、データオブジェクトを作成して、オブジェクト冗長な格納システム内のコーホートに格納するための方法の、高レベルのフローチャートである。

【図6】少なくともいくつかの実施形態による、レプリケーション手法によってデータオブジェクトが格納されるコーホートからデータオブジェクトを抽出するための方法の、高レベルのフローチャートである。

【図7】少なくともいくつかの実施形態による、冗長な符号化手法によってデータオブジェクトが格納されるコーホートからデータオブジェクトを抽出するための方法の高レベルのフローチャートである。

40

【図8】少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームの比較を行うための方法をグラフィカルに図示する。

【図9】少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームの比較を行うための代替の方法をグラフィカルに図示する。

【図10】少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームの比較を行うための方法のフローチャートである。

【図11A - 11C】少なくともいくつかの実施形態による、オブジェクト冗長化格納シ

50

システム上での照合調整プロセスの一部として、コーホートの小ボリュームの比較を行うための代替の方法のフローチャートである。

【図12】少なくともいくつかの実施形態による、ハッシュ木の例を図示する。

【図13】いくつかの実施形態で用いられ得るコンピュータシステムの例を図示するブロック図である。

【0004】

いくつかの実施形態の例示及び例証となる図面の例示を目的として、実施形態が本明細書に記載されているが、当業者においては、記載された実施形態または図面に限定されないことが認識されよう。図面及びその詳細な説明は、実施形態を開示された特定の形態に限定することを意図するものではなく、それどころか、添付の特許請求の範囲によって画定された本質及び範囲に含まれるすべての変形、等価物及び変形をを対象とすることが意図されることが理解されるべきである。本明細書で用いられた表題は、構成的な目的のためのみのものであり、明細書または特許請求の範囲の範囲を限定するために用いられることを意図しない。本出願全体で用いられる場合、用語「may」は、必須の意味（すなわち、mustを意味する）というよりは、許容的な意味（すなわち、可能性を有することを意味する）で用いられる。同様に、用語「include」、「including」、及び「includes」は、含むが限定されないことを意味する。

【発明を実施するための形態】

【0005】

オブジェクト冗長化格納システムにおいてボリュームコーホートを設けるための方法及び機器の様々な実施形態が記載される。オブジェクト冗長化格納システムでは、データオブジェクトは、レプリケーション手法によって複製されることができ、レプリカは、2つ以上の異なる格納位置に格納され得る。あるいはまたはそれに加えて、オブジェクト冗長化格納システムで冗長な符号化手法、たとえばイレージャ符号化を用いて、データオブジェクトから多数のシャードを生成してもよく、多数の異なる格納位置にわたってシャードを格納してもよい。本明細書の目的において、レプリケーション手法及び冗長な符号化手法は、オブジェクト冗長化手法と総称されることがある。本明細書の目的のために、データオブジェクトのレプリカ及びシャードは、オブジェクト冗長化データ要素、または単にデータ要素と総称されることがあり、この場合1つのデータ要素は、所与のデータオブジェクトの1つのレプリカまたは1つのシャードに対応する。さらに、本明細書で用いられる場合、データオブジェクトは、オブジェクト冗長化格納システム内の位置に格納され得る任意のタイプのデータであってもよく、任意のサイズであってもよいことに留意すべきである。また、データオブジェクトは、単一のデータ要素または単一のタイプのデータ、同じタイプのまたは異なるタイプのデータ要素の集合、またさらにはデータオブジェクトの集合を含んでもよい。

【0006】

オブジェクト冗長化格納システム内において、永続データのこれらデータ要素（レプリカまたはシャード）を検索するための従来手法は、たとえばランダム選択、または格納システム内のすべての格納ノードの中から、所与のデータオブジェクトについて生成されたデータ要素の位置のセットを選択するための他の何らかの手法を用いて、データ要素のための位置を、データオブジェクトごとに独立して選択することである。しかしながら、この方法は、一般的には、オブジェクトごとにデータ要素の格納位置を追跡するための大量のメタデータを伴い、また抽出のためにデータオブジェクトを検索する場合、及び/または故障した格納デバイスまたはノードを回復させる場合に、大量のオーバヘッドを伴うことがある。

【0007】

追跡に要するメタデータ量を低減させ、データオブジェクトの抽出におけるオーバヘッドを低減させ得る、上記手法の代替として、格納デバイス群または格納デバイスの一部を、格納システム内に作成してもよい。そして、データオブジェクトは、群に割り当てられることができ、1つのデータ要素が、群中の各メンバーデバイス（またはデバイスの一部

10

20

30

40

50

)に格納された所与のデータオブジェクトから生成される。データオブジェクトを検索するために、データオブジェクトが格納された群が最初に検索され、その後は、群中の当該位置から、データオブジェクトを抽出し得る。この手法を実施する例示の格納システム100が、図1に図示される。格納システム100は、多数の格納ノード110と、それを介して1つ以上のクライアント190が格納システム100にデータオブジェクトを格納し、かつそこからデータオブジェクトを抽出し得るインターフェース(たとえば、アプリケーションプログラミングインターフェース(API))を提供する格納サービス150とを含み得る。図1に示されるように、格納ノード110A~110mの群または格納ノード110A~110mの一部は、ボリューム102を構成するかまたは包含し得る。図1に図示されたもの以外に、格納システム110内にさらなる格納ノード110、格納ノード110の他の群、及び他のボリューム102が存在し得ることに留意すべきである。

10

【0008】

広義には、ボリュームは、本明細書で用いられる場合、データオブジェクトの集合であり、2つ以上の物理的格納ノード110間に分散され得る仮想的格納デバイスと見なされ得る。たとえば、ボリューム102は、図1に示されるように、ノード110A~110mにわたって分散される。ボリューム102は、多数の小ボリューム120を構成すると見なされ得る。小ボリューム120は、一般的には、格納ノード110上の連続した格納ブロックであってもよく、各小ボリューム120は、数千または数百万のデータ要素122を包含し得る。各小ボリューム120は、単一の格納ノード120上に存在するが、しかしながら、ボリューム102の小ボリューム120は、通常は各々異なる格納ノード110上に存在し得る。図1には示されないが、2つ以上の異なるボリューム102から2つ以上の小ボリューム120は、同じ格納ノード110上に共存してもよい。さらに、格納システム100内の2つ以上のボリューム102は、同じ格納ノード110、格納ノード110の異なる群、または格納ノード110の重複する群に及んでもよい。ボリューム102は、多数の小ボリューム120で構成され得る一方で、格納サービス150インターフェースは、ボリューム102を、単一の仮想格納デバイスまたはシステムとしてクライアント190に提示し得る。

20

【0009】

図1に図示された格納システム100では、ボリューム102の各小ボリューム120上のデータ要素122は、同じデータオブジェクトのセットに対応でき、各データオブジェクトは、各小ボリューム120に格納されたデータ要素122(シャードまたはレプリカ)を有する。言い換えると、各小ボリューム120は、ボリューム102に格納されたデータオブジェクトごとにデータ要素122(シャードまたはレプリカ)を含む。従って、小ボリューム120は、ボリューム102の単一のレプリカ、または「シャード」と見なされ得る。これを図示するための図1を用いると、レプリケーション手法を用いて、クライアント190から受信されたデータオブジェクトを永続的に格納すると、各データオブジェクトのレプリカはその後、ボリューム102の小ボリューム120A~120mの各々にデータ要素122として格納され、各小ボリューム120上のデータ要素122のセットは、一般的には同一であるはずである。あるいは、冗長な符号化手法(たとえば、イレージャ符号化)を用いて、クライアント190から受信されたデータオブジェクトを永続的に格納すると、その後各データオブジェクトからM個のシャードが生成され、シャードの異なる1つが、ボリューム102の小ボリューム120A~120mの各々に、データ要素122として格納される。従って、ボリューム102の小ボリューム120内のデータ要素122は、一般的には、同じデータオブジェクトのセットにすべて対応するはずである。

30

40

【0010】

ノード110全体にわたる小ボリューム120でデータオブジェクトが複製される格納システム100では、データオブジェクトを抽出するために、ボリューム102から1つのレプリカのみを抽出することが必要である。しかしながら、格納システム100で冗長な符号化手法、たとえばイレージャ符号化を用いる場合、データオブジェクトは、一般的

50

には、小ボリューム120に格納された生成されたシャードのうち2つ以上であるが、すべてより少ないものから再生され得る。たとえば、図1に図示されるような、データオブジェクトからM個のシャードを生成し、シャードの異なる1つを、小ボリューム120A~120mの各々にデータ要素122として格納するイレージャ符号化手法を用いると、対応するデータオブジェクトを再生するために、m個の小ボリューム120のサブセットのいくつかからシャードを抽出することを要する。非限定的な例として、イレージャ符号化方式は、m個のシャードが再生されて、データオブジェクトを再生するためにシャードの半分を要する場合に使用され得る。従って、データオブジェクトを再生するために要するシャードの(最少)数は、 $m/2$ であり得る。

【0011】

図1を参照して上述したような、オブジェクト冗長化格納システムにデータ要素を格納するための手法は、データ要素(レプリカまたはシャード)が、データ要素ごとに選択された位置に独立して格納される上記の第1の手法と比較した場合、データオブジェクトを追跡するために要するメタデータ量を低減させることができ、所与のデータオブジェクトを抽出するために要するオーバーヘッドを低減させることができるが、図1に図示された手法は、相関エラーと称される場合がある状況の可能性をもたらす。

【0012】

第1または第2の手法を用いると、格納システム内の単一の(あるいは数個の)格納ノードがダウンした場合、そのノード上に格納されたデータオブジェクトはその後、一般的には格納システム内の他の格納ノード上に格納されたデータ要素(レプリカまたはシャード)から回復させることができる。しかしながら、格納システムでのノードの多重故障は、結果としていくらかのデータを損失することになるかもしれない。第1の手法を用いると、格納システムにおけるノードの多重故障は、結果として、個々のデータオブジェクトのいくつかを損失することとなり、格納システムから回復できなくなる場合がある。たとえば、各データ要素(シャードまたはレプリカ)が、格納システム内の任意の多くのノードから選択された4つのノードの各々に独立して格納されると、その後格納システム内の4つのノードの故障によって、結果として別個のデータオブジェクトの比較的小さなサブセットのいくつかを損失することになる場合がある。

【0013】

しかしながら、第2の手法を用いると、ノードの多重故障は、場合によっては、データオブジェクトの全ボリュームを損失する結果となるかもしれない。第2の手法を用いると、一般的には、ノードの多重故障に起因して、ボリュームから何らかのデータオブジェクトを損失すると、その後ボリュームからデータオブジェクトのすべてが失われる。このことは、相関エラーと称されるものである。ボリューム内の小ボリュームにm個のレプリカが格納され、その後格納システム内の格納ノードで損失するレプリケーション方式を用いることは、結果として特定のm個の格納ノードに及ぶ1つ以上のボリューム内に格納されたデータオブジェクトのすべてを損失することになるかもしれない。たとえばデータオブジェクトのM個のシャードがボリューム内の小ボリュームに格納され、データオブジェクトを再生するためにシャードのサブセットを要し、その後格納システム内のデータオブジェクトを再生するために要するm個の格納ノードの分数よりも1多く(たとえば、データオブジェクトを再生するためにm個のシャードの $1/2$ を要する場合、 $(m/2)+1$ 個の格納ノード)を損失するイレージャ符号化等の冗長な符号化手法を用いることは、結果として、特定のm個の格納ノードに及ぶ1つ以上のボリュームに格納されたデータオブジェクトのすべてを損失することになるかもしれない。

【0014】

上述された2つの手法を用いた個々のデータオブジェクトの平均故障間隔(Mean Time Between Failure: MTBF)は、類似しているかまたは同じであるかもしれない。しかし、結果として相関エラーとなり、そのため全ボリュームに影響を及ぼす第2の手法を用いたノードの多重故障は、格納システムのクライアントにとって、一般的には相関していないデータオブジェクトを損失することがある第1の手法を用

10

20

30

40

50

いたノードの多重故障よりもはるかに目立ち、ましてや望ましくない場合がある。

【 0 0 1 5 】

(オブジェクト冗長化格納システムにおけるボリュームコーホート)

オブジェクト冗長化格納システムにボリュームコーホートを設けるための方法及び機器の実施形態が説明され、これらは上記2つの手法の利点を提供する一方で、2つの手法の課題を低減させ得る。ボリュームコーホート手法または方法の実施形態が説明され、これは、オブジェクト冗長化格納システムで実施することができ、第1の手法と比較して、データオブジェクトを追跡するために要するメタデータ量及び/または所与のデータオブジェクトを抽出するために要するオーバーヘッドを低減させると同時に、第2の手法の関連エラーの課題を低減させるかまたは取り除くことができる。

10

【 0 0 1 6 】

実施形態においては、ボリュームコーホート、または単にコーホートは、オブジェクト冗長化格納システム内の格納ノードのセットまたは群に及んで作成され得る。図1に示されたような格納システムについて説明されたボリュームと同様に、オブジェクト冗長化格納システム内のコーホートは、2つ以上の物理的格納ノードにわたって分散し得る仮想的格納デバイスと見なされ得る。しかしながら、図1を参照して説明されたボリュームとは異なり、オブジェクト冗長化手法によってコーホートに格納された所与のデータオブジェクトは、コーホート内の格納ノードのサブセットのみに及ぶ。このように、コーホート内にN個の格納ノードがあると、任意の所与のデータオブジェクトのデータ要素(レプリカまたはシャード)が、コーホート内の格納ノードのM個のみに格納され、ここで、MはN未満である。図1に図示されたような格納システムについて説明されたボリュームと同様に、コーホートの小ボリュームは、コーホート内のN個の格納ノードの各々に位置する。しかしながら、図1に図示されたような格納システムの場合とは異なり、コーホート小ボリュームは同一ではない。すなわち、コーホート小ボリュームは、同じデータオブジェクトのセットについてデータ要素のセットを各々包含することがないが、これは、コーホートに格納された各データオブジェクトが、コーホートの各小ボリュームに格納されたデータ要素(シャードまたはレプリカ)を有していないためである。

20

【 0 0 1 7 】

図2A及び2Bは、少なくともいくつかの実施形態による、ボリュームコーホートを実施するオブジェクト冗長化格納システムを示す。図2Aに図示されるように、格納システム200は、多数の格納ノード210と、それを介して1つ以上のクライアント290が格納システム200にデータオブジェクトを格納し、そこからデータオブジェクトを抽出するインターフェース(たとえば、アプリケーションプログラミングインターフェース(API))を提供する格納サービス250とを含み得る。格納システム200は、一般的には、クライアントに対するオブジェクト冗長な格納を提供する任意の格納システムであってもよいことに留意すべきである。たとえば、格納システム200は、1つ以上のクライアントデバイスに接続されたローカルな格納システム、ローカルネットワークに接続されてローカルネットワーク上の多数のクライアントにアクセス可能であるネットワークベースの格納システム200、またはプロバイダネットワーク上に組み込まれ、多数のクライアントに対する遠隔仮想化格納サービスとして設けられて、APIに従って、及びネットワーク、たとえばインターネットを介してクライアントにアクセス可能である遠隔仮想化格納システムであり得る。

30

40

【 0 0 1 8 】

図2Aに示されるように、コーホート202は、格納システム200内の多数の格納ノード210A~210Nに及び得る。図2Aに示されるように、コーホート202は、ノード210A~210Nにわたって分散する仮想的格納デバイスと見なされ得る。コーホート202は、多数の小ボリューム220A~220Nを含んでもよく、ここで各小ボリューム220は、格納ノード210上の連続した格納ブロックであってもよく、各小ボリューム220は、格納システム200に格納されたデータオブジェクトのデータ要素222(シャードまたはレプリカ)を格納する。各小ボリューム220は、単一の格納ノード

50

220上にあるが、しかしながら、通常は、コーホート202の小ボリューム220は、異なる格納ノード210上にある。さらに、コーホート202は、多数の小ボリューム220で構成され得るが、一方で格納サービス250インターフェースは、クライアント290に、コーホートを単一の仮想的格納デバイスまたはシステムとして提示されることがある。

【0019】

図2Aは、簡潔化のために、N個の格納ノード210にわたって分散する1つのコーホート202だけを示すことに留意すべきである。しかしながら、図2Aの例に示されたものよりも多くの格納ノード210及び格納システム内の多くのコーホート202があってもよい。図2Bに示されるように、格納システム200内の2つ以上のコーホート202は、同じ格納ノード210(図2Bのコーホート202A及び202B)、格納ノード210の異なる群(図2Bのコーホート202A及び202D)、または格納ノード210の重複する群(図2Bのコーホート202A及び202C)に及んでもよい。このように、2つ以上の異なるコーホート202からの2つ以上の小ボリューム220は、同じ格納ノード210上に共存し得る。たとえば、図2Bの格納ノード210Cは、コーホート202Aの小ボリューム220A3、コーホート202Bの小ボリューム220B3、及びコーホート202Cの小ボリューム220C1を含む。このように、格納システム200内の各格納ノード210は、多数のコーホート202に関与してもよく、所与の格納ノード210が関与する2つ以上のコーホート202は、異なる数のノード210(すなわち、それぞれのコーホート202に関与する異なる格納ノード210のセット)を有していてもよい。

【0020】

図2Aを再度参照すると、コーホート202内の所与のN個の格納ノード210を所与とすると、任意の所与のデータオブジェクトのデータ要素222(レプリカまたはシャード)は、コーホート202内の格納ノード210上のM個の小ボリューム220にのみ格納され、ここで、MはN未満である。さらに、データオブジェクトごとのデータ要素222が格納される特定のM個の小ボリューム220は、選択手法(たとえば、ランダム選択手法)によって決定されてもよい。この場合、M個の小ボリューム220がN個の小ボリューム220すべての中から選択されて、データ要素222が、N個の小ボリューム220すべての間に分配される。言い換えると、第1のデータオブジェクトについてのデータ要素222が格納されるM個の小ボリューム220のセットは、一般的には、第2のデータオブジェクトについてのデータ要素222が格納される異なるM個の小ボリューム220のセットであることがある(しかし、必ずしもそうではない)。

【0021】

このように、図2Aに示された例示の格納システム200では、図1に図示された例示の格納システム内のボリューム102とは異なり、コーホート202の各小ボリューム220上のデータ要素222は、同じデータオブジェクトのセットに対応しない。これは、所与のデータオブジェクトについてのデータ要素222が、N個の小ボリューム220のサブセットのみに格納されるためである。たとえば、図2Aでは、データ要素222Aは、小ボリューム220A及び220Nに格納されているが、小ボリューム220Bには格納されていない。そして、データ要素222Bは、小ボリューム220A及び220Bに格納されているが、小ボリューム220Nには格納されていない。

【0022】

図2Aに示されたような格納システム200では、レプリケーション手法を用いて、クライアント290から受信されたデータオブジェクトを永続的に格納すると、その後コーホート202のN個の小ボリュームからM個の小ボリュームが選択され、データオブジェクトのレプリカは、データ要素222としてM個の小ボリュームの各々に格納される。あるいは、冗長な符号化手法(たとえば、イレージャ符号化)を用いて、クライアント290から受信されたデータオブジェクトを永続的に格納すると、その後、各データオブジェクトからM個のシャードが生成され、コーホート202のN個の小ボリュームからM個の

10

20

30

40

50

小ボリュームが選択され、シャードのうち異なる1つが、選択されたM個の小ボリュームの各々にデータ要素222として格納される。冗長な符号化手法、たとえばイレージャ符号化が用いられるシステムでは、当該手法によって生成されたシャードの全数が、Mを決定し得ることに留意すべきである。

【0023】

図2Aに示されたような、データオブジェクトが、コーホート202のN個の小ボリューム220のM個にわたって複製される格納システム200では、データオブジェクトを抽出するために、コーホート202から1つのレプリカのみが抽出されることを要する。しかしながら、格納システム200で冗長な符号化手法、たとえばイレージャ符号化を用いた場合、データオブジェクトは、一般的には、小ボリューム220に格納された、生成されたシャードのうち2つ以上であるがすべてより少ないものから再生されることがある。たとえば、図2Aに示されるように、データオブジェクトからM個のシャードを生成し、異なる1つのシャードを、データ要素222として、選択されたM個の小ボリューム220の各々に格納するイレージャ符号化手法を用いると、対応するデータオブジェクトを再生するために、M個の小ボリューム220のサブセットのいくつかからシャードを抽出することを要する。非限定的な例として、M個のシャードが生成され、データオブジェクトを再生するためにシャードの半分を要するイレージャ符号化方式を用いると、それによって、データオブジェクトを再生するために要するシャードの(最少)数は、 $M/2$ である場合がある。本明細書では、たとえばイレージャ符号化等の冗長な符号化方式を用いてデータオブジェクトを再生するために要するシャードの数を、Rと称することがあり、よって、本例では $R = M/2$ である。単なる1つの具体的な例として、イレージャ符号化方式を用いてもよく、その後データオブジェクトについて20個のシャードが生成され($M = 20$)、データオブジェクトを生成するために10個のシャードを必要とする($R = 10$)。別の例として、イレージャコード化方式を用いてもよく、その後データオブジェクトについて22個のシャードが生成されることがあり($M = 22$)、データオブジェクトを再生するために11個のシャードが必要とされる($R = 11$)。

【0024】

図2Aに示されたような格納システム200内のコーホート202のいくつかの実施では、Mは、Nの2分の1として選択されることがあり、またはNはMの2倍として選択されることがある。たとえば、例示の実施では、 $M = 20$ 及び $N = 40$ 、または $M = 22$ 及び $N = 36$ である。これらの実施では、各データオブジェクトは、コーホート202内の小ボリューム220のちょうど半分に格納される。しかしながら、実施形態において、M及び/またはNについて他の値、及びN対Mについて他の比率を用いてもよい。いくつかの非限定的な例を提供するために、Nは、 $4M$ (たとえば、 $M = 10$ に対して $N = 40$)、または $M = 20$ に対して $N = 22$ として選択されてもよい。

【0025】

M及びNの選択、及びN対Mの比率は、実装固有であってもよく、また用いられた特定の冗長な符号化方式、利用可能な格納ノード数、及び性能オーバーヘッドとデータ損失保護との間のトレードオフを含むがこれらに限定されない因子に基づいてもよい。性能オーバーヘッドとデータ損失保護との間のトレードオフに関しては、より高いN対Mの比率(たとえば、 $3:1$ 、 $4:1$ 、またはそれ以上)は、所与のイベントにおいて損失する予想データオブジェクト数を低下させる一方でさらなる格納ノードが含まれるためにオーバーヘッドを増大させる場合があるが、一方で、より低いN対Mの比率(たとえば、 $2:1$ 、 $3:2$ 、またはそれ以下)は、より少ないオーバーヘッドで、所与のイベントにおいて損失するであろう予想データオブジェクト数を増大させるかもしれないことに留意すべきである。所与の実施のために、当該実施のために、許容可能な数のオーバーヘッドによる適切なレベルのリスク低減を提供するM及びNの値が決定され得る。本明細書において、いくつかの二項計算が後述され、これらは、M及びNの値を評価し、また場合によっては選択する際に用いられ得る。

【0026】

図 2 A に示されたような格納システム 2 0 0 内のボリュームコーホート 2 0 2 が、図 1 に示されたような、各データオブジェクトがボリューム 1 0 2 内の小ボリューム 1 2 0 のすべてにわたって格納される格納システム 1 0 0 に生じ得る相関エラーの問題をどのようにして克服し得るかを例証するために、図 3 A には簡潔な例示のコーホート 3 0 2 A が設けられている。コーホート 3 0 2 A では、 $M = 2$ 及び $N = 4$ であり、4 つの小ボリューム 3 2 0 A ~ 3 2 0 D があり、所与のデータオブジェクトについて 2 つのデータ要素 2 2 2 (レプリカまたはシャード) が作成され、選択手法によって決定された小ボリューム 3 2 0 のいずれか 2 つに格納され得る。図 3 A は、小ボリューム 3 2 0 A 及び 3 2 0 B に格納されたデータ要素 2 2 2 A (データオブジェクトから作成されたデータオブジェクトまたはシャードのレプリカであってもよい) を示す。それらのうちいずれか 2 つが選択されることとなる 4 つの小ボリュームを所与として、データオブジェクトについてのデータ要素 3 2 2 が格納され得る 2 つの小ボリューム 3 2 0 の 6 個の可能な順不同の組み合わせ C がある。

10

[A B、A C、A D、B C、B D、C D]

【 0 0 2 7 】

順序は重要ではないことに留意すべきであり、言い換えると、 $A B = B A$ 及び $C D = D C$ である。コーホート 3 0 2 A 内の 4 つの小ボリューム 3 2 0 間のランダムなデータ要素の分布を前提とすると、複製されたデータオブジェクトに対して、小ボリューム 3 2 0 の上位機器となる 4 つの格納ノードのうち 2 つが故障したとすると、所与のデータオブジェクトが損失する確率は、6 個のうち 1 つである。言い換えると、小ボリュームの上位機器となる 4 つの格納ノード 3 2 0 のうち 2 つを損失した場合、レプリケーション手法によってコーホート 3 0 2 A に格納された 6 個のデータオブジェクトの内およそ 1 つのみが損失するかもしれない。

20

【 0 0 2 8 】

総括すると、 M 及び N の値からの、データオブジェクトが格納され得るコーホート内の小ボリュームの順不同の組み合わせ C の数は、二項係数

【 0 0 2 9 】

【数 1】

$$\binom{n}{k}$$

30

【 0 0 3 0 】

によって与えられ、これは「 n は k に組み合わせる」と解釈される。

【 0 0 3 1 】

【数 2】

$$C = \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

【 0 0 3 2 】

ここで、 $n!$ は分数関数であり、 $k = M$ 、 $n = N$ である。例えば、 $M = 2$ 、 $N = 4$ の場合、図 3 A に示すとおり、データオブジェクトが格納され得る、

40

【 0 0 3 3 】

【数 3】

$$\binom{4}{2} = 6 \text{ 個}$$

【 0 0 3 4 】

の小ボリュームの組み合わせがある。

50

【 0 0 3 5 】

図 3 B は、例示のコーホート 3 0 2 B を示し、8 個の小ボリューム 3 2 0 A ~ 3 2 0 H があり、各データオブジェクトについてのデータ要素 3 2 2 は、小ボリューム 3 2 0 のうち 4 つに格納される。このように、図 3 B のコーホート 3 0 2 B では、 $M = 4$ 及び $N = 8$ である。二項係数を適用して、組み合わせの数を求める。

【 0 0 3 6 】

【数 4】

$$C = \binom{8}{4} = 70$$

10

【 0 0 3 7 】

このように、コーホート 3 0 2 B 内の、データオブジェクトが格納され得る 4 つの小ボリューム / 格納ノードには 7 0 個の可能な組み合わせがある。そして、コーホート 3 0 2 B 内の 8 個の格納ノードのうち 4 つの任意の組み合わせを損失することは、結果として、レプリケーションを用いた場合に 7 0 個のデータオブジェクトのうちおよそ 1 つを損失することになるかもしれない。5 つの例示のデータオブジェクトからのデータ要素 3 2 2 における、コーホート 3 0 2 B 内の小ボリューム 3 2 0 の異なる組み合わせが、図 3 B に示される。たとえば、第 1 のデータオブジェクトについてのデータ要素 3 2 2 A は、3 2 0 A、3 2 0 B、3 2 0 C、及び 3 2 0 D に格納され、一方で第 2 のデータオブジェクトについて

20

【 0 0 3 8 】

異なる値が M 及び N に与えられた場合の、順不同の組み合わせ C のいくつかの他の非限定的な例が、以下に提示される。

【 0 0 3 9 】

【数 5】

$$M = 10, N = 20: \quad C = \binom{20}{10} = 184,756$$

$$M = 18, N = 36: \quad C = \binom{36}{18} = 9,075,135,300$$

$$M = 18, N = 20: \quad C = \binom{20}{18} = 190$$

30

【 0 0 4 0 】

このように、図 2 A に示されたような格納システム 2 0 0 内のボリュームコーホート 2 0 2 は、図 1 に示されたような格納システム 1 0 0 で生じ得る関連エラーの問題をどのようにして克服し得るかの例として、 $M = 10$ 及び $N = 20$ であり、 $C = 184,756$ であるコーホートにおいて、データオブジェクトの M 個のレプリカがコーホート内の M 個の小ボリュームに格納され、 M 個の小ボリュームの異なる組み合わせが、選択手法によって N 個の小ボリュームから選択され、その後格納システム内の M 個の格納ノードを損失するレプリケーション方式を用いることは、結果として特定の M 個の格納ノードを含む N 個の格納ノードに及びコーホートに格納された 1 8 4 , 7 5 6 個のデータオブジェクトのうちおよそ 1 つのみを損失することになり得る。さらに、コーホートから M 個の格納ノードよりも少数が損失すると、その後、一般的には複製されたデータオブジェクトは損失されないが、これは、残余の小ボリュームのうち少なくとも 1 つが、任意の所与のデータオブジェクトについてのレプリカを含んでいるはずであり、データオブジェクトが、単一のレプリカから回復され得るためである。

40

50

【 0 0 4 1 】

一定数の格納ノードの故障を所与として、データオブジェクトの数を求めるための計算は、データオブジェクトのM個のシャードが、選択手法によってコーホート内のN個の小ボリュームから選択されたM個の小ボリュームに格納され、及びデータオブジェクトを再生するためにシャードのサブセットRを要する、イレージャ符号化等の冗長な符号化手法を用いた場合のコーホートにおいて異なる。そのようなシステムでは、 $(M - R) + 1$ シャードが損失すると、データオブジェクトが損失する可能性がある。簡易な例では、 $R = 2$ 、 $M = 4$ 、及び $N = 8$ であるコーホートでは、 $(4 - 2) + 1 = 3$ シャードが損失すると、データオブジェクトが損失する可能性がある。このように、各データオブジェクトが、コーホート内の8個の小ボリュームの内4つに(シャードとして)格納されたとしても、任意の3つの小ボリュームを損失することは、結果として、損失した小ボリューム3つすべてにシャードを回らずも格納させたいずれかのデータオブジェクトを損失することになるかもしれない。このように、

10

【 0 0 4 2 】

【数6】

$$C = \binom{N}{M} = \binom{8}{4} = 70$$

【 0 0 4 3 】

に代えて、レプリケーションを用いてデータオブジェクトを格納するケースのように、以下のような計算が成り立つ。

20

【 0 0 4 4 】

【数7】

$$C = \binom{N}{(M-R)+1} = \binom{8}{3} = 56$$

【 0 0 4 5 】

言い換えると、本例の $R = 2$ 、 $M = 4$ 、及び $N = 8$ であるコーホート内の任意の3つの小ボリュームを損失することは、結果として、冗長な符号化方式によって、シャードとしてコーホートに格納された56個のデータオブジェクトのうち、およそ1つを損失することになる。

30

【 0 0 4 6 】

本例のコーホートにおいて任意の4つの小ボリュームの損失を所与として、4つの小ボリュームのセット内の3つの小ボリュームの、4つの可能な順不同の組み合わせがあるため、

【 0 0 4 7 】

【数8】

$$C = \binom{4}{3} = 4$$

40

となる。

【 0 0 4 8 】

このように、任意の4つの小ボリュームの故障は、結果として、冗長な符号化手法によってコーホートに格納された56個のデータオブジェクトの内およそ4つ、または14個のうち1つを損失することになるかもしれない。

【 0 0 4 9 】

別の例として、 $R = 5$ 、 $M = 10$ 、及び $N = 20$ であるコーホートでは、 $(10 - 5)$

50

+ 1 = 6 個のシャードを損失した場合にデータオブジェクトを損失する場合があります、計算は以下の通りである。

【 0 0 5 0 】

【 数 9 】

$$C = \binom{20}{6} = 38,760$$

【 0 0 5 1 】

言い換えると、本例のコーホートにおいて任意の 6 個の小ボリュームが損失することは、結果として、冗長な符号化方式によって、シャードとしてコーホートに格納された 38,760 個のデータオブジェクトのうちおよそ 1 つを損失することになる。本例のコーホートにおける任意の 10 個の小ボリュームの損失を所与として、10 個の小ボリュームのセットにおいて 6 個の小ボリュームの 210 個の可能な順不同の組み合わせがあるため、

【 0 0 5 2 】

【 数 10 】

$$C = \binom{10}{6} = 210$$

となる。

【 0 0 5 3 】

このように、任意の 10 個小ボリュームの故障は、結果として、冗長なコード化手法によってコーホートに格納された 38,760 個のデータオブジェクトのうちおよそ 210、または 185 個の内およそ 1 つを損失することになるかもしれない。

【 0 0 5 4 】

イレージャ符号化等の冗長な符号化方式が用いられるコーホートから、 $(M - R) + 1$ 個よりも少ない格納ノードを損失すると、その後、一般的にはデータオブジェクトを損失しないが、これはコーホート内の残余の小ボリュームが、任意の所与のデータオブジェクトを再生するために十分なシャードを含んでいるはずであるためであることに留意すべきである。

【 0 0 5 5 】

上記の計算は、たとえば特定の冗長な符号化方式による特定のコーホート構成における予想損失率を求めるために用いてもよく、それによってボリュームコーホートを採用しているオブジェクト冗長な格納システムでの R、M、及び N の異なる値のリスクを評価してもよい。その結果は、たとえば性能オーバーヘッドとデータ損失保護との間のトレードオフの評価、冗長な符号化方式の選択、及び特定のボリュームコーホートの実施における R、M、及び N の値の選択に用いてもよい。

【 0 0 5 6 】

(小ボリュームのセットの選択)

実施形態においては、データオブジェクトごとのデータ要素が格納されるコーホート内の特定の M 個の小ボリュームは、M 個の小ボリュームを N 個の小ボリュームすべての中から選択して、データ要素が N 個の小ボリュームすべての間に概ね均等に分配されるようにする選択手法によって決定され得る。言い換えると、第 1 のデータオブジェクトについてのデータ要素が格納される M 個の小ボリュームのセットは、一般的には、第 2 のデータオブジェクトについてのデータ要素が格納される M 個の小ボリュームの異なるセットである場合がある (しかし、必ずしもそうではない)。いくつかの実施形態では、ランダムなまたは擬似ランダムな手法を用いて、所与のデータオブジェクトについてのデータ要素格納されることになっている小ボリュームの特定のセットを選択し得る。しかしながら、いくつかの実施形態では、データオブジェクトを格納するコーホートから M 個の小ボリューム

10

20

30

40

50

を選択する場合には、1つ以上の要素が考慮され得る。これらの要素は、コーホート内の格納ノード上の1つ以上の利用可能な格納スペース、格納ノード現在の利用可能度、及び格納ノードに対する帯域上の問題点を含み得るが、これらに限定されない。

【0057】

いくつかの実施形態では、所与のデータオブジェクトについてのデータ要素が格納されることになっているコーホート内の特定のM個の小ボリュームを選択することは、データオブジェクトの識別情報に、または識別情報の変換式(たとえばハッシュ)に、少なくとも部分的に基づいてもよい。データオブジェクトの識別情報の一例として、少なくともいくつかの実施形態では、各データオブジェクトは、格納システム内のデータオブジェクトを一意的に識別し得るオブジェクト識別子(オブジェクトID)によって示され得る。オブジェクトIは、任意の好適なタイプ(英数字列、数字等)、及び任意の好適な長さまたはサイズ(32ビット、64ビット、128ビット等)であってもよい。

10

【0058】

データオブジェクト識別情報に基づいて小ボリュームを選択することの一例として、コーホート内のN個の格納ノードは、N/2対の格納ノードに分割してもよく、データオブジェクトの識別情報のハッシュ(または他の変換式)の所与のビットは、所与の対の格納ノードに対応してもよく、及びビットを用いて、を所与の対の内のいずれの格納ノードが、このデータオブジェクトについてのデータ要素を包含するかを示してもよい。簡潔な例証として、A~Pとして指定された16個の格納ノード/小ボリュームを有するコーホートでは、格納ノードは、以下のように対をなしてもよい。

20

対1	対2	対3	対4	対5	対6	対7	対8
A~B	C~D	E~F	G~H	I~J	K~L	M~N	O~P

【0059】

0ビットは、1つの対の中の第1のノードを指定することがあり、1ビットは、1つの対の中の第2のノードを指定することがある。例示のデータオブジェクトの識別情報のハッシュの8ビット部分は、たとえば、

10011100

であることがある。

【0060】

左端のビットが対1に対応することを前提として、このデータオブジェクトのデータ要素(レプリカまたはシャード)は、8個の対のノードB、C、E、H、J、L、M、及びOに格納される。他のデータオブジェクトの識別情報のハッシュが、ビットの比較的ランダムな組み合わせを生成するはずであるため、データ要素は、各対内のノードに、及びノード全体に比較的均等に分配される傾向がある。

30

【0061】

オブジェクトの識別情報のハッシュを用いて、ノードの対間にデータ要素を格納する上述の方式は、例示を目的として提供され、限定することを意図しないことに留意すべきである。データオブジェクト識別情報、またはその変換式は、データオブジェクトから生成されたデータ要素を格納するためのコーホート内のノード間で選択すると説明されたものとは異なる方法で用いられてもよい。たとえば、オブジェクトIDハッシュまたは他の変換式は、それぞれのデータオブジェクトから生成されたデータ要素が格納されることになっているコーホート内のM個の格納ノードの特定のサブセットを、確定的に示し得る。さらに、コーホート内の格納ノードを選択するためにデータオブジェクトの識別情報を用いることに加えて、少なくともいくつかの実施形態では、データオブジェクトの識別情報(たとえば、データオブジェクトの識別情報のハッシュまたは他の変換式)を、コーホートからデータオブジェクトを抽出すること、たとえば識別情報のハッシュによってコーホート内に予め格納されたデータオブジェクトのデータ要素(シャードまたはレプリカ)を検索することに用いてもよいことに留意すべきである。

40

【0062】

(メタデータによるデータ要素のタグ付け)

50

ボリュームコーホートを組み込んだオブジェクト冗長化格納システムの少なくともいくつかの実施形態では、コーホート内のN個の小ボリュームから選択されたM個の小ボリュームのセットに格納されたデータオブジェクトから生成されたM個のデータ要素のうち少なくとも1つが、格納サービスによって、データオブジェクトから生成され、コーホートに格納されたデータ要素とは他のものを検索することに用いられ得るメタデータにタグ付けされ得る。所与の小ボリューム上の所与のデータ要素についてのメタデータは、データ要素が格納されているコーホートのすべて位置（格納ノード/小ボリューム）を示し得る。格納サービスが、小ボリュームからデータ要素を抽出すると、メタデータ（存在する場合）もまた抽出され、そして必要であれば、データ要素の他のものを抽出するために用いられ得る。

10

【0063】

位置を示しているメタデータによるコーホート内のデータ要素のこのタグ付けは、レプリケーション手法を用いたオブジェクト冗長化格納システム内のレプリカであるデータ要素に対して行われ得るが、メタデータは、冗長な符号化手法、たとえばイレージャ符号化によってデータオブジェクトから生成されたシャードであるデータ要素のために特に有用であり得る。これは、データオブジェクトを再生するために、データオブジェクトについて生成されたM個のシャードのいくつかの最小数R（たとえば、M/2個のシャード）を要するためである。格納サービスがシャードを抽出するとき、存在する場合にはメタデータもまた抽出され、格納サービスによって、それぞれのデータオブジェクトを再生するために要するその他のシャードを検索するために用いられ得る。

20

【0064】

図4A~図4Cは、実施形態による、ボリュームコーホート内のメタデータによるデータ要素のタグ付けを示す。これらの図では、図3Bからコーホート302Bを例として用いている。図4Aでは、コーホート302Bからの小ボリューム320Bが示される。図3Bに示されるように、データ要素322A及び322Eは、小ボリューム320Bに格納される。データ要素322Aは、データ要素322Aが小ボリューム320A、320C、及び320D上でもまた検索されることを示すメタデータ324Aによってタグ付けされる。データ要素322Eは、データ要素322Aが小ボリューム320C、320E、及び320G上でもまた検索されることを示すメタデータ324Eによってタグ付けされる。いくつかの実施形態では、一貫性を保つために、小ボリューム320上の所与のデータ要素322についてのメタデータ324は、データ要素322がその小ボリューム320に格納されることをさらに示すことがあることに留意すべきである。たとえば、メタデータ324Aは、データ要素322Aが小ボリューム320A、320B、320C、及び320D上で検索されることを示すことがある。

30

【0065】

データ要素322に対するメタデータ324は、様々な実施形態において、数多くの形式または表現のいずれかで格納されもよい。しかしながら、少なくともいくつかの実施形態では、圧縮形式を用いて、格納及びデータ抽出オーバーヘッドを低減させ得る。単なる1つの非限定的な例として、メタデータ324のための簡易な表現は、ビットフィールドであり、この場合、各ビットは、コーホート内の格納ノード/小ボリュームのうち1つに対応する。たとえば、図3Bに例示のコーホート302Bに対するメタデータ324は、左端のビットが小ボリューム320Aに対応している等である8ビットフィールドであり得る。本例の表現を用いて、データ要素322Aに対するメタデータ324Aを、

40

1 1 1 1 0 0 0 0

と表現することができ、データ要素322Aがコーホート302B内の小ボリューム320A、320B、320C、及び320D上で検索されることを示している。データ要素322Eに対するメタデータ324Eは、

0 1 1 0 1 0 1 0

と表現することができ、データ要素322Eがコーホート302B内の小ボリューム320B、320C、320E、及び320Gで検索されることを示している。

50

【 0 0 6 6 】

いくつかの実施形態では、メタデータ 3 2 4 は、データオブジェクトについて生成された各データ要素 3 2 2 とともに格納される場合がある。図 4 B は、図 3 B のコーホート 3 0 2 B の第 1 の 4 つの小ボリューム 3 2 0 A ~ 3 2 0 D を示し、小ボリューム 3 2 0 A ~ 3 2 0 D の各々に格納されたデータ要素 3 2 2 A が、データ要素 3 2 2 A もまた格納されるコーホート 3 0 2 B 上の他の位置（たとえば、他の小ボリューム 3 2 0 ）を、データ要素 3 2 2 A ごとに示すそれぞれのメタデータ 3 2 4 A 1 ~ 3 2 4 A 4 によってタグ付けされることを示す。データ要素 3 2 2 A のうちいずれか 1 つが、格納サービスによって抽出されると、その対応するメタデータ 3 2 4 A もまた抽出され、コーホート 3 0 2 B に格納されたその他のデータ要素 3 2 2 A の 1 つ以上の検索及び抽出に用いられ得る。

10

【 0 0 6 7 】

メタデータ 3 2 4 がデータオブジェクトについての各データ要素 3 2 2 とともに格納される実施形態では、小ボリューム 3 2 0 の任意の所与の対に対して、対の中の各小ボリューム 3 2 0 は、データオブジェクトが、(a) 小ボリューム 3 2 0 に格納されたデータ要素 3 2 2 (シャードまたはレプリカ) を有し、及び (b) 対の中のその他の小ボリューム 3 2 0 上に格納されたデータ要素 3 2 2 (シャードまたはレプリカ) を有するはずであることが分かっている（または、小ボリューム 3 2 0 上のメタデータ 3 2 4 から判断することができる）ことに留意すべきである。たとえば、図 4 B では、小ボリューム 3 2 0 A がデータ要素 3 2 2 A を格納しているはずであることは、小ボリューム 3 2 0 B 上のメタデータ 3 2 4 A 2 から判断することができ、小ボリューム 3 2 0 B がデータ要素 3 2 2 A を有するはずであることは、小ボリューム 3 2 0 A 上のメタデータ 3 2 4 A 1 から判断することができる。この情報は、たとえば、コーホート内でデータ要素（シャードまたはレプリカ）が欠落しているため、修復されることを要するデータオブジェクトを識別するプロセスで用いられ得る。

20

【 0 0 6 8 】

冗長な符号化手法、たとえばイレージャ符号化によってデータオブジェクトから M 個のシャードが生成されて、コーホート内の N 個の小ボリュームのうち M 個に、シャードごとにメタデータとともに格納される、オブジェクト冗長化格納システムでのメタデータの使用の一例として、格納サービスは、データオブジェクトから M 個のシャードを生成し、データオブジェクトのシャードを格納するために十分な利用可能スペースを有するコーホート内の N 個の格納ノード / 小ボリュームの中から M 個の小ボリュームをランダムに選択し得る。各シャードが格納されると、コーホート内の M 個のシャードすべての位置が、メタデータとして、たとえば圧縮形式でシャードとともに格納される。コーホートからデータオブジェクトを抽出するために、格納サービスは、データオブジェクトを要求している（及び、要求されたデータオブジェクトに識別情報を提供している）コーホート内の N 個の小ボリュームのうち少なくとも R にアクセスすることができ、ここで R は、用いられている冗長な符号化方式によってデータオブジェクトを再生するために要するシャードの最小数である。場合によっては、R 個のノードのすべてがシャードを返すことがあり、この場合、データオブジェクトは、抽出されたシャードから再生されることができ、抽出が行われる。他のケースでは、R 個のノードはいずれもシャードを返さないかもしれず、この場合、R 個のノードの重複していないセットが、データオブジェクトを要求するために格納サービスによってアクセスされるかもしれない。しかしながら、ほとんどの場合、R 個のノードの 1 つ以上が、それらが指定されたデータオブジェクトについてのシャードを格納していないことを示し得る。一方で、R 個のノードの 1 つ以上の他のものが、オブジェクトについてのシャードが格納されたコーホート内の他の位置を示すメタデータとともに、シャードを返すことがある。ノードのいずれか（しかし全てではない）が、メタデータとともにシャードを返さないと、その後格納サービスは、コーホート内のさらなるノードにアクセスして、さらなるシャードを得ることを要することを知らず、返されたメタデータを用いて、アクセスされるコーホート内のノードを知的に選択することができる。

30

40

50

【0069】

いくつかの実施形態では、メタデータ324をデータオブジェクトについて生成された各データ要素322とともに格納することに代えて、データオブジェクトについて生成されたデータ要素322のサブセットのみが、メタデータによってタグ付けされる場合がある。図4Cは、図3Bのコホート302Bの第1の4つの小ボリューム320A~320Dを示し、データ要素322Aが小ボリューム320A~320Dの各々に格納されることを示す。しかしながら、小ボリューム320Aに格納されたデータ要素322Aのみが、データ要素322Aもまた格納されるコホート302B上の他の位置（たとえば、他の小ボリューム320）を示すメタデータ324A1によってタグ付けされる。データ要素322Aが格納サービスによって小ボリューム320Aから抽出されると、その対応するメタデータ324A1もまた抽出され、コホート302B上に格納されたその他のデータ要素322Aつ以上の検索及び抽出に用いられ得る。

10

【0070】

メタデータがデータオブジェクトから生成されてコホートに格納されるM個のシャードの一部のみとともに格納されるオブジェクト冗長化格納システムの一例として、格納サービスは、データオブジェクトからM個のシャードを生成して、データオブジェクトのシャードを格納するために十分な利用可能なスペースを有するコホート内のN個の格納ノード/小ボリュームの中からM個の小ボリュームをランダムに選択し得る。コホート内のM個のシャードすべての位置が、メタデータとしてシャードのサブセットのみとともに格納される。一例として、40個の小ボリューム(N=40)を有するコホートでは、M=20、及びR=10である場合、メタデータは、シャードのうち5つとものみ格納され得る。例示のコホートからデータオブジェクトを抽出するために、格納サービスは、データオブジェクトを要求しているコホート内のN個の小ボリュームのうち20個にアクセスし得る。シャードと、他のシャードのすべての位置を示すメタデータを含む5個の小ボリュームがあるため、これらの5個の小ボリュームのうち少なくとも1つが、アクセスされた20個の小ボリュームの中にある可能性は、およそ98%である。(別の例として、20に代えてサンプルサイズ10では、これらの5個の小ボリュームのうち少なくとも1つが、アクセスされた10個の小ボリュームの中にある可能性は、およそ78%である)。このように、メタデータは、小ボリュームのサブセットのみとともに格納されることができ、一方で十分なサンプルサイズを所与として、メタデータが最初のアクセスで得られる高い可能性を依然として提供する。いったん格納サービスがメタデータを得ると、当該サービスは、メタデータを用いて、アクセスされるコホート内のさらなるノードを選択して、必要であればさらなるシャードを得ることができる。

20

30

【0071】

いくつかの実施形態は、コホートにデータ要素を格納し、そこからデータ要素を抽出するためのハイブリッドな方法を用いてもよく、これはコホート内の他のデータ要素を検索するために、メタデータとともにコホートに格納されたデータ要素をタグ付けするための方法と、データオブジェクト(またはその変換式)の識別情報を用いて、コホート内のデータ要素のための格納位置を選択する方法との組み合わせである。たとえば、データオブジェクトの識別情報のハッシュを用いて、データオブジェクトから生成された1つ以上のデータ要素のための小ボリュームを選択し得ると同時に、その他のデータ要素のための小ボリュームをランダムに選定し得る。データ要素のすべてを検索するためのメタデータは、識別情報のハッシュによって位置が決定されるデータ要素とともに格納され得る。このように、データオブジェクトを抽出すると、識別情報のハッシュは、格納サービスを、そのメタデータが格納される1つ以上のデータ要素の位置に向かわせることができ、抽出されたメタデータは、必要であれば、データオブジェクトについてのさらなるデータ要素を知的に抽出する際に用いられ得る。

40

【0072】

図5は、少なくともいくつかの実施形態による、データオブジェクトを作成して、オブジェクト冗長化格納システム内のコホートに格納するため方法の高レベルなフローチャ

50

ートである。本方法は、たとえば1つ以上のデバイス上に組み込まれた格納サービスによって、またはそれを介して行われ得る。格納サービスの実施形態が組み込まれ得るシステムの例を、図13に示す。

【0073】

500に示されるように、N個の小ボリュームを含むコーホートを作成し得る。少なくともいくつかの実施形態では、N個の小ボリュームの各々が、異なる格納ノードまたはデバイス上に存在し得る。いくつかの実施形態では、コーホートは、コーホートにデータオブジェクトの基本または初期セットを格納することによって初期化され得る。502に示されるように、コーホートに格納されることになっているデータオブジェクトは、たとえば格納サービスの1つ以上のクライアントのうち1つから受信され得る。

10

【0074】

504に示されるように、格納サービスは、データオブジェクトが格納されることになっているN個の小ボリュームのうちM個を選択することができ、ここで、MはN未満である。実施形態では、異なる手法を用いて、M個の小ボリュームを選択してもよい。たとえば、いくつかの実施形態では、ランダムなまたは擬似ランダムな選択手法を用いてもよい。いくつかの実施形態では、M個の小ボリュームを選択するときに、1つ以上の要素、たとえば利用可能な格納スペースを考慮してもよい。いくつかの実施形態では、データオブジェクト（またはその変換式、たとえばハッシュ）の識別情報を、M個の小ボリュームのうち少なくともいくつかを決定する際に用いてもよい。いくつかの実施形態においては、これらの手法の組み合わせまたは変形を用い得る。

20

【0075】

506に示されるように、M個のデータ要素は、データオブジェクトについて、またはそこから生成され得る。データ要素は、たとえばレプリケーション手法によって生成されたデータオブジェクトのレプリカであってもよい。あるいは、データ要素は、冗長な符号化方式、たとえばイレージャ符号化によって生成されたデータオブジェクトのシャードであってもよい。

【0076】

508に示されるように、M個のデータ要素のうち1つが、コーホート内のM個の選択された小ボリュームの各々に格納され得る。少なくともいくつかの実施形態では、M個のデータ要素のうち少なくとも1つの他のものの格納位置（たとえば、小ボリューム）を示すメタデータが、M個のデータ要素のうち少なくとも1つとともに格納され得る。いくつかの実施形態では、メタデータは、コーホート内のM個のデータ要素の各々とともに格納される。

30

【0077】

図6は、少なくともいくつかの実施形態による、データオブジェクトがレプリケーション手法によって格納される、コーホートからデータオブジェクトを抽出するための方法の高レベルのフローチャートである。本方法は、たとえば1つ以上のデバイス上に組み込まれた格納サービスによって、またはそれを介して行われ得る。格納サービスの実施形態が組み込まれ得るシステムの例を、図13に示す。レプリケーション手法では、データオブジェクトを抽出するために、1つのレプリカのみがコーホートから抽出されることを要することに留意すべきである。

40

【0078】

600に示されるように、たとえば格納サービスの1つ以上のクライアントのうち1つから、コーホートに格納されたデータオブジェクトについての要求を受信し得る。データオブジェクトは、データオブジェクトのM個のレプリカを生成し、コーホート内のN個の小ボリュームから選択されたM個の小ボリュームの各々に、レプリカのうち1つを格納するレプリケーション手法によって、予めコーホートに格納されていてもよい。

【0079】

602に示されるように、コーホート内のN個の小ボリュームのうちP個が選択され得る。異なる手法を用いて、P個の小ボリュームを選択してもよい。たとえば、いくつかの

50

実施形態では、ランダムなまたは擬似ランダムな選択手法を用いてもよい。いくつかの実施形態では、データオブジェクトの識別情報（またはその変換式、たとえばハッシュ）を、P個の小ボリュームのうち少なくともいくつかを決定する際に用いてもよい。一般的に、PはMと等しいかまたはそれより小さいことがあることに留意すべきである。しかしながら、Pは、1から最大Nまでであり、Nを含む任意の数であることができる。

【0080】

604に示されるように、データオブジェクトのレプリカが、選択されたP個の小ボリュームの各々から要求され得る。606において、P個の小ボリュームのうち少なくとも1つからレプリカが返されると、その後608に示されるように、要求者にデータオブジェクトを提供し得る。そうでなければ、604において、本方法は、602に戻って、P個の（または他の何らかの数の）小ボリュームの重複していないセットを選択し、小ボリュームの新しいセットからのデータオブジェクトを要求することができる。

10

【0081】

図7は、少なくともいくつかの実施形態による、データオブジェクトが冗長な符号化手法によって格納されるコーホートからデータオブジェクトを抽出するための方法の高レベルのフローチャートである。本方法は、たとえば1つ以上のデバイス上に組み込まれた格納サービスによって、またはそれを介して行われ得る。格納サービスの実施形態が組み込まれ得るシステムの例を、図13に示す。冗長な符号化手法では、データオブジェクトを再生するために、データオブジェクトから作成されたシャードのほぼ最小数が必要とされることに留意すべきである。本明細書では、Rは、データオブジェクトを再生するために要するシャードの最小数を表現するために用いる。一方、Mは、データオブジェクトから作成され、コーホート内のN個の小ボリュームから選択されたM個の小ボリュームのサブセットに格納されたシャードの総数を表現する。必ずしもそうではないが、一般的には、 $R = M / 2$ であることに留意すべきである。たとえば、非限定的な例のイレージャ符号化方式では、 $M = 20$ 、及び $R = 10$ である。

20

【0082】

700に示されるように、たとえば格納サービスの1つ以上のクライアントのうち1つから、コーホートに格納されたデータオブジェクトについての要求を受信し得る。データオブジェクトは、データオブジェクトのM個のシャードを生成し、コーホート内のN個の小ボリュームから選択されたM個の小ボリュームの各々に、シャードのうち1つを格納する冗長な符号化手法によって、予めコーホートに格納されていてもよい。

30

【0083】

702に示されるように、コーホート内のN個の小ボリュームのうちR個が選択され得る。異なる手法を用いて、R個の小ボリュームを選択してもよい。たとえば、いくつかの実施形態では、ランダムなまたは擬似ランダムな選択手法を用いてもよい。いくつかの実施形態では、データオブジェクトの識別情報（またはその変換式、たとえばハッシュ）を、R個の小ボリュームのうち少なくともいくつかを決定する際に用いてもよい。本例の実施では、Rは、用いられている冗長な符号化方式によってデータオブジェクトを再生するために要するシャードの最小数であり、R個の小ボリュームは、シャードについて問い合わせるために、少なくとも初期に選択されることに留意すべきである。しかしながら、他の実施においては、問い合わせのためにR個よりも多いかまたはそれよりも少ない小ボリュームが選択されてもよい。

40

【0084】

704に示されるように、データオブジェクトのシャードが、選択されたR個の小ボリュームの各々から要求され得る。要求されたR個の小ボリュームの各々は、データオブジェクトについてのシャードを格納しているか、またはしていない。要求されたR個の小ボリュームのうち1つが、データオブジェクトについてのシャードを有している場合、小ボリュームは、シャードを返す。小ボリュームが、コーホート内の他のシャードを検索するためのメタデータもまた格納している場合、メタデータも返されることがある。

【0085】

50

706において、選択されたR個の小ボリュームからシャードが返されない場合、その後本方法は、702に戻り、R個（または他の何らかの数）の小ボリュームの重複していないセットを選択して、データオブジェクトについての小ボリュームの新しいセットについて問い合わせることができる。

【0086】

706において、選択されたR個の小ボリュームから少なくとも1つのシャードが戻されると、その後格納サービスは、データオブジェクトを再生するために十分なシャード（すなわち、少なくともR個のシャード、ここでRは、用いられている冗長な符号化方式によってデータオブジェクトを再生するために要するシャードの最小数である）が得られているかを判断し得る。708において、十分なシャード（すなわち、少なくともR個のシャード）が得られている場合、本方法はその後、712に進む。

10

【0087】

708において、少なくとも1つであるが十分ではない（すなわち、R個よりも少ない）シャードが得られている場合、少なくともいくつかの実施形態では、格納サービスはその後、コーホート内の他の小ボリュームから1つまたはさらなるシャードを得ることがある。すなわち、抽出されたシャードの少なくとも1つとともに格納され、R個の要求された小ボリュームから得られたシャードとともに格納サービスに戻されたメタデータに従って、その他の小ボリュームを格納サービスによって検索し得る。シャードの所与の1つとともに格納されたこのメタデータは、少なくとも1つの他のシャードのコーホート内の位置を示し、かついくつかの実施形態では、シャードごとの位置を示すことに留意すべきである。

20

【0088】

あるいは、710に示されるように、シャードとともに抽出されたメタデータを用いてさらなるシャードを検索することに代えて、いくつかの実施形態では、本方法は、代わりに702に戻り、R個の（または他の何らかの数）小ボリュームの重複していないセットを選択して、さらなるシャードを抽出するために小ボリュームの新しいセットについて問い合わせることができる。

【0089】

712において、いったん少なくともR個のシャードが、コーホートから抽出されると、その後、用いられている冗長な符号化方式、たとえばイレージャ符号化方式によって、抽出されたシャードからデータオブジェクトを再生し得る。そして、再生されたデータオブジェクトは、要求者、たとえば格納サービスの1つ以上のクライアントのうち1つに提供され得る。

30

【0090】

（ボリュームコーホート内の小ボリュームの照合調整）

オブジェクト冗長化格納システム内のボリュームコーホートの小ボリュームを照合調整するための方法及び機器の様々な実施形態が説明される。図2A～図7を参照して説明されたようなコーホートを用いるオブジェクト冗長化格納システムでは、図1に示されたような格納ノードコーホートを用いるオブジェクト冗長化格納システムのように、ボリューム内部の小ボリュームのコンテンツを、たとえばデータオブジェクトについてのデータ要素（レプリカまたはシャード）が、データ要素が格納されることが想定される小ボリュームのセット内の小ボリュームのすべてに実際に格納されることを確実にする照合調整プロセスの一部として、周期的に比較することが必要であるかまたは望ましいかもしれない。図1の格納システムにおいては、このセットは、ボリューム102内のM個の小ボリューム120全てを含むが、これは、各データオブジェクトのレプリカまたはシャードが、コーホートの中の各小ボリュームに格納されているはずであり、群の中の各小ボリュームが、同一のコンテンツを含んでいるはずであるためである。しかしながら、図2A～図7を参照して説明されたようなボリュームコーホートを用いる格納システムでは、データオブジェクトごとに、セットは、この特定のデータオブジェクトについてのデータ要素（レプリカまたはシャード）が格納されることになっているコーホート202内のN個の小ボリ

40

50

ューム 2 2 0 の中から選択された M 個の小ボリューム 2 2 0 を含む。故に、コーホート内の任意の 2 つの所与の小ボリュームのコンテンツは、通常は同一ではないことに留意すべきである。

【 0 0 9 1 】

いずれのタイプのオブジェクト冗長化格納システムにおいても、この小ボリュームのコンテンツの比較は、中央部またはシステムに、各小ボリュームのコンテンツの完全な一覧またはリストをダウンロードして、中央システムに比較をさせることによって行うことができる。しかしながら、全ての小ボリュームから中央部に完全かつ詳細なリストをダウンロードすることは、著しい量のネットワーク帯域幅を用いることとなる。

【 0 0 9 2 】

オブジェクト冗長化格納システム内の小ボリュームを比較して照合調整するとき、様々な手法を採用して、ネットワーク帯域幅を低減させ得る。たとえば、図 1 に示されるように、ボリューム内の小ボリュームごとのコンテンツが同一であるべきであるオブジェクト冗長化格納システムでは、ハッシュ手法を用いてもよい。たとえばハッシュ木手法を用いてもよく、これは、群の中の各小ボリュームが、その小ボリューム内に格納されたデータオブジェクトの識別子（オブジェクト ID と称される）のリストを取得し、オブジェクト ID のリストを、たとえばオブジェクト ID のハッシュまたは他の変換式に基づいてサブリストに分け、（分類された）サブリストごとにハッシュを生成し、そしてこれらのハッシュのハッシュ値を計算してルートハッシュを生成する。各小ボリュームから結果として得られたルートハッシュは、中央部またはシステムに送られる。中央システムでは、すべての小ボリュームが同じルートハッシュを報告した場合、さらなる照合調整の必要はない。2 つ以上の小ボリュームについてのルートハッシュが異なる場合、その後サブリストハッシュを比較して、異なっているサブリストハッシュを識別し得る。異なるサブリストハッシュを用いて、ハッシュ木のレベルに応じて異なる小ボリュームからのデータオブジェクトのサブリストを識別し得る。いったんデータオブジェクトの異なるサブリストが識別されると、異なるサブリストについてのみのオブジェクト ID のリストが、中央システムにダウンロードされて比較され、ボリューム内での照合調整（たとえば、レプリケーションまたはイレージャ符号化再構築）を必要とするデータオブジェクトのセットを識別し得る。

【 0 0 9 3 】

しかしながら、図 2 A ~ 図 7 を参照して説明されたような、ボリュームコーホートを用いるオブジェクト冗長化格納システムでは、各小ボリュームは、一般的には、唯一のデータオブジェクトのセットを含むため、図 1 を参照して上述されたようなハッシュ木手法は機能しない。これは、コーホート内の任意の 2 つの所与の小ボリュームのルートハッシュが、通常は異なるためである。

【 0 0 9 4 】

オブジェクト冗長化格納システム内のボリュームコーホートの小ボリュームを比較して照合調整するための方法の実施形態が説明される。この実施形態は、比較的小さなハッシュを用いて、格納システム内で比較及び照合調整を行い、比較及び照合調整プロセスの間、ネットワーク帯域幅をセーブすることを可能にする。いくつかの実施形態では、これらの方法は、たとえば図 4 A 及び図 4 B に示されたような、小ボリュームの各々にデータ要素とともに格納されるメタデータ 3 2 4 を活用し得る。図 3 A、図 3 B、図 4 A、及び図 4 B を参照すると、所与のデータ要素 3 2 2 についてのメタデータ 3 2 4 は、それぞれのデータオブジェクトについてのデータ要素が格納される（またはされるはずである）コーホート 3 0 2 B 内の小ボリューム 3 2 0 のすべてを示す。このメタデータ 3 2 4 を介して、図 4 B に示されたように、全てのデータ要素 3 2 2 とともに格納されると、所与のデータオブジェクトが属する（及び、そのために、オブジェクトについてのデータ要素が格納されるはずである）M 個の小ボリューム 3 2 0（コーホート内の N 個の小ボリュームのうち）のフルセットが、その特定のデータオブジェクトを所有する全ての小ボリューム 3 2 0 に知られるか、またはそれらによって決定されることができる。

10

20

30

40

50

【 0 0 9 5 】

いくつかの実施形態では、データ要素とともに格納されたメタデータを活用して、他の小ボリュームを有する共通のオブジェクトリストを決定することに代えて、小ボリュームは、データオブジェクトの識別情報（たとえば、オブジェクトID）を用いて、共通のオブジェクトリストを決定し得る。たとえば、所与のデータ要素に関連するオブジェクトIDのハッシュまたは他の変換を用いて、それぞれのデータオブジェクトから生成されたデータ要素が格納されたコーホート内のM個の格納ノードのサブセットを確定的に示し得る。このように、格納ノード/小ボリュームは、この変換を、そのデータ要素を格納するデータオブジェクトのオブジェクトIDに適用して、データオブジェクトもまた格納されているはずであるその他の格納ノード/小ボリュームを判断し得る。

10

【 0 0 9 6 】

図8は、少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームの比較を行うための方法をグラフィカルに図示する。本方法は、たとえば図2Aに示されたようなオブジェクト冗長化格納システム200において実施され得る。図8を参照すると、本方法は、「ゴシップ」法として見なされることがあり、これは、コーホート内の各小ボリューム820が、コーホート内の他の小ボリューム820と互いに通信して、小ボリューム820の各々の対が共通して有しているはずであるコンテンツ（共通のオブジェクトリスト826）を比較し、その後、照合調整モジュール890によって、図示されたような中央部に、任意の検出された差分830を通信する。照合調整モジュール890は、小ボリューム820から差分830を収集し、収集された情報を用いて、オブジェクト冗長のためにレプリケーション手法を用いる格納システムにおいて、必要に応じてデータ要素のレプリケーションを行うか、またはイレージャ符号化等の冗長な符号化手法を用いる格納システムにおいて、必要に応じて冗長なコード（シャード）再構築を行う。少なくともいくつかの実施形態では、照合調整モジュール890は、照合調整プロセスの少なくとも一部を実施してもよく、コンポーネント、モジュール、または図2Aに示されたような格納システム200及び/または格納サービス250の一部であってもよいことに留意すべきである。

20

【 0 0 9 7 】

小ボリューム比較方法の実施形態では、コーホート内の各小ボリューム820は、互いにコーホート内の他の小ボリューム820と互いに周期的または非周期的に通信して、小ボリュームコンテンツを比較する。図8は、2つの例示の小ボリューム820A及び820Bについての本プロセスを示す。小ボリューム820Aは、関連付けられたメタデータ824Aを有するデータ要素822A（シャードまたはレプリカ）のセットを含む。小ボリューム820Bは、関連づけられたメタデータ824Bを有するデータ要素822B（シャードまたはレプリカ）のセットを含む。コーホート格納システムでは、一般的には、任意の2つの小ボリューム820内のデータ要素822のセットは、ある程度重複するが、同一ではない。言い換えると、2つの小ボリューム820は、同じデータオブジェクトのセットから生成された何らかのデータ要素を各々包含するが、小ボリューム820は両方とも、その他の小ボリュームと共有していない他のデータオブジェクトについての他のデータ要素もまた包含している。いくつかの実施形態では、各小ボリューム820上のメタデータ824が、それぞれのデータオブジェクトについてのデータ要素822が格納されているはずであるM個の小ボリューム820の完全なセットをデータ要素822ごとに示す。所与の小ボリューム820は、それがメタデータ824に従って、コーホート内の任意の他の小ボリューム820と共通して有している（または有しているはずである）データオブジェクトのリスト（共通のオブジェクトリスト826）を包含し得るか、または生成し得る。あるいは、各小ボリューム820は、それが格納しているデータ要素822のオブジェクトIDに変換式（たとえば、ハッシュ）を適用して各データ要素822もまた格納されているはずであるその他の小ボリューム820を判定してもよく、またこの情報を用いて、共通のオブジェクトリスト826を生成してもよい。

30

40

【 0 0 9 8 】

50

図8に示されるように、小ボリューム820Aは、それが小ボリューム820Bと共通して有しているかまたは有しているはずであるデータオブジェクトのすべてをリストアップする共通のオブジェクトリスト826Aを包含するかまたは生成する。同様に、小ボリューム820Bは、それが小ボリューム820Aと共通して有しているかまたは有しているはずであるデータオブジェクトのすべてをリストアップする共通のオブジェクトリスト826Bを包含するかまたは生成する。少なくともいくつかの実施形態では、各データオブジェクトは、格納システム内のデータオブジェクトを一意的に識別し得るオブジェクト識別子(オブジェクトID)によって示され得る。オブジェクトIDは、任意の好適なタイプ(英数字列、数字等)、及び任意の好適な長さまたはサイズ(32ビット、64ビット、128ビット等)であってもよい。少なくともいくつかの実施形態では、リスト826A及び826B内のオブジェクトIDは、同じ配列方式に従って配列されて、2つのリスト826が全く同じオブジェクトIDのセットを包含する場合、2つのリスト826が同一であるようにされ得る。しかしながら、リスト826Aは、少なくとも初期に順不同とされてもよく、オブジェクトIDの配列は、以下に説明されるようなハッシュ手法の間に行われる。

10

【0099】

そして、各小ボリューム820は、その共通のオブジェクトリスト826それぞれにハッシュ手法850を適用して、リスト826についてのハッシュ値を生成する。小ボリューム820の両方が、基本的に同じハッシュ手法850を用いて、2つのリスト826が同一である場合にハッシュ値が同一であるようにすることに留意すべきである。

20

【0100】

少なくともいくつかの実施形態では、ハッシュ手法850は、ハッシュ木手法であってもよい。ハッシュ木手法の少なくともいくつかの実施形態では、共通のオブジェクトリスト826は、たとえばオブジェクトIDのハッシュまたは他の変換式に基づいて、2つ以上の格納されたサブリストに分けられる。ハッシュは、各サブリストから生成される。いくつかの実施形態では、ハッシュは、サブリスト内のオブジェクトIDごとに生成される。あるいは、2つ以上のオブジェクトIDが組み合わせられてもよく、ハッシュは、サブリスト内の2つ以上のオブジェクトIDのそのような組み合わせごとに生成されてもよい。サブリストごとに、ハッシュを組み合わせるとハッシュ値を計算し、サブリストハッシュを生成してもよい。そして、サブリストハッシュを組み合わせると(たとえば、連結させて)ハッシュ値を計算し、共通のオブジェクトリスト826についてのルートハッシュを生成してもよい。いくつかの実施形態では、オブジェクトIDに代えてまたはそれに加えて、ハッシュ手法において、オブジェクトID以外のデータオブジェクトに関する情報を用いてもよい。

30

【0101】

図12は、オブジェクト識別子(オブジェクトID)1200のセット(たとえば、オブジェクトリスト)から生成された非限定的な例のハッシュ木を示し、これは、オブジェクトID1200から生成されたハッシュ1202、オブジェクトIDハッシュ1202のセットのハッシュとして生成されたサブリストハッシュ1204、及びサブリストハッシュ1204のハッシュとして生成されたルートハッシュ1206を示す。図12は、3つのレベルを有するハッシュ木の例を示すが、いくつかの実施形態では、ハッシュ木にさらなるレベルがあってもよいことに留意すべきである。いくつかの実施形態では、ハッシュ1202は、オブジェクトID1200ごとに生成され得る。あるいは、2つ以上のオブジェクトIDを組み合わせるともよく、ハッシュ1202を、2つ以上のオブジェクトIDのそのような組み合わせごとに生成してもよい。たとえば、ハッシュ1202とサブリストハッシュ1204との間にさらなるハッシュ木レベルがあり、2つ以上のハッシュ1202が組み合わせられてハッシュ値が計算され、中間レベルのハッシュを生成してもよく、そして、中間レベルのハッシュを組み合わせるとハッシュ値を計算し、サブリストハッシュ1204を生成してもよい。

40

【0102】

50

図 8 を再度参照すると、共通のオブジェクトリスト 8 2 6 についてルートハッシュが生成された後、2 つの小ボリューム 8 2 0 は、生成されたハッシュ値（共通のオブジェクトリスト 8 2 6 についてのルートハッシュ）を交換する。各小ボリューム 8 2 0 上では、ハッシュ比較 8 6 0 機能またはモジュールが、2 つの共通のオブジェクトリスト 8 2 6 についての 2 つのハッシュ値を比較する。

【 0 1 0 3 】

2 つのハッシュ値が同じである場合、そのとき 2 つの共通のオブジェクトリスト 8 2 6 は同じであり、2 つの小ボリューム 8 2 0 は、的確に重複しているデータオブジェクトのセットを有すると考えられる。その場合、2 つの小ボリューム 8 2 0 はその後、比較を終える。しかしながら、2 つの小ボリューム 8 2 0 が、周期的にまたは非周期的に比較を繰り返す得ることに留意すべきである。いくつかの実施形態では、小ボリューム 8 2 0 は、中央システムまたは中央部、たとえば照合調整モジュール 8 9 0 に通知して、モジュール 8 9 0 に、調査が所与の他の小ボリューム 8 2 0 で進行し、2 つの小ボリューム 8 2 0 が、的確に重複しているデータオブジェクトのセットを有することを知らせ得ることに留意すべきである。

10

【 0 1 0 4 】

2 つのハッシュ値が同じではない場合、そのとき 2 つのオブジェクトリスト 8 2 6 間には何らかの差分がある。そして、各小ボリューム 8 2 0 は、共通のオブジェクトリスト 8 2 6 に対して差分判定 8 7 0 を行い、両方の小ボリューム 8 2 0 上にあるはずであるがそこにはないデータオブジェクトについての 1 つ以上のオブジェクト ID を判定し得る。少なくともいくつかの実施形態では、ハッシュ手法 8 5 0 によって 2 つの小ボリューム上に生成されたハッシュ木を用いて、具体的な差分を判定してもよい。たとえば、いくつかの実施形態では、2 つのハッシュ木のサブリストハッシュを比較して、オブジェクト ID の異なる特定のサブリストを判定してもよく、そのように識別されたサブリストを比較して、差分を判定することができる。

20

【 0 1 0 5 】

少なくともいくつかの実施形態では、両方の小ボリューム 8 2 0 A 及び 8 2 0 B 上の差分判定 8 7 0 の結果は、それぞれ差分リスト 8 3 0 A 及び 8 3 0 B として、中央システムまたは中央部、たとえば照合調整モジュール 8 9 0 に送られることがある。少なくともいくつかの実施形態では、各差分リスト 8 3 0 は、リスト 8 3 0 が生成された 2 つの小ボリューム 8 2 0 を示してもよく、識別された小ボリューム 8 2 0 の両方にあるはずであるがそこにはないデータオブジェクトについての 1 つ以上のオブジェクト ID をリストアップし得る。

30

【 0 1 0 6 】

少なくともいくつかの実施形態では、コーホート内の小ボリューム 8 2 0 の各々が、周期的にまたは非周期的に上記の比較方法を、コーホート内の他の小ボリューム 8 2 0 各々に実施する。このように、いくつかの実施形態では、各小ボリューム 8 2 0 は、照合調整モジュール 8 9 0 に、特定の他の小ボリュームとの比較専用の 1 つ、2 つ、またはそれ以上の差分リスト 8 3 0 を送り得る。あるいは、小ボリューム 8 2 0 は、2 つ以上の他の小ボリューム 8 2 0 との比較を行い、検出された差分を収集して、照合調整モジュール 8 9 0 に、当該小ボリューム 8 2 0 と 2 つ以上の他の小ボリューム 8 2 0 との間の差分を示す組み合わせられた差分リスト 8 3 0 を周期的にまたは非周期的に送ってもよい。

40

【 0 1 0 7 】

中央システムまたは位置（たとえば、照合調整モジュール 8 9 0 ）は、コーホート内の N 個の小ボリューム 8 2 0 のいくつかまたはすべてから、差分リスト 8 3 0 を収集する。周期的にまたは非周期的に、または必要または所望のように、照合調整モジュール 8 9 0 は、コーホート内の小ボリューム 8 2 0 の 1 つ、2 つ、またはそれ以上で照合調整を行い、小ボリューム上の実際のデータオブジェクトを、小ボリューム 8 2 0 間で行われた比較によって判定されたように、小ボリューム 8 2 0 上にあるはずであるデータオブジェクトのセットと照合調整し得る。照合調整は、レプリケーション手法を用いる格納システム内

50

の1つ以上の小ボリューム820に対するデータオブジェクトのレプリケーションを伴う場合があり、または冗長な符号化手法、たとえばイレージャ符号化を用いる格納システムのデータオブジェクトについてのシャードの再生を伴う場合がある。

【0108】

中央システムまたは中央部、たとえば照合調整モジュール890は、冗長な符号化手法を用いて、コーホート内に格納されたデータオブジェクトについてのシャードを回復させるかまたは再構築する格納システムにおいて、本質的に必要であるかもしれないことに留意すべきである。中央照合調整モジュール890は、レプリケーション手法を用いる格納システムで、そのようなシステムの代替として用いられ得るが、2つの小ボリューム820は、図8に示されたような比較を行い、一方の小ボリューム820上で欠落しているレプリカを判定し得、その他の小ボリューム820は、第1の小ボリューム820に、欠落しているレプリカを直接提供し得る。

10

【0109】

いくつかの実施形態では、図8のハッシュ手法850のようなハッシュ木手法を用いることに代えて、共通のオブジェクトリスト826のかなり簡易なハッシュまたは他の変換式を生成して、小ボリューム820によって交換させてもよい。これは、ハッシュを生成するプロセスを簡易にするかもしれないが、ハッシュ木のレベルは、生成されたハッシュ値が異なる場合に差分判定870中に使用できない。そのため、2つの小ボリューム820間の差分を判定するために、さらなる情報を交換することを要する場合がある。

【0110】

いくつかの実施形態では、図8に示されたようなオブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームを比較するための方法の代替または変形として、各小ボリューム820は、図8に示されるように、共通のオブジェクトリスト826を生成し、リスト826に対するハッシュ値(たとえば、ルートハッシュ)を生成し、ハッシュ値(たとえば、ルートハッシュ)を他の小ボリューム820と交換して、ハッシュ値(たとえば、ルートハッシュ)を比較し得る。しかしながら、小ボリューム820で差分判定870を行い、図8に示されたような照合調整モジュール890に、もし差分があればその判定された差分を報告することに代えて、各小ボリューム820が代わりに、共通のオブジェクトリストについてのハッシュ値(たとえば、ルートハッシュ)が、別の小ボリューム820から受信されたそれぞれのハッシュ値と一致しないことを、照合調整モジュール890に報告する。そして、照合調整モジュール890は、必要に応じて差分判定を行い、照合調整を要するかもしれない2つの小ボリューム820間の特定の差分を判定する。これは、照合調整モジュール890が小ボリューム820からのさらなる情報を要求することを必要とするかもしれないことに留意すべきである

20

30

【0111】

図10は、少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームの比較を行うための方法のフローチャートである。本方法は、たとえば図2Aに示されたようなオブジェクト冗長化格納システム200において示されたような格納ノード及び/または小ボリューム上で行われ得る。本方法は、「ゴシップ」法として見なされることがあり、これは、コーホート内の各小ボリュームが、コーホート内の他の小ボリュームと互いに通信して、小ボリュームの各対が共通して有しているはずであるコンテンツを比較し、その後中央部、たとえば照合調整モジュールに、いずれかの検出された差分を通信する。中央部は、小ボリュームから差分を収集し、収集された情報を用いて、必要であればオブジェクト冗長のためにレプリケーション手法を用いる格納システム内でデータ要素のレプリケーションを行うか、または必要であれば、冗長な符号化手法、たとえばイレージャ符号化を用いる格納システムでシャード再構築を行う。

40

【0112】

図10の1000に示されるように、小ボリュームは、コーホート内の1つ以上の他の小ボリュームごとに共通のオブジェクトリストを判定するかまたは生成し得る。その他の

50

小ボリュームの所与の1つについての共通のオブジェクトリストは、この小ボリュームがその他の小ボリュームと共通して有しているかまたは有しているはずであるデータオブジェクトを示す。少なくともいくつかの実施形態では、各データオブジェクトは、オブジェクト識別子（オブジェクトID）によるリストで示されてもよい。

【0113】

1002に示されるように、小ボリュームは、共通のオブジェクトリストのハッシュを生成してもよい。少なくともいくつかの実施形態では、図8を参照して説明されたハッシュ木手法を用いて、共通のオブジェクトリストについてのハッシュ値（すなわち、ルートハッシュ）を生成してもよい。しかしながら、いくつかの実施形態では、他のハッシュ手法を用いてもよい。

10

【0114】

1004に示されるように、小ボリュームは、共通のオブジェクトリストを、1つ以上の他の小ボリュームの各々と交換し得る。そして要素1006～1014が、1004でハッシュが交換された他の小ボリュームごとに行われる。

【0115】

1006において、1つ以上の他の小ボリュームのうち特定の1つの共通のオブジェクトリストについて生成されたハッシュ値は、特定の他の小ボリュームから受信されたハッシュ値と比較され得る。1008において、2つのハッシュが一致する場合、2つの小ボリュームは整合性があり、本方法は1014に飛ぶ。1008において、2つのハッシュが一致しない場合、2つの小ボリューム上の共通のオブジェクトリスト間の差分が、1010において判定され得る。少なくともいくつかの実施形態では、差分を判定することは、両方の小ボリュームにあるはずであるがそこにはないデータオブジェクトについての1つ以上のオブジェクトIDを判定することを伴う。少なくともいくつかの実施形態では、ハッシュ木手法によって2つの小ボリューム上に生成されたハッシュ木を用いて、小ボリューム間の特定の差分を判定し得る。たとえば、いくつかの実施形態では、2つのハッシュ木の1つ以上の異なるレベルにおけるハッシュ値を比較して、異なるオブジェクトIDの特定のサブリストを判定してもよく、そして、識別されたサブリストを比較して、特定の差分を判定することができる。

20

【0116】

1012に示されるように、1010で判定された任意の差分は、中央部、たとえば照合調整モジュールに報告され得る。少なくともいくつかの実施形態では、差分リストは、リストが生成された2つの小ボリュームを示すことがある中央部に送られる。中央部は、識別された小ボリュームの両方にあるはずであるがそこにはないデータオブジェクトについての1つ以上のオブジェクトIDをリストアップすることがある。

30

1014において、比較されるべきさらなるハッシュがある場合、本方法はその後、要素1006に戻る可能性がある。そうでなければ、次回当該方法が起動されるまで、比較方法は、この小ボリュームで行われる。

【0117】

図9は、少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームを比較するための代替方法をグラフィカルに示す。この代替方法は、たとえば図2Aに示されたようなオブジェクト冗長化格納システム200で実施され得る。図9を参照すると、図8に示された「ゴシップ」方法を採用することに代えて、コーホート内の各小ボリュームがコーホート内の他の小ボリュームと互いに通信して、たとえば照合調整モジュールのような中央部と通信する小ボリューム間の差分を判定する。各小ボリューム920は、2つ以上の共通のオブジェクトリスト926のセットを、コーホート内のその他の小ボリューム920のうち1つに対応する各々共通のオブジェクトリスト926とともに、周期的にまたは非周期的に生成する。そして、ハッシュ機能950が共通のオブジェクトリスト926の各々に適用され、たとえば照合調整モジュール990のような中央部に、ハッシュ値を提供する。

40

【0118】

50

少なくともいくつかの実施形態では、共通のオブジェクトリスト 926 は、コーホート内の他の N 個の小ボリューム 920 の各々について、コーホート内の各小ボリューム 920 上に生成され、ハッシュ値は、共通のオブジェクトリスト 926 ごとに N 個の小ボリューム 920 の各々に生成され、すべてのハッシュ値が、照合調整モジュール 990 に提供される。

【0119】

図 9 は、小ボリューム 920 A が一例として示される。小ボリューム 920 A の要素 922 及びメタデータ 924 から、互いにコーホート内の他の小ボリューム 920 の各々についての 1 つの共通のオブジェクトリストとともに、共通のオブジェクトリスト 926 のセットが生成され得る。あるいは、小ボリューム 920 A は、変換式（たとえば、ハッシュ）を、それが格納しているデータ要素 922 のオブジェクト ID に適用して、各データ要素 922 もまた格納されているはずであるその他の小ボリューム 920 を判定してもよく、またこの情報を用いて、共通のオブジェクトリスト 926 を生成してもよい。ハッシュ機能 950 は、各共通のオブジェクトリスト 926 に適用され、結果として得られたハッシュ値は、照合調整モジュール 990 に送られる。少なくともいくつかの実施形態では、図 8 を参照して説明されたハッシュ木手法を用いて、共通のオブジェクトリスト 926 についてのハッシュ値を生成してもよい。しかしながら、いくつかの実施形態では、他のハッシュ手法を用いてもよい。

10

【0120】

コーホート内の他の小ボリューム 920 B ~ 920 N の各々は、小ボリューム 920 A が、それらに共通のオブジェクトリストについてのハッシュ値を照合調整モジュール 990 に提供するために示されたものと同様の方法を行ってもよい。

20

【0121】

図 9 に示されるように、照合調整モジュール 990 は、小ボリューム 920 から受信されたハッシュ値のハッシュ比較 960 をローカルに行ってもよく、それぞれの共通のオブジェクトリスト 926 についての異なるハッシュ値を有することが判定された任意の 2 つの小ボリューム 920 に対する差分判定 970 を行ってもよい。差分判定 970 は、所与の 2 つの小ボリューム 920 に対して行われ、両方の小ボリューム 920 にあるはずであるがそこにはないデータオブジェクトについての 1 つ以上のオブジェクト ID を判定し得る。少なくともいくつかの実施形態では、差分判定 970 を行うために、照合調整モジュール 990 は、小ボリューム 920 の一方または両方からのさらなるハッシュ情報（たとえば、サブリストハッシュ）及び/またはさらなるオブジェクト ID 情報（たとえば、オブジェクト ID のリストまたは部分的なリスト）を要求する必要がある。あるいは、いくつかの実施形態では、照合調整モジュール 990 が、2 つの小ボリューム 920 に対して異なるハッシュ値を検出した場合、照合調整モジュール 990 は、小ボリューム 920 の一方または両方が、差分判定の少なくともいくつかを行うことを要求する必要がある。

30

【0122】

ハッシュ比較 960 及び差分判定 970 の結果に従って、照合調整モジュール 990 は、コーホート内の小ボリューム 920 の 1 つ、2 つ、またはそれ以上で照合調整を行い、図 9 に示されたような比較方法によって判定された、小ボリューム 920 上にあるはずであるデータオブジェクトのセットを有する小ボリューム上の実際にあるデータオブジェクトを照合調整し得る。

40

【0123】

図 11 A ~ 図 11 C は、少なくともいくつかの実施形態による、オブジェクト冗長化格納システム上での照合調整プロセスの一部として、コーホートの小ボリュームを比較するための代替方法のフローチャートである。図 11 A 及び図 11 C は、たとえば図 2 A に示されたようなオブジェクト冗長化格納システム 200 に示されたような格納ノード及び/または小ボリュームの各々で行われることがある。一方、図 11 B は、たとえば図 2 A に示されたような格納サービス 250 の照合調整プロセスまたはモジュールの中央部で、またはそれによって行われ得る。図 10 の方法に示されたような、コーホート内の各小ボリ

50

ュームが、コーホート内の他の小ボリュームと各々通信して、図 1 1 A ~ 図 1 1 C に示されたような方法で、中央部たとえば照合調整モジュールと通信する小ボリューム間の差分を判定する「ゴシップ」方法を採用することに代えて、小ボリュームは、1 つ以上のまたはそれ以上の共通のオブジェクトリストを、コーホート内の他の小ボリュームのうち 1 つに対応する各共通のオブジェクトリストとともに、周期的にまたは非周期的に生成する。そして、ハッシュ機能は、共通のオブジェクトリスト各々に適用され、ハッシュ値は、コーホート内の小ボリュームからハッシュ値を収集する中央部、たとえば照合調整モジュールに提供され、必要に応じてハッシュ比較、差分判定、及び照合調整を行う。中央部は、小ボリュームからさらなる情報を要求することがあるが、必ずしもそうではない。このように、小ボリュームは、中央部からさらなる情報についての要求を受信することがあるが、必ずしもそうではない。

10

【 0 1 2 4 】

図 1 1 A は、少なくともいくつかの実施形態による、コーホートの各小ボリュームで行われ得る、共通のオブジェクトリストを判定してハッシュを生成するための方法の高レベルのフローチャートである。図 1 1 A の 1 1 0 0 に示されるように、小ボリュームは、コーホート内の他の小ボリュームのうち少なくとも 1 つについての共通のオブジェクトリストを判定または生成し得る。その他の小ボリュームのうち所与の 1 つについての共通のオブジェクトリストは、この小ボリュームがその他の小ボリュームと共通して有しているかまたは有しているはずであるデータオブジェクトを示す。少なくともいくつかの実施形態では、各データオブジェクトは、オブジェクト識別子 (オブジェクト ID) によるリストに示されることがある。1 1 0 2 に示されるように、小ボリュームは、共通のオブジェクトリストの各々のハッシュを生成し得る。少なくともいくつかの実施形態では、図 8 及び 9 を参照して説明されたようなハッシュ木手法を用いて、共通のオブジェクトリストについてのハッシュ値 (すなわち、ルートハッシュ) を生成してもよい。しかしながら、いくつかの実施形態では、他のハッシュ手法を用いてもよい。1 1 0 4 に示されるように、小ボリュームは、たとえば照合調整モジュールのような中央部、に、生成されたハッシュを送ることがある。1 1 0 4 から 1 1 0 0 への戻り矢印によって示されるように、各小ボリュームは、図 1 1 A に示されるように、本方法を周期的にまたは非周期的に繰り返すことがある。

20

【 0 1 2 5 】

図 1 1 B は、少なくともいくつかの実施形態による、たとえば照合調整モジュールのような中央部、またはプロセスで行われるか、またはそれによって行われる方法の高レベルのフローチャートである。1 1 1 0 に示されるように、中央部は、コーホート内の小ボリュームのすべてからハッシュ値を収集してもよく、周期的にまたは非周期的に (たとえば、ハッシュ値がコーホート内のすべての小ボリュームから受信されたときに) ハッシュ比較を行ってもよい。ハッシュ比較に基づいて、差分判定及び照合調整が必要に応じて行われてもよい。少なくともいくつかの実施形態では、差分判定の一部及び/または照合調整の一部として、中央部は、1 つ以上の小ボリュームからさらなる情報が必要とされることを判断し得る。1 1 1 2 では、さらなる情報を要する場合、中央部はその後、1 1 1 4 に示されるように、小ボリュームからのさらなる情報を要求することがある。たとえば、照合調整モジュールは、1 つ以上の小ボリュームからのさらなるハッシュ情報 (たとえば、サブリストハッシュ) 及び/またはさらなるオブジェクト ID 情報 (たとえば、オブジェクト ID のリストまたは部分的なリスト) を要求することがある。1 1 1 2 及び 1 1 1 4 から 1 1 1 0 への戻り矢印によって示されるように、中央部は、図 1 1 B に示されたような方法を、周期的にまたは非周期的に繰り返すことがある。

30

40

【 0 1 2 6 】

図 1 1 C は、少なくともいくつかの実施形態による、コーホートの各小ボリュームで行われ得る、さらなる情報に対する要求を扱うための方法の高レベルのフローチャートである。図 1 1 C の方法は、図 1 1 A の方法と非同期的に行われ得ることに留意すべきである。図 1 1 C の 1 1 2 0 では、小ボリュームがさらなる情報についての要求を受信すると、

50

その後 1 1 2 2 において、要求された情報が判定され、1 1 2 4 において照合調整モジュールに報告される。

【 0 1 2 7 】

(例証的なシステム)

少なくともいくつかの実施形態では、本明細書に記載されたオブジェクト冗長化格納システム内のボリュームコーホートのための方法及び機器の一部またはすべてを組み込むサーバは、1つ以上のコンピュータアクセス可能媒体にアクセスすることを含むかまたはそのように構成された汎用コンピュータシステム、たとえば図 1 3 に示されたコンピュータシステム 2 0 0 0 を含む得る。図示された実施形態では、コンピュータシステム 2 0 0 0 は、入出力 (I / O) インターフェース 2 0 3 0 を介してシステムメモリ 2 0 2 0 に接続された1つ以上のプロセッサ 2 0 1 0 を含む。コンピュータシステム 2 0 0 0 は、I / O インターフェース 2 0 3 0 に接続されたネットワークインターフェース 2 0 4 0 をさらに含む。

10

【 0 1 2 8 】

様々な実施形態では、コンピュータシステム 2 0 0 0 は、1つのプロセッサ 2 0 1 0 を含むユニプロセッサシステム、または数個の (たとえば、2、4、8、または別の好適な数の) プロセッサ 2 0 1 0 を含むマルチプロセッサシステムであり得る。プロセッサ 2 0 1 0 は、命令を実行することが可能である任意の好適なプロセッサであり得る。たとえば、様々な実施形態では、プロセッサ 2 0 1 0 は、x 8 6、Power PC、SPARC、または MIPS ISA 等の多様な命令セットアーキテクチャ (ISA) のいずれか、または任意の他の好適な ISA を実施する汎用または組み込み型プロセッサであってもよい。マルチプロセッサシステムでは、プロセッサ 2 0 1 0 の各々が、共通して同じ ISA を実施し得るが、必ずしもそうではない。

20

【 0 1 2 9 】

システムメモリ 2 0 2 0 は、プロセッサ 2 0 1 0 によってアクセス可能である命令及びデータを格納するように構成され得る。様々な実施形態では、システムメモリ 2 0 2 0 は、任意の好適なメモリ技術、たとえばスタティックランダムアクセスメモリ (SRAM)、シンクロナスダイナミック RAM (SDRAM)、不揮発性 / フラッシュ型メモリ、または任意の他のタイプのメモリを用いて実施し得る。図示された実施形態では、オブジェクト冗長化格納システム内のボリュームコーホートに対して、1つ以上の所望の機能、たとえば上述された方法、手法、及びデータを実施するプログラム命令及びデータが、システムメモリ 2 0 2 0 内部にコード 2 0 2 5 及びデータ 2 0 2 6 として格納されて示される。

30

【 0 1 3 0 】

1つの実施形態では、I / O インターフェース 2 0 3 0 は、プロセッサ 2 0 1 0 と、システムメモリ 2 0 2 0 と、ネットワークインターフェース 2 0 4 0 または他の周辺インターフェースを含むデバイス内の任意の周辺機器との間の I / O トラフィックを調整するように構成され得る。いくつかの実施形態では、I / O インターフェース 2 0 3 0 は、必要に応じてプロトコル、タイミング、または他のデータ変換式を行い、データ信号を1つのコンポーネント (たとえば、システムメモリ 2 0 2 0) から、別のコンポーネント (たとえば、プロセッサ 2 0 1 0) で用いられるために好適な形式に変換し得る。いくつかの実施形態では、I / O インターフェース 2 0 3 0 は、様々なタイプの周辺バスを通して接続されたデバイスのためのサポート、たとえば周辺コンポーネント相互接続 (PCI) バス規格またはユニバーサルシリアルバス (USB) 規格の変形等を含む得る。いくつかの実施形態では、I / O インターフェース 2 0 3 0 の機能は、2つ以上の別個のコンポーネント、たとえばノースブリッジ及びサウスブリッジ等に分けることがある。また、いくつかの実施形態では、I / O インターフェース 2 0 3 0 の機能性のいくつかまたはすべて、たとえばシステムメモリ 2 0 2 0 に対するインターフェースが、プロセッサ 2 0 1 0 内に直接組み入れられてもよい。

40

【 0 1 3 1 】

50

ネットワークインターフェース2040は、コンピュータシステム2000と、ネットワーク2050に接続された他のデバイス2060、たとえば図1～図12に示されたような他のコンピュータシステムまたはデバイス等との間で、データが交換されることが可能になるように構成され得る。たとえば、様々な実施形態では、ネットワークインターフェース2040は、任意の好適な有線または無線汎用データネットワーク、たとえば一種のイーサネットネットワーク等を介した通信をサポートし得る。加えて、ネットワークインターフェース2040は、アナログ音声ネットワークまたはデジタルファイバ通信ネットワーク等の電気通信/電話通信ネットワークを介した、ファイバチャネルSAN等のストレージエリアネットワークを介した、または任意の他の好適なタイプのネットワーク及び/またはプロトコルを介した通信をサポートし得る。

10

【0132】

いくつかの実施形態では、システムメモリ2020は、図1～図12について上述されたようなオブジェクト冗長化格納システム内のボリュームコーホートの実施形態を実施するためのプログラム命令及びデータを格納するように構成されたコンピュータアクセス可能媒体の1つの実施形態であってもよい。しかしながら、他の実施形態では、プログラム命令及び/またはデータは、異なるタイプのコンピュータアクセス可能媒体で受信され、送信され、または格納されてもよい。一般に、コンピュータアクセス可能媒体は、磁気媒体または光媒体等の非一時的記憶媒体またはメモリ媒体、たとえば、I/Oインターフェース2030を介してコンピュータシステム2000に接続されたディスクまたはDVD/CDを含み得る。非一時的コンピュータアクセス可能記憶媒体は、いくつかの実施形態

20

【0133】

(結論)

様々な実施形態は、コンピュータアクセス可能媒体上において、前述の説明に従って実施される命令及び/またはデータを受信すること、送ること、または格納することをさらに含み得る。一般には、コンピュータアクセス可能媒体は、磁気媒体または光媒体等の記憶媒体またはメモリ媒体、ディスクまたはDVD/CD-ROM、RAM(たとえばSDRAM、DDR、RDRAM、SRAM等)、ROM等の揮発性または不揮発性媒体に加えて、ネットワーク及び/または無線リンク等の通信媒体を介して伝達される、電気信号、電磁信号またはデジタル信号等の伝送媒体または信号を含み得る。

30

【0134】

図面に示され本明細書に記載されたような様々な方法は、方法の代表的な実施形態を表す。本方法は、ソフトウェア、ハードウェア、またはそれらの組み合わせにおいて実施されてもよい。方法の順序が変更されてもよく、様々な要素の付加し、並べ替え、組み合わせ、省略、改変等がなされてもよい。

40

【0135】

本開示の利益を有する当業者には明白であるとおり、様々な改変及び変更がなされてもよい。したがって、すべてのそのような改変及び変更を包含し、上述の説明が制限的な意味ではなく例示的な意味で解釈されることが意図される。

【0136】

本開示の実施形態は、以下の項目を考慮して説明されることができる。

1. 複数の格納ノードと、
- 1つ以上のプロセッサと、

オブジェクト冗長化格納システムを組み込むために、1つ以上のプロセッサの少なくとも

50

も1つによって実行可能であるプログラム命令を格納するメモリであって、当該システムが、

複数の格納ノードのN個におよぶコーホートを確立し、

格納システムに格納されるデータオブジェクトを受信するように構成され、

データオブジェクトごとに、

オブジェクト冗長化手法によって、データオブジェクトからN未満のM個のオブジェクト冗長化データ要素を生成し、

選択手法によって、コーホート内のN個の格納ノードの中からM個の格納ノードを選択し、

選択されたM個の格納ノードにM個のデータ要素を格納することであって、M個のデータ要素の1つが、M格納ノードの各々に格納される、メモリと、

を含み、

生成すること、選択すること、及び格納することが、受信されたデータオブジェクトについて生成されたデータ要素を、コーホート内のN個の格納ノードにわたって分配し、N個の格納ノードの所与の2つが、異なるデータ要素のセットを含む、システム。

2．オブジェクト冗長化手法がレプリケーション手法であり、オブジェクト冗長化データ要素がデータオブジェクトのレプリカである、項目1に記載のシステム。

3．オブジェクト冗長化手法が、イレージャ符号化手法であり、オブジェクト冗長化データ要素が、データオブジェクトから生成されたシャードであり、冗長な符号化手法によって所与のデータオブジェクトについて生成された、M個のシャードの少なくともR個のサブセットが、それぞれのデータオブジェクトを再生するために必要とされる、項目1に記載のシステム。

4．オブジェクト冗長化格納システムが、データオブジェクトごとに、

オブジェクトについて生成されたM個のデータ要素のうち少なくとも1つについて、M個のデータ要素の1つ以上の他のものが格納されたコーホート内の1つ以上の位置を示すメタデータを生成し、

選択されたM個の格納ノード上の少なくとも1つのデータ要素それぞれとともにメタデータを格納する

ようにさらに構成される、項目1に記載のシステム。

5．オブジェクト冗長化格納システムが、

格納システムに格納されたデータオブジェクトについての要求を受信し、

コーホート内のN個の格納ノードのサブセットを選択し、

要求されたデータオブジェクトに対応するデータ要素について、格納ノードの選択されたサブセットに問い合わせる、

ようにさらに構成される、項目1に記載のシステム。

6．オブジェクト冗長化格納システムが、

N個の格納ノードコーホート内の別のサブセットを選択し、

要求されたデータオブジェクトに対応するデータ要素について、格納ノードのその他のサブセットに問い合わせる、

ようにさらに構成される、項目5に記載のシステム。

7．格納ノードの問い合わせを受けたサブセットから、要求されたデータオブジェクトに対応する少なくとも1つのデータ要素を受信するようにさらに構成される、項目5に記載のシステム。

8．オブジェクト冗長化格納システムが、

要求されたデータオブジェクトに対応する少なくとも1つのさらなるデータ要素を要することを判定し、

コーホート内の少なくとも1つのさらなる格納ノードに、少なくとも1つのさらなるデータ要素について問い合わせる

ようにさらに構成される、項目7に記載のシステム。

9．オブジェクト冗長化格納システムが、

10

20

30

40

50

格納ノードの問い合わせを受けたサブセットから受信された少なくとも1つのデータ要素とともに、要求されたデータオブジェクトについてのM個のデータ要素の他の1つ以上が格納された、コーホート内の1つ以上の位置を示すメタデータを受信するようにさらに構成され、

コーホート内の少なくとも1つのさらなる格納ノードに、少なくとも1つのさらなるデータ要素について問い合わせるために、オブジェクト冗長化格納システムが、受信されたメタデータから少なくとも1つのさらなる格納ノードを判定するようにさらに構成される、項目8に記載のシステム。

10 . 選択手法が、コーホート内のN個の格納ノードの中から、M個の格納ノードの少なくとも1つをランダムに選択する、項目1に記載のシステム。

11 . 選択手法が、それぞれのデータオブジェクトの識別情報によって、M個の格納ノードの少なくとも1つを選択する、項目1に記載のシステム。

12 . 1つ以上のコンピューティングデバイス上に組み込まれた格納サービスによって、オブジェクト冗長化格納システムに格納されるデータオブジェクトを受信することであって、オブジェクト冗長化格納システムが、N個の格納ノードにおよぶコーホートを含む、受信することを含む、

受信されたデータオブジェクトごとに、

オブジェクト冗長化手法によって、データオブジェクトからN未満のM個のオブジェクト冗長化データ要素を生成し、

コーホート内のN個の格納ノードの中からM個の格納ノードを選択し、

選択されたM個の格納ノードにM個のデータ要素を格納することであって、M個のデータ要素の1つは、M個の格納ノードの各々に格納され、

M個のデータ要素の1つ以上が格納されるコーホート内の1つ以上の位置を示すメタデータを格納する、方法。

13 . M個のデータ要素の1つ以上が格納されるコーホート内の1つ以上の位置を示すメタデータを格納することは、メタデータが、M個のデータ要素の1つ以上の他のデータ要素が格納されるコーホート内の1つ以上の位置を示すM個のデータ要素少なくとも1つとともに格納されることを含む、項目12に記載の方法。

14 . オブジェクト冗長化手法が、レプリケーション手法またはイレージャコード化手法の一方である、項目12に記載の方法。

15 . 選択手法が、M個の格納ノードの少なくとも1つをランダムに選択する手法か、またはそれぞれのデータオブジェクトの識別情報によって、M個の格納ノードの少なくとも1つを選択する手法かの一方である、項目12に記載の方法。

16 . 格納システムに格納されたデータオブジェクトについての要求を受信することと、

コーホート内のN個の格納ノードのサブセットを選択することと、

要求されたデータオブジェクトに対応するデータ要素のための格納ノードの選択されたサブセットについて問い合わせることと、

格納ノードの問い合わせを受けたサブセットから、要求されたデータオブジェクトに対応するデータ要素の少なくとも1つを受信することと、

をさらに含む、項目12に記載の方法。

17 . 要求されたデータオブジェクトに対応する少なくとも1つのさらなるデータ要素を要することを判定することと、

1つ以上のM個のデータ要素が格納されたコーホート内の1つ以上の位置を示すメタデータから、要求されたデータオブジェクトに対応するデータ要素を含む少なくとも1つのさらなる格納ノードを判定することと、

コーホート内の少なくとも1つのさらなる格納ノードに、少なくとも1つのさらなるデータ要素について問い合わせることと、

をさらに含む、項目16に記載の方法。

18 . 1つ以上のコンピュータ上で実行可能なプログラム命令を格納して、N個の格納

10

20

30

40

50

ノードのセットに冗長に格納されるデータオブジェクトを受信するように構成されたオブジェクト冗長化格納システムを実施する、非一時的なコンピュータアクセス可能記憶媒体であって、

受信されたデータオブジェクトごとに、

オブジェクト冗長化手法によって、データオブジェクトから N 未満の M 個のオブジェクト冗長化データ要素を生成し、

N 個の格納ノードの中から M 個の格納ノードを選択し、

選択された M 個の格納ノードに M 個のデータ要素を格納し、 M 個のデータ要素の 1 つは M 個の格納ノードの各々に格納され、

M 個のデータ要素のうち少なくとも 1 つとともに、 M 個のデータ要素ごとに、それぞれのデータ要素が格納された位置を示すメタデータを格納する、非一時的なコンピュータアクセス可能記憶媒体。 10

19. オブジェクト冗長化手法が、レプリケーション手法またはイレージャコード化手法の一方である、項目 18 に記載の非一時的なコンピュータアクセス可能記憶媒体。

20. 選択手法が、 M 個の格納ノードの少なくとも一方をランダムに選択する手法か、またはそれぞれのデータオブジェクトの識別情報によって、 M 個の格納ノードのうち少なくとも 1 つを選択する手法かのいずれか一方である、項目 18 に記載の非一時的なコンピュータアクセス可能記憶媒体。

21. オブジェクト冗長化格納システムが、

格納システムに格納されたデータオブジェクトについての要求を受信し、 20

N 個の格納ノードのサブセットを選択し、

要求されたデータオブジェクトに対応するデータ要素のための格納ノードの選択されたサブセットについて問い合わせ、

要求された格納ノードのサブセットから、要求されたデータオブジェクトに対応するデータ要素の少なくとも 1 つを受信する、

ようにさらに構成される、項目 18 に記載の非一時的なコンピュータアクセス可能記憶媒体。

22. オブジェクト冗長化格納システムが、

要求されたデータオブジェクトに対応する少なくとも 1 つのさらなるデータ要素を要することを判断し、 30

データ要素とともに格納されたメタデータから、要求されたデータオブジェクトに対応するデータ要素を含む少なくとも 1 つのさらなる格納ノードを判定し、

少なくとも 1 つのさらなる格納ノードに、少なくとも 1 つのさらなるデータ要素について問い合わせる、

ようにさらに構成される、項目 21 に記載の非一時的なコンピュータアクセス可能記憶媒体。

【0137】

本開示のさらなる実施形態は、以下の項目を考慮して説明されることができる。

1. 複数の格納ノードと、

格納サービスを実施するための 1 つ以上のデバイスであって、 40

複数の格納ノードにおよぶコーホートを確立し、

コーホートにデータオブジェクトを格納することであって、各データオブジェクトから生成されたオブジェクト冗長化データ要素のセットは、コーホート内の格納ノードの選択されたサブセットに格納される、

ように構成されたデバイスとを含み、

コーホート内の格納ノードの各々が、

コーホート内の格納ノードのうち別の 1 つと共通して格納ノードが有しているはずである格納ノード上のデータオブジェクトのリストを生成し、

データオブジェクトのリストについてのハッシュ値を生成し、

その他の格納ノード上のデータオブジェクトのリストについてのハッシュ値を受信し、 50

格納ノード上のデータオブジェクトのリストについてのハッシュ値が、その他の格納ノードから受信されたハッシュ値と一致しないことを判定し、

ハッシュ値が一致しないことの判定に回答して、格納サービスの照合調整プロセスに、2つの格納ノードについてのハッシュ値が一致しないことを通知する、ように構成される、システム。

2. ハッシュ値が一致しないことの判定に回答して、格納ノードが、格納ノードの両方にあるはずであるがそこにはない1つ以上のデータオブジェクトを判定し、格納サービスの照合調整プロセスに、判定された1つ以上のデータオブジェクトを報告するようにさらに構成される、項目1に記載のシステム。

3. 通知に回答して、照合調整プロセスが、格納ノードの両方にあるはずであるがそこにはない1つ以上のデータオブジェクトを判定するように構成される、項目1に記載のシステム。

4. 照合調整プロセスが、コーホート内の格納ノードの1つ以上から受信された通知に回答して、コーホート内の格納ノードに格納されたデータオブジェクトを調整するように構成され、各通知が、ハッシュ値が一致しない格納ノードの特定された2つを示す、項目1に記載のシステム。

5. コーホート内の格納ノードのうち別の1つと共通して有しているはずである格納ノード上のデータオブジェクトのリストを生成するために、格納ノードは、格納ノード上のデータ要素に対応する情報によってリストを生成するように構成され、情報が、格納ノード上のデータ要素ごとに、それぞれのデータオブジェクトについて生成されたその他のデータ要素が格納されたコーホート内の1つ以上の位置を示す、項目1に記載のシステム。

6. データオブジェクトのリストについてのハッシュ値を生成するために、格納ノードが、データオブジェクトのリストからハッシュ木を生成するハッシュ手法によって、データオブジェクトのリストからハッシュ値を生成するように構成され、ハッシュ値が、ハッシュ木のルートハッシュである、項目1に記載のシステム。

7. 格納ノードの両方にあるはずであるがそこにはない1つ以上のデータオブジェクトを判定するために、格納ノードが、ハッシュ木内の情報を分析して、その他の格納ノード上の特定のサブセットとは異なる、格納ノード上のデータオブジェクトの特定のサブセットを検索するように構成される、項目6に記載のシステム。

8. データオブジェクトのリストが、各データオブジェクトを識別するオブジェクト識別子を含み、データオブジェクトのリストについてのハッシュ値を生成するために、格納ノードが、

オブジェクト識別子によって、リストを2つ以上の分類されたサブリストに分け、

分類されたサブリストごとに1つ以上のハッシュ値を生成し、各ハッシュ値は、それぞれの分類されたサブリスト内のデータオブジェクトのうち1つ以上のオブジェクト識別子から生成され、

分類されたサブリストごとに1つ以上のハッシュ値を組み合わせるハッシュ値を計算して、2つ以上のサブリストハッシュを生成し、

組み合わせられたサブリストハッシュを組み合わせるハッシュ値を計算して、データオブジェクトのリストについてのルートハッシュを生成し、データオブジェクトのリストについてのハッシュ値は、ルートハッシュである、

ように構成される、項目1に記載のシステム。

9. オブジェクト冗長化データ要素のセットが格納されたコーホート内の格納ノードの各サブセットが、コーホート内の複数の格納ノードのなかから、サブセット内の格納ノードの少なくとも1つをランダムに選択する選択手法に従って選択される、項目1に記載のシステム。

10. 1つ以上のデバイスに組み込まれた格納サービスによって、データオブジェクトから生成されたデータ要素についてのセットを、コーホート内の複数の格納ノードに組み込まれた複数の小ボリュームの選択されたサブセットに格納することと、

2つの小ボリュームの各々において、その他の小ボリュームについての共通のオブジェ

10

20

30

40

50

クトリストを生成することであって、小ボリュームのうち1つでの共通のオブジェクトリストが、その他の小ボリュームにも格納されているはずである、小ボリューム上のデータオブジェクトを示す、生成することと、

2つの小ボリュームの各々において、小ボリュームにおける共通のオブジェクトリストについてのハッシュ値を生成することと、

2つの小ボリューム上に生成されたハッシュ値が一致していないことを判定することと、

判定に応答して、小ボリュームの両方に格納されているはずであるが格納されていない1つ以上のデータオブジェクトを識別することと、

を含む方法。

10

11. 判断すること及び識別することが、2つの小ボリュームの各々で行われ、方法が、格納サービスの照合調整プロセスに、識別された1つ以上のデータオブジェクトを報告することをさらに含む、項目10に記載の方法。

12. 方法が、格納サービスの照合調整プロセスに、生成されたハッシュ値を提供することをさらに含み、照合調整プロセスが、判断すること及び識別することを行う、項目10に記載の方法。

13. 2つの小ボリュームを照合調整して、2つの小ボリュームが各々識別された1つ以上のデータオブジェクトから生成されたデータ要素を格納するようにする、格納サービスの照合調整プロセスをさらに含む、項目10に記載の方法。

14. オブジェクト冗長化データ要素のセットが、オブジェクト冗長化手法によってデータオブジェクトから生成され、オブジェクト冗長化手法が、レプリケーション手法またはイレージャ符号化手法のいずれか一方である、項目10に記載の方法。

20

15. 小ボリュームにおいて共通のオブジェクトリストを生成することが、小ボリューム上にデータ要素とともに格納されたメタデータによって、共通のオブジェクトリストを生成することを含み、小ボリューム上の所与のデータ要素についてのメタデータが、データ要素の生成されたセットにおける他のデータ要素が格納される1つ以上の他の小ボリュームを示す、項目10に記載の方法。

16. 小ボリュームにおいて共通のオブジェクトリストを生成することが、格納ノード上のデータオブジェクトのオブジェクト識別子によって、共通のオブジェクトリストを生成することを含み、オブジェクト識別子の変換式が、それぞれのデータオブジェクトについて生成された他のデータ要素が格納されるコーホート内の1つ以上の位置を示す、項目10に記載の方法。

30

17. 小ボリュームにおいて共通のオブジェクトリストについてのハッシュ値を生成することが、共通のオブジェクトリスト内のオブジェクト識別子からハッシュ木を生成することを含み、ハッシュ値が、ハッシュ木のルートハッシュである、項目10に記載の方法。

18. 1つ以上のコンピュータ上で実行可能なプログラム命令を格納する非一時的なコンピュータアクセス可能記憶媒体であって、

データオブジェクトから生成されたデータ要素のセットを、コーホート内の複数の格納ノード上に組み込まれた複数の小ボリュームの選択されたサブセットに格納することと、

40

小ボリュームの対についての共通のオブジェクトリストを生成することであって、対の中の1つの小ボリューム上の共通のオブジェクトリストが、対の中のその他の小ボリューム上にも格納されているはずである、小ボリューム上のデータオブジェクトを示す、生成することと、

対の中の小ボリュームの各々において、共通のオブジェクトリストについてのハッシュ値を生成することと、

2つの小ボリューム上で生成されたハッシュ値が一致しないことを判定することと、

判定に応答して、小ボリュームの両方に格納されているはずであるがそこにはない1つ以上のデータオブジェクトを識別することと、

を実施する、非一時的なコンピュータアクセス可能記憶媒体。

50

19．小ボリュームで共通のオブジェクトリスト生成する際に、小ボリューム上のデータ要素とともに格納されたメタデータに従って、共通のオブジェクトリストの生成を実施するためのプログラム命令をさらに実行可能であり、小ボリューム上の所与のデータ要素についてのメタデータが、データ要素の生成されたセット内の他のデータ要素が格納された1つ以上の他の小ボリュームを示す、項目18に記載の非一時的なコンピュータアクセス可能記憶媒体。

20．共通のオブジェクトリストが、データオブジェクトについてのオブジェクト識別子を含み、対の中の小ボリュームの各々において、共通のオブジェクトリストについてのハッシュ値を生成する際に、

オブジェクト識別子によって、共通のオブジェクトリストを2つ以上の分類されたサブリストに分けることと、

分類されたサブリストごとに1つ以上のハッシュ値を生成することであって、各ハッシュ値が、それぞれの分類されたサブリスト内の1つ以上のデータオブジェクトのオブジェクト識別子から生成される、生成することと、

分類されたサブリストごとに1つ以上のハッシュ値を組み合わせるハッシュ値を計算して、2つ以上のサブリストハッシュを生成することと、

組み合わせられたサブリストハッシュを組み合わせるハッシュ値を計算して、データオブジェクトのリストについてのルートハッシュを生成し、共通のオブジェクトのリストについてのハッシュ値が、ルートハッシュであることと、

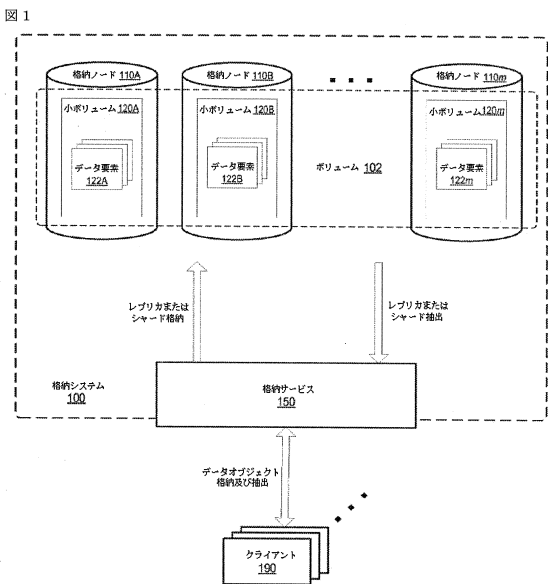
を実施するように、プログラム命令をさらに実行可能である、項目18に記載の非一時的なコンピュータアクセス可能記憶媒体。

21．オブジェクト識別子の分類されたサブリストについてのハッシュ値、サブリストハッシュ、及びルートハッシュが、2つ以上のレベルを有するハッシュ木を形成し、小ボリュームの両方に格納されているはずであるがそこにはない1つ以上のデータオブジェクトを識別する際に、ハッシュ木による1つ以上のデータオブジェクトの検索を実施するために、プログラム命令をさらに実行可能である、項目20に記載の非一時的なコンピュータアクセス可能記憶媒体。

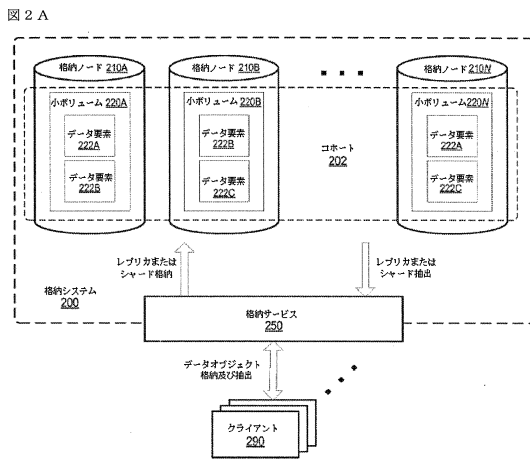
10

20

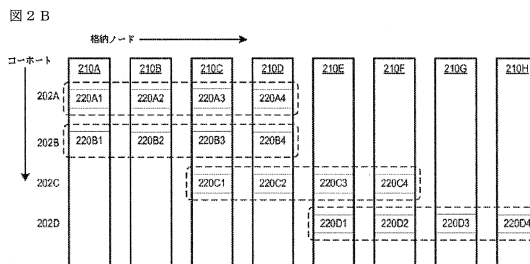
【図1】



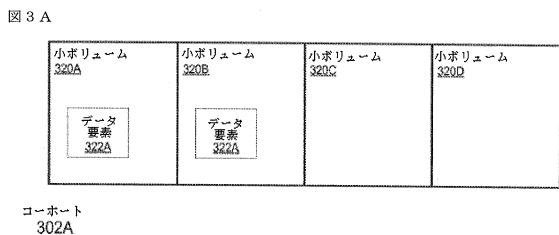
【図2A】



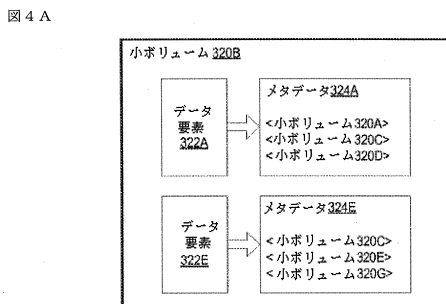
【図2B】



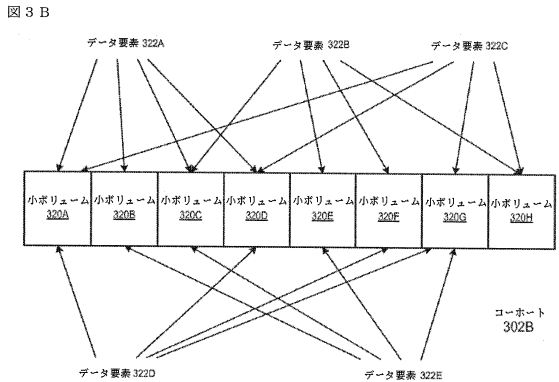
【図3A】



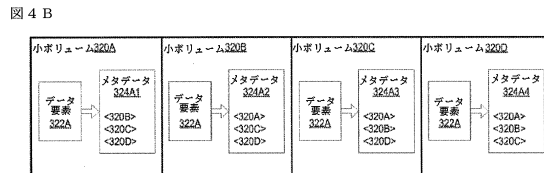
【図4A】



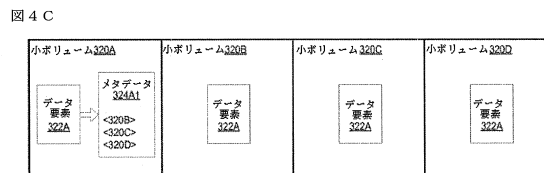
【図3B】



【図4B】

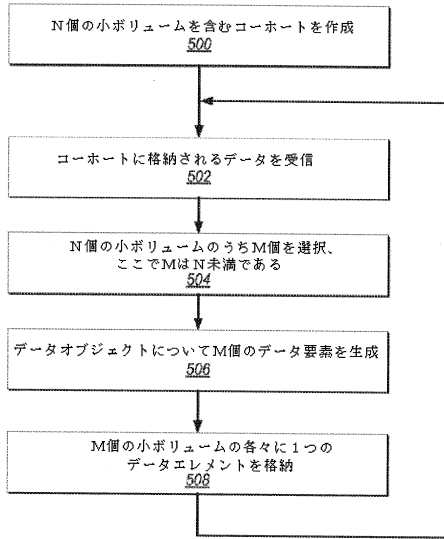


【図4C】



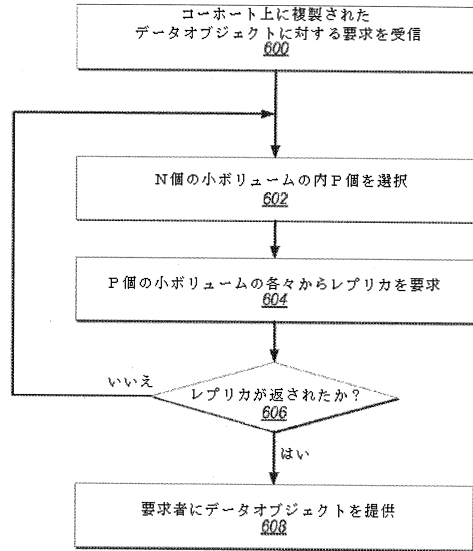
【図5】

図5



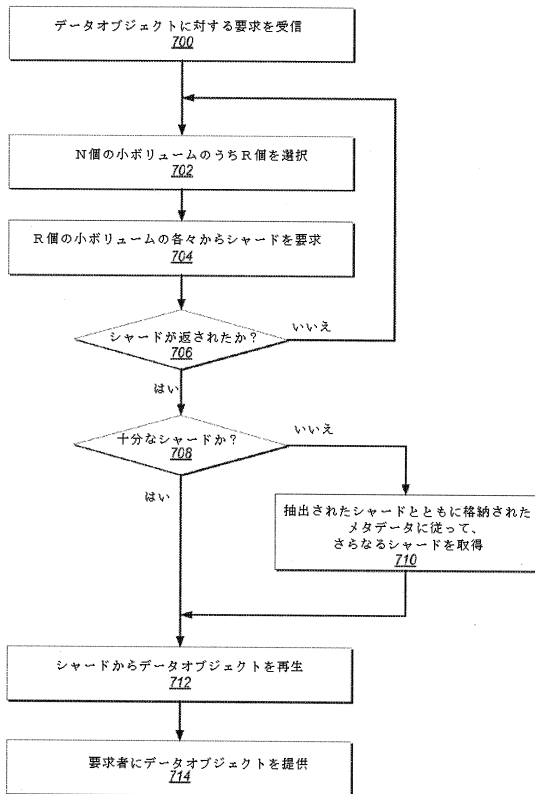
【図6】

図6



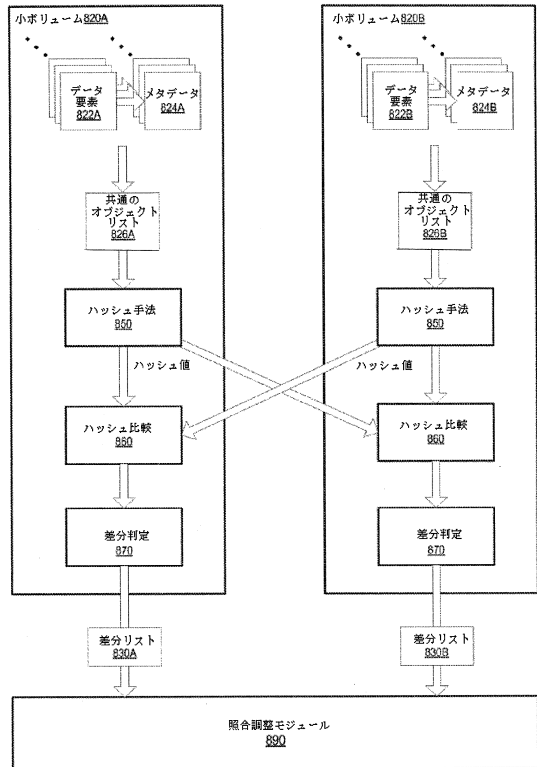
【図7】

図7



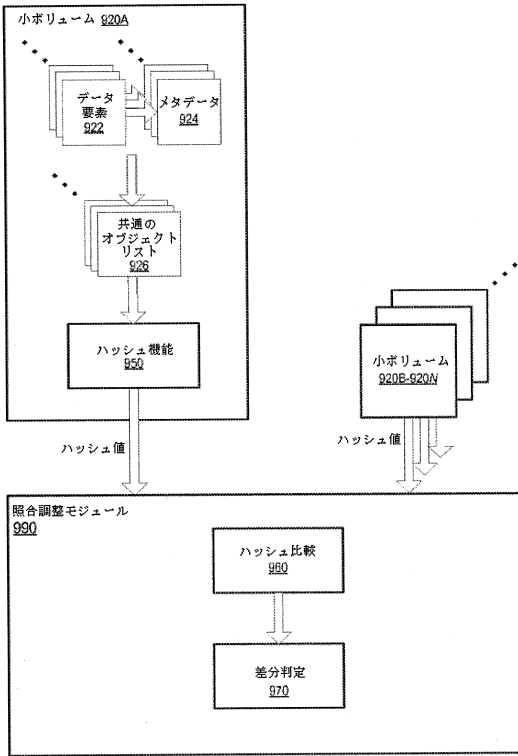
【図8】

図8



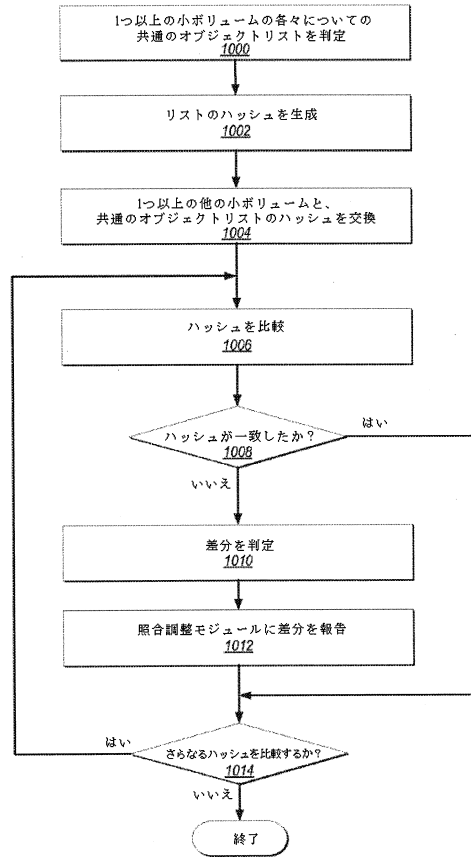
【図9】

図9



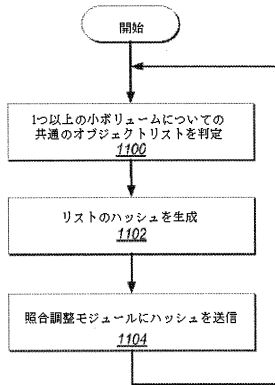
【図10】

図10



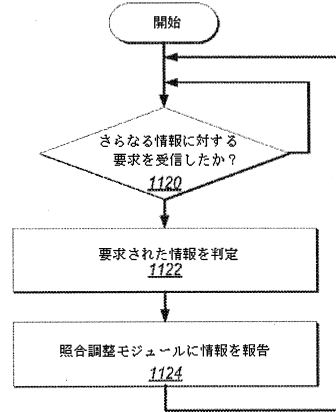
【図11A】

図11A



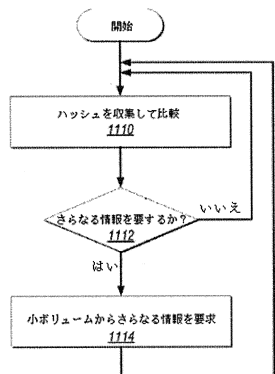
【図11C】

図11C



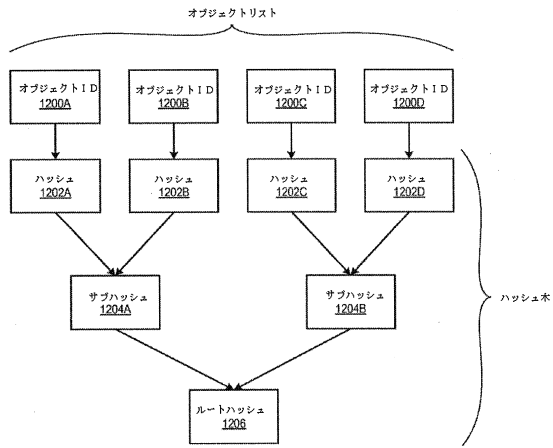
【図11B】

図11B



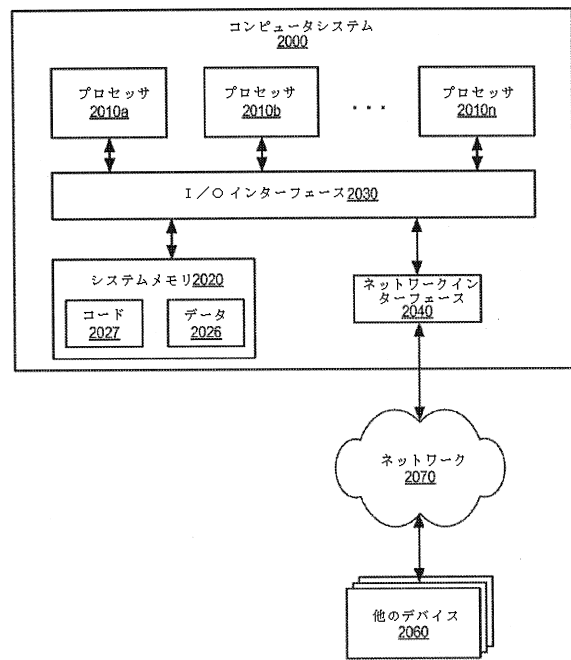
【 図 1 2 】

図 1 2



【 図 1 3 】

図 1 3



フロントページの続き

(72)発明者 フランクリン, ポール・デイヴィッド
アメリカ合衆国・98109-5210・ワシントン州・シアトル・テリー アヴェニュー ノース
・410

審査官 笠田 和宏

(56)参考文献 特開平06-332782(JP,A)
特開2003-029933(JP,A)
特表2012-531644(JP,A)
特開2006-268740(JP,A)
特開2009-245089(JP,A)
特開2009-187141(JP,A)
米国特許出願公開第2012/0271795(US,A1)
米国特許出願公開第2013/0138607(US,A1)
米国特許出願公開第2011/0029840(US,A1)

(58)調査した分野(Int.Cl., DB名)

IPC G06F 11/14
12/00
17/30