

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5918243号
(P5918243)

(45) 発行日 平成28年5月18日 (2016. 5. 18)

(24) 登録日 平成28年4月15日 (2016. 4. 15)

(51) Int. Cl. F I
G 0 6 F 12/00 (2006.01)
 G 0 6 F 12/00 5 3 1 R
 G 0 6 F 12/00 5 4 5 A

請求項の数 27 (全 28 頁)

(21) 出願番号	特願2013-530183 (P2013-530183)	(73) 特許権者	510071448
(86) (22) 出願日	平成23年9月13日 (2011. 9. 13)		ヒタチ データ システムズ コーポレー ション
(65) 公表番号	特表2013-544386 (P2013-544386A)		HITACHI DATA SYSTEM S CORPORATION
(43) 公表日	平成25年12月12日 (2013. 12. 12)		アメリカ合衆国 カリフォルニア州 95 050 サンタ・クララ ラファイエット ストリート 2845
(86) 国際出願番号	PCT/US2011/051313		
(87) 国際公開番号	W02012/039988	(74) 代理人	110000279
(87) 国際公開日	平成24年3月29日 (2012. 3. 29)		特許業務法人ウィルフォート国際特許事務 所
審査請求日	平成26年9月11日 (2014. 9. 11)	(72) 発明者	ブライアント, アラン ジー. アメリカ合衆国 02032 マサチュセ ッツ州 イーストウォルポール, ユニオン ストリート, 228
(31) 優先権主張番号	12/889, 744		
(32) 優先日	平成22年9月24日 (2010. 9. 24)		
(33) 優先権主張国	米国 (US)		

最終頁に続く

(54) 【発明の名称】 分散型データベースにおいてインテグリティを管理するためのシステム及び方法

(57) 【特許請求の範囲】

【請求項1】

(A) 1以上のメタデータを含んだ単位であるリージョンの正式なコピーと前記リージョンのバックアップコピーであり保護レベルと同数のバックアップコピーとで構成され複数のノードに格納されている複数のリージョンコピーのうちの前記正式なコピー又は前記バックアップコピーの損失が検出された場合、前記リージョンの不完全なコピーを生成し、

(B) ソースコピーから、前記不完全なコピーに、メタデータをコピーし、

(C) 前記不完全なコピーから、前記ソースコピーからコピーされたメタデータのうち、保留の更新でありターゲットコピーに無いメタデータを、前記ターゲットコピーにコピーし、

(D) (B)の完了前に、前記ターゲットコピーの内容が前記ソースコピーの内容と同じになった場合、前記ターゲットコピーを、前記リージョンのバックアップコピーとし、

(E) (C)の完了前に、前記不完全なコピーの内容が前記ソースコピーの内容と同じになった場合、前記不完全なコピーをバックアップコピーとする、

ことをコンピュータに実行させ、

前記損失が、いずれかのバックアップコピーの消失の場合、前記ソースコピーは、前記正式なコピーであり、前記ターゲットコピーは、前記消失の後に復元したバックアップコピーである、

ことを特徴とするコンピュータプログラム。

【請求項 2】

前記消失の後に復元したバックアップコピーを、(B)の完了まで、部分的なコピーとし、

(D)において、前記部分的なコピーを、前記リージョンのバックアップコピーとし、且つ、前記不完全なコピーを削除する、
ことをコンピュータに実行させることを特徴とする請求項1記載のコンピュータプログラム。

【請求項 3】

前記損失が、前記正式なコピーのダウンの場合、(A)の実行前に、前記1以上のバックアップコピーのうちのいずれかのバックアップコピーを、前記リージョンの正式なコピーに昇格する、

ことを更にコンピュータに実行させ、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、前記昇格された正式なコピーであり、前記ターゲットコピーは、前記ダウンした正式なコピーである、
ことを特徴とする請求項1又は2記載のコンピュータプログラム。

【請求項 4】

前記損失が、前記正式なコピーのダウンの場合、前記ダウンした正式なコピーを部分的なコピーとする、

ことを更にコンピュータに実行させ、

前記損失が、前記正式なコピーのダウンの場合、前記ターゲットコピーは、前記部分的なコピーであり、

(D)において、前記部分的なコピーを、前記リージョンのバックアップコピーとし、且つ、前記不完全なコピーを削除する、

ことを特徴とする請求項3記載のコンピュータプログラム。

【請求項 5】

前記複数のノードにわたる前記正式なコピー、前記バックアップコピー及び前記不完全なコピーの位置を特定するリージョンマップを提供する、

ことをコンピュータに実行させることを特徴とする請求項1乃至4のうちのいずれか1項に記載のコンピュータプログラム。

【請求項 6】

(A)1以上のメタデータを含んだ単位であるリージョンの正式なコピーと前記リージョンのバックアップコピーであり保護レベルと同数のバックアップコピーとで構成され複数のノードに格納されている複数のリージョンコピーのうちの前記正式なコピー又は前記バックアップコピーの損失が検出された場合、前記リージョンの不完全なコピーを生成し、

(B)ソースコピーから、前記不完全なコピーに、メタデータをコピーし、
(C)前記不完全なコピーから、前記ソースコピーからコピーされたメタデータのうち、

保留の更新でありターゲットコピーに無いメタデータを、前記ターゲットコピーにコピーし、

(D)(B)の完了前に、前記ターゲットコピーの内容が前記ソースコピーの内容と同じになった場合、前記ターゲットコピーを、前記リージョンのバックアップコピーとし、

(E)(C)の完了前に、前記不完全なコピーの内容が前記ソースコピーの内容と同じになった場合、前記不完全なコピーをバックアップコピーとする、

ことをコンピュータに実行させ、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、いずれかのバックアップコピーから昇格された正式なコピーであり、前記ターゲットコピーは、前記ダウンした正式なコピーである、

ことを特徴とするコンピュータプログラム。

【請求項 7】

前記損失が、前記正式なコピーのダウンの場合、(A)の実行前に、前記1以上のバック

10

20

30

40

50

クアップコピーのうちのいずれかのバックアップコピーを、前記リージョンの正式なコピーに昇格する、
ことを更にコンピュータに実行させ、
前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、前記昇格された正式なコピーである、
ことを特徴とする請求項 6 記載のコンピュータプログラム。

【請求項 8】

前記損失が、前記正式なコピーのダウンの場合、前記ダウンした正式なコピーを部分的なコピーとする、
ことを更にコンピュータに実行させ、
前記損失が、前記正式なコピーのダウンの場合、前記ターゲットコピーは、前記部分的なコピーであり、
(D)において、前記部分的なコピーを、前記リージョンのバックアップコピーとし、
且つ、前記不完全なコピーを削除する、
ことを特徴とする請求項 7 記載のコンピュータプログラム。

10

【請求項 9】

前記複数のノードにわたる前記正式なコピー、前記バックアップコピー及び前記不完全なコピーの位置を特定するリージョンマップを提供する、
ことをコンピュータに実行させることを特徴とする請求項 6 乃至 8 のうちのいずれか 1 項に記載のコンピュータプログラム。

20

【請求項 10】

1 以上のメタデータオブジェクトを含んだ単位であるリージョンの正式なコピーと前記リージョンのバックアップコピーであり保護レベルと同数のバックアップコピーとで構成された複数のリージョンコピーを格納する複数のノードを含んだシステムであって、
前記複数のリージョンコピーのうちの前記正式なコピー又は前記バックアップコピーの損失が検出された場合、前記リージョンの不完全なコピーを生成する生成手段と、
ソースコピーから、前記不完全なコピーに、メタデータをコピーする第 1 処理を実行する第 1 実行手段と、
前記不完全なコピーから、前記ソースコピーからコピーされたメタデータのうち、保留の更新であり前記ターゲットコピーに無いメタデータを、ターゲットコピーにコピーする第 2 処理を実行する第 2 実行手段と、
前記第 1 処理の完了前に、前記ターゲットコピーの内容が前記ソースコピーの内容と同じになった場合、前記ターゲットコピーを、前記リージョンのバックアップコピーとする第 1 変更手段と、
前記第 2 処理の完了前に、前記不完全なコピーの内容が前記ソースコピーの内容と同じになった場合、前記不完全なコピーをバックアップコピーとする第 2 変更手段と
を備え、
前記損失が、いずれかのバックアップコピーの消失の場合、前記ソースコピーは、前記正式なコピーであり、前記ターゲットコピーは、前記消失の後に復元したバックアップコピーである、
ことを特徴とするシステム。

30

40

【請求項 11】

前記消失の後に復元したバックアップコピーを、前記第 1 処理の完了まで、部分的なコピーとし、
前記第 1 変更手段は、前記部分的なコピーを、前記リージョンのバックアップコピーとし、且つ、前記不完全なコピーを削除する、
ことを特徴とする請求項 10 記載のシステム。

【請求項 12】

前記損失が、前記正式なコピーのダウンの場合、前記不完全なコピーの生成前に、前記 1 以上のバックアップコピーのうちのいずれかのバックアップコピーを、前記リージョン

50

の正式なコピーに昇格する手段

を更に備え、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、前記昇格された正式なコピーであり、前記ターゲットコピーは、前記ダウンした正式なコピーである、
ことを特徴とする請求項 10 又は 11 記載のシステム。

【請求項 13】

前記損失が、前記正式なコピーのダウンの場合、前記ダウンした正式なコピーを部分的なコピーとする手段

を更に備え、

前記損失が、前記正式なコピーのダウンの場合、前記ターゲットコピーは、前記部分的なコピーであり、

前記第 1 変更手段は、前記部分的なコピーを、前記リージョンのバックアップコピーとし、且つ、前記不完全なコピーを削除する、

ことを特徴とする請求項 12 記載のシステム。

【請求項 14】

前記正式なコピー、前記バックアップコピー及び前記不完全なコピーの位置を特定するリージョンマップが、前記複数のノードの各々に格納される、

ことを特徴とする請求項 10 乃至 13 のうちのいずれか 1 項に記載のシステム。

【請求項 15】

1 以上のメタデータオブジェクトを含んだ単位であるリージョンの正式なコピーと前記リージョンのバックアップコピーであり保護レベルと同数のバックアップコピーとで構成された複数のリージョンコピーを格納する複数のノードを含んだシステムであって、

前記複数のリージョンコピーのうちの前記正式なコピー又は前記バックアップコピーの損失が検出された場合、前記リージョンの不完全なコピーを生成する生成手段と、

ソースコピーから、前記不完全なコピーに、メタデータをコピーする第 1 処理を実行する第 1 実行手段と、

前記不完全なコピーから、前記ソースコピーからコピーされたメタデータのうち、保留の更新であり前記ターゲットコピーに無いメタデータを、ターゲットコピーにコピーする第 2 処理を実行する第 2 実行手段と、

前記第 1 処理の完了前に、前記ターゲットコピーの内容が前記ソースコピーの内容と同じになった場合、前記ターゲットコピーを、前記リージョンのバックアップコピーとする第 1 変更手段と、

前記第 2 処理の完了前に、前記不完全なコピーの内容が前記ソースコピーの内容と同じになった場合、前記不完全なコピーをバックアップコピーとする第 2 変更手段と
を備え、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、いずれかのバックアップコピーから昇格された正式なコピーであり、前記ターゲットコピーは、前記ダウンした正式なコピーである、

ことを特徴とするシステム。

【請求項 16】

前記損失が、前記正式なコピーのダウンの場合、前記不完全なコピーの生成前に、前記 1 以上のバックアップコピーのうちいずれかのバックアップコピーを、前記リージョンの正式なコピーに昇格する手段

を更に備え、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、前記昇格された正式なコピーである、

ことを特徴とする請求項 15 記載のシステム。

【請求項 17】

前記損失が、前記正式なコピーのダウンの場合、前記ダウンした正式なコピーを部分的なコピーとする手段

を更に備え、

前記損失が、前記正式なコピーのダウンの場合、前記ターゲットコピーは、前記部分的なコピーであり、

前記第1変更手段は、前記部分的なコピーを、前記リージョンのバックアップコピーとし、且つ、前記不完全なコピーを削除する、
ことを特徴とする請求項16記載のシステム。

【請求項18】

前記正式なコピー、前記バックアップコピー及び前記不完全なコピーの位置を特定するリージョンマップが、前記複数のノードの各々に格納される、
ことを特徴とする請求項15乃至17のうちのいずれか1項に記載のシステム。

10

【請求項19】

1以上のメタデータオブジェクトを含んだ単位であるリージョンの正式なコピーと前記リージョンのバックアップコピーであり保護レベルと同数のバックアップコピーとで構成された複数のリージョンコピーを格納する複数のノードを含んだシステムの制御方法であって、

(A) 複数のリージョンコピーのうちの前記正式なコピー又は前記バックアップコピーの損失が検出された場合、前記リージョンの不完全なコピーを生成し、

(B) ソースコピーから、前記不完全なコピーに、メタデータをコピーし、

(C) 前記不完全なコピーから、前記ソースコピーからコピーされたメタデータのうち、保留の更新であり前記ターゲットコピーに無いメタデータを、ターゲットコピーにコピーし、

20

(D) (B)の完了前に、前記ターゲットコピーの内容が前記ソースコピーの内容と同じになった場合、前記ターゲットコピーを、前記リージョンのバックアップコピーとし、

(E) (C)の完了前に、前記不完全なコピーの内容が前記ソースコピーの内容と同じになった場合、前記不完全なコピーをバックアップコピーとし、

前記損失が、いずれかのバックアップコピーの消失の場合、前記ソースコピーは、前記正式なコピーであり、前記ターゲットコピーは、前記消失の後に復元したバックアップコピーである、

ことを特徴とする制御方法。

【請求項20】

30

前記消失の後に復元したバックアップコピーを、(B)の完了まで、部分的なコピーとし、

(D)において、前記部分的なコピーを、前記リージョンのバックアップコピーとし、且つ、前記不完全なコピーを削除する、

ことを特徴とする請求項19記載の制御方法。

【請求項21】

前記損失が、前記正式なコピーのダウンの場合、(A)の実行前に、前記1以上のバックアップコピーのうちいずれかのバックアップコピーを、前記リージョンの正式なコピーに昇格する、

ことを更に実行し、

40

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、前記昇格された正式なコピーであり、前記ターゲットコピーは、前記ダウンした正式なコピーである、

ことを特徴とする請求項19又は20記載の制御方法。

【請求項22】

前記損失が、前記正式なコピーのダウンの場合、前記ダウンした正式なコピーを部分的なコピーとする、

ことを更に実行し、

前記損失が、前記正式なコピーのダウンの場合、前記ターゲットコピーは、前記部分的なコピーであり、

(D)において、前記部分的なコピーを、前記リージョンのバックアップコピーとし、

50

且つ、前記不完全なコピーを削除する、
ことを特徴とする請求項 2 1 記載の制御方法。

【請求項 2 3】

前記複数のノードにわたる前記正式なコピー、前記バックアップコピー及び前記不完全なコピーの位置を特定するリージョンマップを提供する、
ことを特徴とする請求項 1 9 乃至 2 2 のうちのいずれか 1 項に記載の制御方法。

【請求項 2 4】

1 以上のメタデータオブジェクトを含んだ単位であるリージョンの正式なコピーと前記リージョンのバックアップコピーであり保護レベルと同数のバックアップコピーとで構成された複数のリージョンコピーを格納する複数のノードを含んだシステムの制御方法であって、

(A) 複数のリージョンコピーのうちの前記正式なコピー又は前記バックアップコピーの損失が検出された場合、前記リージョンの不完全なコピーを生成し、

(B) ソースコピーから、前記不完全なコピーに、メタデータをコピーし、

(C) 前記不完全なコピーから、前記ソースコピーからコピーされたメタデータのうち、保留の更新であり前記ターゲットコピーに無いメタデータを、ターゲットコピーにコピーし、

(D) (B) の完了前に、前記ターゲットコピーの内容が前記ソースコピーの内容と同じになった場合、前記ターゲットコピーを、前記リージョンのバックアップコピーとし、

(E) (C) の完了前に、前記不完全なコピーの内容が前記ソースコピーの内容と同じになった場合、前記不完全なコピーをバックアップコピーとし、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、いずれかのバックアップコピーから昇格された正式なコピーであり、前記ターゲットコピーは、前記ダウンした正式なコピーである、

ことを特徴とする制御方法。

【請求項 2 5】

前記損失が、前記正式なコピーのダウンの場合、(A) の実行前に、前記 1 以上のバックアップコピーのうちいずれかのバックアップコピーを、前記リージョンの正式なコピーに昇格する、

ことを更に実行し、

前記損失が、前記正式なコピーのダウンの場合、前記ソースコピーは、前記昇格された正式なコピーである、

ことを特徴とする請求項 2 4 記載の制御方法。

【請求項 2 6】

前記損失が、前記正式なコピーのダウンの場合、前記ダウンした正式なコピーを部分的なコピーとする、

ことを更に実行し、

前記損失が、前記正式なコピーのダウンの場合、前記ターゲットコピーは、前記部分的なコピーであり、

(D) において、前記部分的なコピーを、前記リージョンのバックアップコピーとし、
且つ、前記不完全なコピーを削除する、

ことを特徴とする請求項 2 5 記載の制御方法。

【請求項 2 7】

前記複数のノードにわたる前記正式なコピー、前記バックアップコピー及び前記不完全なコピーの位置を特定するリージョンマップを提供する、
ことを特徴とする請求項 2 4 乃至 2 6 のうちのいずれか 1 項に記載の制御方法。

【発明の詳細な説明】

【関連出願の相互参照】

【0 0 0 1】

この出願は、また、次の出願と関連している：

10

20

30

40

50

整理番号12/889,762、「部分的なデータベース停電中に分散オブジェクトストレージシステムの有効性を増強するためのシステム及び方法」というタイトルで2010年9月24日に出願された。

【背景技術】

【0002】

本発明は、一般に、分散コンピューターネットワークにおける可用性、信頼性、及び持続性の高いデータストレージのための技術に関する。

【関連技術の説明】

【0003】

従来のテープ及び光学ストレージ解決法に取って代わるかまたはそれを補う、可用性、信頼性、及び持続性の高い方式による「固定コンテンツ」のアーカイブストレージ (archival storage) に対する必要性が高まってきた。用語「固定コンテンツ」とは、通常、参照または他の目的のためにそのまま保存されることが予想される任意のタイプのデジタル情報を指す。このような固定コンテンツの例としては、多くの例の中でもとりわけ、電子メール、文書、診断画像、チェック画像、音声録音、フィルム及びビデオなどが挙げられる。従来の独立ノード冗長アレイ (RAIN) ストレージ手法は、このような固定コンテンツの情報資産をストレージするための大規模なオンラインアーカイブを生み出すために選択されるアーキテクチャとして登場した。RAINアーキテクチャは、ノードが必要に応じてクラスタに結合すること及びクラスタから退出することを可能にすることにより、ストレージクラスタを1つ以上のノードの障害から隔離する。RAINタイプのアーカイブは、データを複数のノード上で複製することにより、ノードの障害または除去を自動的に補償することができる。通常、RAINシステムは、閉システム内の同一コンポーネントから設計されたハードウェア機器として広く提供される。

従来知られているアーカイブストレージシステムは、典型的に、ファイル毎にメタデータをそのコンテンツと同様に格納する。メタデータは、データを説明する、そのデータのコンポーネントである。メタデータは、典型的に、内容、質、条件、及びシステムに保存されている実際のデータのその他の特性について説明する。分散ストレージのコンテキストにおいて、ファイルに関するメタデータは、例えば、ファイルの名前、どこにファイルの断片が格納されているか、ファイルの作成日及び保持値を含む。信頼できるファイルストレージは、ファイルのストレージシステム信頼性及び有効性を達成するのに必要だが、メタデータのインテグリティもまたシステムの重要な部分である。しかし、先行技術では、潜在的に信頼性の低いノードの分散システムにわたってメタデータを分散することが可能ではなかった。本発明は、当技術分野でこの要求に対処する。

【0004】

改善されたアーカイバルストレージシステムは、一般に所有された米国特許第7,155,466号明細書、第7,657,581号明細書及び第7,657,586号明細書に記載されている。このシステムは、ノードの分散されたセットにわたって分散オブジェクトストアを提供する。米国特許第7,657,581号明細書によれば、対称なノードのアーカイバルストレージクラスタは、できればメタデータオブジェクトのかたちでメタデータへのアクセスを組織し提供する「メタデータ管理」システムを含む。メタデータオブジェクトは、それぞれユニークな名前を持っていて、メタデータオブジェクトは、リージョンへ組織される。1つの実施例において、リージョンは、1つ以上のオブジェクト属性 (例えばオブジェクトの名前) のハッシュ値を計算し、生じるハッシュ値のビットの所与の数を抽出することにより選択される。ビットの数は、配置パラメータによってコントロールされてもよい。このスキームでは、リージョンはそれぞれ重複して格納され、リージョンは、リージョンコピーのセットを含む。具体的には、リージョンの1つの正式なコピー及び0以上のバックアップコピーがある。記述されているように、コピーの数は、多くのメタデータ保護レベル (「MDPL」) と呼ばれることがある配置パラメータによってコントロールされてもよい。従って、例えばこのスキームの1つの実施例の中で、リージョンは正式なリージョンコピー及びそのMDPL-1バックアップコピーを含む。リー

10

20

30

40

50

ジョンコピーは、1つのノード当たりの合計のリージョンコピーの数と同様に1つのノード当たりの正式なリージョンコピーの数が均等になるようクラスタのノードにわたって分散される。

【0005】

上記メタデータマネージャシステムの別の態様は、各リージョンの各コピーの原因であるノードを識別するリージョン「マップ」と呼ばれる。リージョンマップは、メタデータ管理システムを有するプロセスによってアクセス可能である。リージョンマップ中のリージョンは、ハッシュ値のセットを表し、すべてのリージョンのセットはあらゆるハッシュ値をカバーする。リージョンは、ハッシュ値のビットの数を抽出することにより引き出される数によって識別される。ネームスペース分割スキームは、リージョンマップ中でリージョンを定義し、かつ所与のリージョンのオーナー権をコントロールするために使用される。この分割スキームは、データベース中で実装される。スキームでは、リージョンコピーは、3つの段階：「正式」、「バックアップ」そして「不完全」、のうちの1つをもつ。リージョンコピーが正式の場合、リージョンへのすべての要求はこのコピーに行き、各リージョンにつき1つの正式なコピーがある。リージョンコピーがバックアップ（あるいは不完全）である場合、コピーは更新要求（正式なリージョンマネージャプロセスからの）を受信する。メタデータがロードされているが、コピーがまだ同期されない（典型的に正式なリージョンコピーに関して）場合、リージョンコピーは不完全である。同期が完了するまで、不完全なリージョンコピーは、別の段階への昇格の資格を有さない、すなわち、その段階ではコピーはバックアップコピーになる。

【0006】

上記メタデータ管理スキームの別の態様は、バックアップリージョンコピーが正式なリージョンコピーと同期した状態で保たれることである。同期は、更新要求が処理されている場合に、プロトコルあるいは正式なリージョンコピーとそのMDPL-1バックアップコピーの間の「コントラクト」を強化することにより保証される。例えば、更新を局所的に委ねた後、正式なリージョンマネージャプロセスは、そのMDPL-1バックアップコピー（それらは典型的に、他のノードに置かれる）の各々への更新要求を出す。この通常の経過の中で、更新要求の受取において、所与のバックアップコピーに関連したリージョンマネージャプロセスは、確認を出すあるいは出すことを試みる。正式なリージョンマネージャプロセスは、更新が成功したという表示を提供する前にMDPL-1バックアップコピーのすべてからの確認を待つ。この更新処理が失敗することがあるいくつかの状況がある、例えば、正式なリージョンマネージャ（確認を待っている間）は、バックアップマネージャプロセスがしたことを示す例外に遭遇するかもしれない、あるいは、たとえそれが確認を出したとしても、バックアップマネージャプロセスは、更新要求を局所的に処理することを失敗するかもしれない、あるいは、確認を出す間のバックアップリージョンマネージャプロセスは、正式なリージョンマネージャプロセスが故障したことを示す例外に遭遇するかもしれない、等等。バックアップリージョンマネージャが更新を処理することができなければ、それはサービスから退出する。バックアップリージョンマネージャプロセスあるいは正式なマネージャプロセスのいずれかが故障の場合、新しいリージョンマップが出される。このように同期を保証することによって、バックアップコピーはそれぞれ正式なコピー用の「ホットスタンバイ」である。そのようなバックアップコピーは、正式なリージョンコピーが失われる場合あるいは現在の正式なリージョンコピーが降格されるべきである（またいくつかのバックアップリージョンコピーが昇格されるべきである）と負荷分散要求が指示するために必要になるかもしれない正式なコピーになることへの昇格の資格を有する。

【0007】

上記設計は、多くの同時のノード故障であってもそれがメタデータの高い有効性を保証するという点で有利である。ノード故障がある場合、1つ以上のリージョンが失われる、すなわち、その後、システムはそれをリカバリする必要がある。リカバリ過程は失われたリージョンのためのリージョンコピーを作成することを含む。その修理は、ハードウェア

10

20

30

40

50

とデータベースのリソースを消費し、従って、性能コストがある。大きなクラスタにおいては、修理が長時間かかる場合があり、その時間の間、クラスタはM D P L未満かもしれない。

【簡易概要】

【0008】

クラスタリカバリタイムは、ここに記載される増分リフレッシュ技術によって短縮される。その技術の目標は、以来データベースのその部分に生じる更新だけの増分リフレッシュを行なうことにより失われた（例えば障害中に）冗長分散データベースの部分を復旧させることである。

【0009】

1つの実施例では、増分リフレッシュは、ネットワーク化された独立ノードの冗長アレイにおいて実装され、メタデータオブジェクトは、アレイにわたって分散されたリージョンのセットに格納される。リージョンマップは、リージョンのセットの位置を識別し、マップは、いくつかの段階：正式な（A）、バックアップ（B）、そして不完全（I）、のうちの1つにあるリージョンのコピーと、どこに不完全なリージョンコピーがリージョンへの保留の更新を格納するかとを、典型的に含む。その技術によれば、その段階は、更に、保留の更新から追加される部分的な段階を含んでいて、クラスタに復帰するリージョンを修理するための復元の間、保留の更新が使用される。具体的には、その修理は、次のステップを含んでいる。最初に、リージョンの部分的な（P）コピーが、好ましくは（A）または（B）リージョンが期限切れであることを決定する場合に、（P）リージョンへ（A）または（B）リージョンコピーを降格すること（あるいは変換すること）により作成される。その後、部分的な（P）リージョンコピーは、適切な不完全な（I）リージョン上で既に待ち行列に入れられた保留の更新の適用により最新にされる。その後、（P）リージョンは、バックアップ（B）リージョンコピーに変換される。それが変換されるとすぐに、それに更新を送ることにより（P）リージョンが、バックアップリージョンとして扱われる。ある状況では、部分的な（P）コピーは、その結果、最後の更新を取消すためにそのリージョンの正式な（A）コピーからのその部分を同期させることにより最新にされる。

【0010】

代替実施例では、上記方法は、プロセッサと、プロセッサによって実行されたときに上記方法を実行するコンピュータプログラム指示を保持するコンピュータメモリとを有する装置によって実行される。

【0011】

別の代替実施例では、上記方法は、それ自体がバッキングストアに置かれているかもしれないバージョンファイルシステムに関連する使用のためのコンピュータ読取り可能な媒体中のコンピュータプログラム製品によって実行される。そのコンピュータプログラム製品は、プロセッサによって実行されたとき上記方法を実行するコンピュータプログラム指示を保持する。

前述のものは、発明のより適切な特徴のいくつかを概説した。これらの特徴は単に例となるために解釈されるべきである。異なるやり方で記載された発明を適用すること、あるいは今後記載される発明の変更により、他の多くの有益な成果が成し遂げられるかもしれない。

【図面の簡単な説明】

【0012】

本発明及びその利点のより完全な理解については、参考文献が、添付図面と併用される次の記載にある。

【図1】ここでの本件をその中で実施できる固定コンテンツストレージャークライプの簡略化されたブロック図である。

【図2】各独立ノードが対称でありアーカイブクラスタアプリケーションをサポートする独立ノード冗長アレイの簡略化された説明図である。

10

20

30

40

50

【図3】所与のノード上で実行されるアーカイブクラスタアプリケーションの様々なコンポーネントの高レベル説明図である。

【図4】クラスタの所与のノード上のメタデータ管理システムのコンポーネントを示す図である。

【図5】例示的なリージョンマップを示す図である。

【図6】クラスタがサイズを増すとともにリージョンマップ変更を容易にするためにネームスペース分割スキームがどのように使用されるかを示す図である。

【図7】(B)リージョンの損失の次に行われるリストア中に、どのように(P)リージョンコピー(Pリージョン)が作成され、そして、ここで記載される技術によって(B)リージョンへ昇格されるかの例を示す図である。

10

【図8】(A)リージョンの損失の次に行われるリストア中に、どのように(P)リージョンコピー(Pリージョン)が作成され、そして、ここで記載される技術によって(B)リージョンへ昇格されるかの例を示す図である。

【例示的な実施例の詳細な記載】

【0013】

以下に記載する技術は、スケーラブルなディスクベースのアーカイブストレージ管理システム、好ましくは、独立ノード冗長アレイに基づくシステムアーキテクチャの中で好ましくは実施される。各ノードは、異なるハードウェアを備える場合があり、したがって「異種である」と考えることができる。ノードは通常、1つまたは複数のストレージディスクへのアクセスを有するが、これらのストレージディスクは、実際の物理ストレージディスクとすることもでき、またはストレージエリアネットワーク(SAN)におけるように仮想ストレージディスクとすることもできる。各ノード上でサポートされるアーカイブクラスタアプリケーション(及び任意選択で、このアプリケーションが実行される、基礎を成すオペレーティングシステム)は、同じかまたはほぼ同じとすることができる。例示的な一実施形態では、各ノード上のソフトウェアスタック(オペレーティングシステムを含むことができる)は対称だが、ハードウェアは異種である場合がある。このシステムを使用して、図1に示すように、企業が、とりわけ文書、電子メール、衛星画像、診断画像、チェック画像、音声録音、ビデオなど、多くの異なるタイプの固定コンテンツ情報のための永続的ストレージを生み出すことができる。当然、これらのタイプは例示に過ぎない。独立サーバまたはいわゆるストレージノード上でデータを複製することによって、高レベルの信頼性が達成される。各ノードは、そのピアと対称であることが好ましい。したがって、好ましくはどんな所与のノードも全ての機能を実施することができるので、いずれか1つのノードの障害はアーカイブの可用性にほとんど影響を及ぼさない。

20

30

【0014】

一般に所有される米国特許第7,155,466号明細書、第7,657,581号明細書、及び第7,657,586号明細書で述べられているように、各ノード上で実行される分散ソフトウェアアプリケーションが、デジタル資産を取り込み、保存し、管理し、取り出す。図2の例示する一実施形態では、個別アーカイブの物理的境界が、クラスタと呼ばれる。通常、クラスタは、単一のデバイスではなく、デバイスの集合である。デバイスは、同種または異種である場合がある。典型的なデバイスは、Linuxなどのオペレーティングシステムを実行するコンピュータまたはマシンである。コモディティハードウェア上でホストされるLinuxベースのシステムのクラスタは、少数のストレージノードサーバから、何千テラバイトものデータを記憶する多くのノードまでスケールできるアーカイブを提供する。このアーキテクチャにより、記憶容量が、組織のますます増加するアーカイブ要件に常に遅れを取らずにいられることが確実になる。アーカイブが常にデバイス障害から保護されるように、データがクラスタ全体で複製されることが好ましい。ディスクまたはノードに障害が発生した場合、クラスタは自動的に、クラスタ中の、同じデータのレプリカを維持する他のノードにフェイルオーバーする。

40

【0015】

例示的なクラスタは、以下の一般的なコンポーネントカテゴリ、すなわちノード202

50

、ネットワークスイッチ204のペア、電力分散ユニット(PDU)206、及び無停電電源(UPS)208を好ましくは有する。ノード202は、典型的に1つ以上のコモディティサーバを有し、CPU(例えばインテルx86、適切なランダムアクセスメモリ(RAM)、1つ以上のハードドライブ(例えば標準のIDE/SATA、SCSIなど)及び2枚以上のネットワークインターフェイス(NIC)カード)を有する。典型的なノードは、2.4GHzのチップ、512MB RAM及び6つの(6)200GBハードドライブを備えた2Uラックマウントユニットである。しかし、これに限られない。ネットワークスイッチ204は、典型的にノード間のピアツーピア通信を可能にする内部スイッチ205、及び、各ノードへの追加のクラスタアクセスを許可する外部スイッチ207を有する。各スイッチは、クラスタ中の全て潜在的なノードを扱うことを十分なポートに要求する。イーサネットまたはGigEのスイッチは、この目的に使用されてもよい。PDU206は、全てのノード及びスイッチに動力を提供するために使用され、UPS208は、全てのノード及びスイッチを保護するために使用される。制限しているつもりではないが、典型的にクラスタは、公衆インターネット、企業イントラネットあるいは他の広域か、ローカルエリアネットワークのようなネットワークに連結可能である。例となる実施例において、クラスタは、企業環境内で実施される。それは、例えばサイトの企業ドメイン名前システム(DNS)ネームサーバによってナビゲートすることにより達するかもしれない。従って、例えば、クラスタのドメインは、既存のドメインの新しいサブドメインかもしれない。代表的な実施において、サブドメインは、企業のDNSサーバの中でクラスタ自体の中のネームサーバに委託される。エンドユーザは、任意の従来のインターフェースあるいはアクセスツールを使用して、クラスタにアクセスする。従って、例えば、クラスタへのアクセスは、任意のIPベースのプロトコル(HTTP、FTP、NFS、AFS、SMB、ウェブサービスなど)上に、APIを経由して、あるいは他の既知か、その後発展したアクセス方式、サービス、プログラムあるいはツールを通じて実行されるかもしれない。

10

20

【0016】

クライアントアプリケーションは、標準UNIXファイルプロトコルのような1つ以上のタイプの外部ゲートウェイによるクラスタ、またはHTTP APIにアクセスする。アーカイブは、好ましくは、オプションで任意の標準UNIXファイルプロトコル系設備の下に位置できる仮想ファイルシステムを通じて露出される。これらは、NFS、FTP、SMB/CIFSなどを含む。

30

【0017】

1つの実施例において、アーカイブクラスタアプリケーションは、クラスタとして(例えば、イーサネット経由で)ネットワーク化される独立ノード(H-RAIN)の冗長アレイ上で作動する。所与のノードのハードウェアは異種かもしれない。最大の信頼性のために、しかし、好ましくは、各ノードは、今、図3に示されるようないくつかのランタイムコンポーネントから成る分散アプリケーション(すなわち、同じインスタンス、あるいは本質的に同じインスタンスかもしれない)のインスタンス300を実行する。従って、ハードウェアは、異種かもしれないが、ノード(少なくともそれが本発明に係る)上のソフトウェアスタックは、同じである。これらのソフトウェアコンポーネントは、ゲートウェイプロトコルレイヤ302、アクセスレイヤ304、ファイルトランザクション及び管理レイヤ306、及び、コアコンポーネントレイヤ308を有する。機能が他の意味のある方法で特徴づけられるかもしれないことを通常のスキルのうちの1つが認識するので、「レイヤ」指定は、説明目的に提供される。レイヤ(あるいはその点でコンポーネント)の1つ以上は、統合されるかそうでないかもしれない。いくつかのコンポーネントは、レイヤに渡って共有されるかもしれない。

40

【0018】

ゲートウェイプロトコルレイヤ302の中のゲートウェイプロトコルは、既存のアプリケーションに透明性を提供する。具体的には、ゲートウェイは、カスタムアプリケーションを構築するためのウェブサービスAPI同様に、NFS310及びSMB/CIFS3

50

12のようなネイティブファイルサービスを提供する。HTTPサポート314も提供される。アクセスレイヤ304は、アーカイブへのアクセスを提供する。具体的には、発行によれば、固定コンテンツファイルシステム(FCSS)316は、アーカイブオブジェクトへフルアクセスを提供するためにネイティブファイルシステムをエミュレートする。あたかもそれらが通常のファイルかのように、FCFSは、アーカイブコンテンツにアプリケーションダイレクトアクセスを与える。好ましくは、メタデータがファイルとして露出されている一方、アーカイブコンテンツは、そのオリジナルフォーマットでレンダリングされる。FCFS316は、管理者らが、使いやすい方法で固定コンテンツのデータをセットアップできるように、ディレクトリと許可についての従来のビュー及びルーチンファイルレベルコールを提供する。ファイルアクセス呼び出しは、好ましくは、ユーザスペースデーモンによって傍受され、ダイナミックに呼び出しのアプリケーションへの適切な表示を作成する適切なコアコンポーネント(レイヤ308に)に送られる。FCFS呼び出しは、好ましくは、自律的なアーカイブ管理を促進するアーカイブポリシーによって抑制される。従って、一例において、管理者がアプリケーションは、保存期間(所与のポリシー)がまだ有効のアーカイブオブジェクトを削除することができない。

10

【0019】

アクセスレイヤ304は、好ましくは、また、ウェブユーザインターフェース(UI)318及びSNMPゲートウェイ320を含む。ウェブユーザインターフェース318は、ファイルトランザクション及び管理レイヤ306での管理エンジン322への対話型のアクセスを提供する管理者コンソールとして好ましくは実施される。管理上のコンソール318は、好ましくは、アーカイブオブジェクト及び個々のノードを含むアーカイブの動的考察を提供するパスワードで保護されウェブベースのGUIである。SNMPゲートウェイ320は、安全にストレージ管理アプリケーションがクラスターアクティビティを監視し制御することを可能にしながらストレージ管理アプリケーションに管理エンジン322への容易なアクセスを提供する。管理エンジンモニタは、システムとポリシーイベントを含むアクティビティをクラスタする。ファイルトランザクションと管理レイヤ306は、また、要求マネージャプロセス324を含む。要求マネージャ324は、コアコンポーネントレイヤ308の中のポリシーマネージャ326からの内部要求と同様に外界(アクセスレイヤ304を通じて)からの全ての要求を統合する。

20

【0020】

コアコンポーネントは、ポリシーマネージャ326に加えて、メタデータマネージャ328及びストレージマネージャ330の1つ以上のインスタンスを含む。メタデータマネージャ328は、各ノードに好ましくはインストールされる。クラスタ中のメタデータマネージャは、集散的に、全てのアーカイブオブジェクトを管理しながら、分散型データベースとして作動する。所与のノードにおいて、メタデータマネージャ328は、好ましくは、アーカイブオブジェクトのサブセットを管理し、各オブジェクトが、好ましくは、外部ファイル(「EF」(ストレージ用にアーカイブに入ったデータ))とアーカイブデータが物理的に検索される内部ファイル(各々「IF」)のセット間をマップする。同じメタデータマネージャ328は、また、他のノードから複製されたアーカイブオブジェクトのセットを管理する。従って、全ての外部ファイルの現状は、いくつかのノード上の複数メタデータマネージャに常に利用可能である。ノード障害の場合には、他のノード上のメタデータマネージャが、機能不全のノードによって以前管理されたデータへのアクセスを提供し続ける。このオペレーションは、以下により詳細に説明される。ストレージマネージャ330は、分散アプリケーション中の他の全てのコンポーネントに利用可能なファイルシステムレイヤを提供する。好ましくは、それはノードのローカルファイルシステムにデータオブジェクトを格納する。所与のノード中の各ドライブは、好ましくは、それぞれ自身のストレージマネージャを持っている。これは、ノードが個別のドライブを削除し、処理能力を最適化することを可能にする。ストレージマネージャ330は、また、システム情報、データの一貫性チェック及び直接ローカル構造をトラバースする能力を提供する。

30

40

【0021】

50

さらに図3で示されるように、クラスタは、通信ミドルウェアレイヤ332及びDNSマネージャ334を通じて内部及び外部通信を管理する。インフラストラクチャ332は、アーカイブコンポーネント中の通信を可能にする効率的で信頼できるメッセージベースのミドルウェアレイヤである。図で示した実施例において、レイヤは、マルチキャストとポイントツーポイント通信をサポートする。DNSマネージャ334は、企業サーバに全てのノードを接続する分散型ネームサービスを行う。好ましくは、DNSマネージャ(単独であるいはDNSサービスと共に)ロードバランスは、最大のクラスタ処理能力及び有効性を保証することを全てのノードに渡って要求する。

【0022】

図で示した実施例において、ARCアプリケーションインスタンスは、Red Hat Linux 9.0、Fedora Core 6などのような基礎オペレーティングシステム336上で実行する。通信ミドルウェアは、任意の便利な分散型通信メカニズムである。他のコンポーネントは、固定コンテンツファイルシステム(FCSS)316のために使用されてもよいFUSE(USErspaceの中のファイルシステム)を含むかもしれない。NFSゲートウェイ310は、標準のnfsd LinuxカーネルNFSドライバによって実施されるかもしれない。各ノード中のデータベースは、実施されるかもしれない、例えば、オブジェクト関係データベース管理システム(ORDBMS)であるPostgreSQL(またここにPostgresとして引用される)である。ノードは、Java HTTPサーバ及びservletコンテナであるジェティのようなウェブサーバを含むかもしれない。もちろん、上記のメカニズムは、単に例となる。

【0023】

所与のノード上のストレージマネージャ330は、物理的な記憶デバイスを管理する責任がある。好ましくは、各ストレージマネージャインスタンスは、全てのファイルがその配置アルゴリズムによって入れられる単一のルートディレクトリの責任がある。複数のストレージマネージャインスタンスは、ノード上で同時に作動することができ、各々は、通常、システムで異なる物理的なディスクを表わす。ストレージマネージャは、システムの残りから使用されているドライブ及びインターフェース技術を抽象する。ストレージマネージャインスタンスがファイルを書くように依頼される場合、それはそのために責任を負う表現用のフルパス及びファイル名を生成する。代表的な実施例において、ストレージマネージャ上に格納される各オブジェクトは、それが異なるタイプの情報を追跡するデータを格納する場合、ファイルにそれ自身のメタデータを加えて、そのときストレージマネージャと共に保存されるローデータとして受信される。例として、このメタデータは、限定無しで含む：EF長さ(バイトでの外部ファイルの長さ)、IFセグメントサイズ(内部ファイルのこの部分のサイズ)、EF保護表現(EF保護モード)、IF保護役割(この内部ファイルの表現)、EF生成タイムスタンプ(外部ファイルタイムスタンプ)、シグネチャ(シグネチャタイプを含む書き込み(PUT)の時間の内部ファイルのシグネチャ)及びEFファイル名(外部ファイルファイル名)。内部ファイルデータでこの追加のメタデータを格納することは、追加のレベルの保護を提供することである。具体的には、スカビンジンは、内部ファイルに保存されたメタデータからデータベースに外部ファイルレコードを作成することができる。他のポリシは、内部ファイルが元の状態のままになることを有効にするために内部ファイルに対する内部ファイルハッシュを有効にすることができる。

【0024】

前述のように、内部ファイルは、アーカイブオブジェクト中でオリジナルの「ファイル」の一部を表わす、データの「チャンク」かもしれない、また、それらはストライピングと保護ブロックを達成するために異なるノードに置かれるかもしれない。典型的に、1つの外部ファイルエントリは各アーカイブオブジェクトのためのメタデータマネージャの中にあり、その一方で、個々の外部ファイルエントリのための多くの内部ファイルエントリがあるかもしれない。典型的に、内部ファイルレイアウトは、システムに依存する。所与の実施において、ディスク上のこのデータの実際の物理フォーマットは一連の可変長レコードに格納される。

10

20

30

40

50

【 0 0 2 5 】

要求マネージャ 3 2 4 は、システム内の他のコンポーネントとのやりとりによりアーカイブアクションを行なうために必要とされるオペレーションのセットを実行する責任がある。要求マネージャは、異なるタイプの多くの同時のアクションをサポートし、機能不全のトランザクションをロールバックすることができ、実行するのに長い時間かかるトランザクションをサポートする。要求マネージャは、更に、アーカイブの読取り書き込みオペレーションが適切に扱われる事を保証し、全ての要求がいつでも既知の状態である事を保証する。更に、それは、所与のクライアント要求を満たすためにノードに渡って複数の読取り書き込みオペレーションを調整するためにトランザクション制御を提供する。更に、要求マネージャは、最近使われたファイルのためのメタデータマネージャエントリをキャッシュに格納し、データブロックと同様にセッションのためのバッファリングを提供する。

10

【 0 0 2 6 】

クラスタの主要な責任は、ディスク上に無制限のファイルを実際に格納することである。それが何らかの理由で手が届かないか、そうでなければ利用不可能かもしれないという意味で、所与のノードは「信頼性が低い」と見なされるかもしれない。そのような潜在的に信頼性の低いノードのコレクションは、確実に、高度に利用可能なストレージを作成することに協力する。一般に、格納される必要のある 2 つのタイプの情報がある：ファイル自体及びファイルに関するメタデータ。

【 0 0 2 7 】

メタデータ管理

20

【 0 0 2 8 】

米国特許第 7 , 6 5 7 , 5 8 1 号明細書（その開示は参考文献によってここに組込まれる）に記載のように、メタデータマネジメントシステムは、システムメタデータのような所与のメタデータへのアクセスを組織し提供する責任がある。このシステムメタデータは、構成情報、管理 UI に表示された情報、メトリクス、復元不能なポリシー違反についての情報などに表示された情報等と同様にアーカイブに置かれたファイルについての情報を含む。詳細に示されていないが、他のタイプのメタデータ（例えばアーカイブしたファイルに関連したユーザメタデータ）も今説明されるメタデータ管理システムを使用して管理されるかもしれない。

30

【 0 0 2 9 】

代表的な実施例では、メタデータ管理システムは、次のオブジェクトタイプ（それらは単に例となる）の 1 つ以上を含んでいるかもしれないメタデータオブジェクトのセットのための持続性を提供する。

- ・ ExternalFile : アーカイブのユーザによって知覚されるようなファイル、
- ・ InternalFile : ストレージマネージャによって格納されたファイル。典型的には、外部ファイルと内部ファイルの間に一対多数の関係があるかもしれない、
- ・ ConfigObject : クラスタを構成するのに使われる名前 / 値ペア、
- ・ AdminLogEntry : 管理者 UI に表示されるメッセージ、
- ・ MetricsObject : ある時点でのアーカイブ（例えばファイルの数）のある測定を表わす、タイムスタンプされたキー / 値ペア、そして
- ・ PolicyState : あるポリシーの違反。

40

【 0 0 3 0 】

各メタデータオブジェクトは、好ましくは、変わらないユニークな名前を持っているかもしれない。メタデータオブジェクトは、リージョンに組織される。リージョンは、正式なリージョンコピーとメタデータ保護レベル（MPDL）数（0 以上のセット）バックアップリージョンコピーを有する。0 のコピーで、メタデータマネジメントシステムは、計量可能であるが、高度に利用可能ではないかもしれない。リージョンは、1 つ以上のオブジェクト属性（例えばフルパス名やその一部のようなオブジェクトの名前）をハッシュ及びハッシュ値のビットの所与数の抽出により選択される。これらのビットは、リージョン

50

番号から成る。選択されたビットは、低位ビット、高位ビット、中位ビットあるいは個々のビットの任意のコンビネーションかもしれない。代表的な実施例において、所与のビットはハッシュ値の低位ビットである。オブジェクトの属性が属性（複数）は、任意の便利なハッシュ関数を使用してハッシュされるかもしれない。これらは制限なしで、java.lang.string.hashCode等のようなJavaベースのハッシュ関数を含む。好ましくは、リージョン番号から成るビットの数は、ここでregionMapLevelと呼ばれ、構成パラメータによってコントロールされる。この構成パラメータが6にセットされる場合、例えば、これは $2^6 = 64$ リージョンが得られる。もちろん、以下に説明されるように、多くのリージョンは許され、リージョンの数はネームスペース分割スキームを使用して自動的に調節されるかもしれない。

10

【0031】

米国特許第7,657,581号明細書に記載通り、各リージョンは、重複して格納されるかもしれない。上記の通り、リージョンの1つの正式なコピー及び0以上のバックアップコピーがある。前述のように、バックアップコピーの数は、メタデータ・データ保護レベル（あるいは「MDPL」）構成パラメータによってコントロールされる。好ましくは、リージョンコピーは、1つのノード当たりの正式なリージョンコピーの数の平衡を保ち、かつ1つのノード当たりの合計のリージョンコピーの数の平衡を保つようにクラスタの全てのノードに渡って分散される。

【0032】

メタデータ管理システムは、各ノード上で作動するデータベースにメタデータオブジェクトを格納する。このデータベースは、リージョンマップをサポートするために使用される。典型的なデータベースは、オープンソースとして利用可能であるPostgreSQLを使用して実施される。好ましくは、各リージョンコピーのスキーマがあり、各スキーマでは、各タイプのメタデータオブジェクト用のテーブルがある。スキーマは、単にテーブル、インデックス、手順及び他のデータベースオブジェクトを所有することができるネームスペースである。各リージョンは、好ましくは、それ自身のスキーマを持っている。各スキーマは、テーブル式、すなわち各メタデータオブジェクトに1つ持っている。これらのテーブルのうちの1つの列は、単一のメタデータオブジェクトに相当する。Postgresが好ましいデータベースであると同時に、任意の便利なリレーショナルデータベース（例えばオラクル、IBM DB/2など）が使用されてもよい。

20

30

【0033】

図4で示されるように、各ノード400は、プロセスあるいはコンポーネント、すなわち、1つ以上のリージョンマネージャ（RGM）402a-n、メタデータマネージャ（MM）404、少なくとも1つのメタデータマネージャクライアント（MMC）406、及び1つ以上のスキーマ410a-nがある1つのデータベース408、のセットを有する。RGM(s)、MM及びMMCコンポーネントは、Java仮想マシンのようなパーチャルマシン412で実行する。各リージョンコピーにつき1つのRGMがある。従って、正式なリージョンコピー用のRGM、各バックアップリージョンコピー用のRGM及びそれぞれ不完全なリージョンコピー用のRGMがある。RGM402のスキーマを管理する各RGM402用のデータベーススキーマ410もある。データベースは、また、リージョンマップ405を格納する。上述の特許の開示によれば、各ノードは、好ましくは、同期スキームによって強化されている要求と共に、リージョンマップの同じ全体的な見解を持っている。リージョンマネージャRGM402は、リージョンコピー（それが正式な、バックアップ、あるいは不完全な場合によっては）上で作動し、メタデータマネージャクライアント406、及び他のリージョンマネージャ402によって提出された要求の実行に責任がある。要求は、図3で示された通信ミドルウェアあるいは他のメッセージングレイヤのような任意の便利な手段を通じて所与のRGMに提供される。リージョンマネージャは、これらの要求が実行する実行環境を提供する、例えば、スキーマのRGMによって管理されているスキーマ上で作動するように構成されているデータベースへの接続を提供することによって。各リージョンマネージャは、データベース408にそのデータを格納す

40

50

る。メタデータマネージャ404は、ノード上のメタデータ管理の責任があるトップレベルのコンポーネントである。それは、リージョンマネージャ(RGM)を作成し破壊し、そして、RGM、例えばクラスタ構成情報、データベース接続のプールによって必要とされるリソースを組織する責任がある。好ましくは、所与のメタデータマネージャ(所与のノード中の)は、リーダーとして働き、どのメタデータマネージャ(ノードのセット又はサブセットに渡った)がどのリージョンコピーに責任を負うかを定める責任がある。賛成アルゴリズム又はその変形のようなリーダー選挙アルゴリズムは、メタデータマネージャリーダーを選ぶために使用されるかもしれない。好ましくは、1つのノード当たり複数のMMを実行することは可能であるが、各ノードは、1つのメタデータマネージャを持っている。一旦リージョンオーナー権がネームスペース分割スキーム(下記に述べられるように)によって確立されたならば、各メタデータマネージャは、1つ以上のリージョンマネージャのそのセットに従って調節することに責任がある。システムコンポーネント(例えば管理エンジン、ポリシマネージャなど)は、メタデータマネージャクライアントを通じてメタデータマネージャMMとやりとりをする。MMCは、所与の要求を実行するためにRGMを見つける事、選択されたRGMに要求を出す事、及び選択されたRGMが利用不可能な場合に(例えば、ノードが機能しなくなったので)要求を再試行することに責任がある。後者の場合は、新しいリージョンマップがノードで受信される場合、再試行要求が成功するであろう。

10

【0034】

上記の通り、リージョンマップは、各リージョンの各コピーに責任のあるノードを識別する。パーチャルマシン412(またその中での各RGM、MM、及びMMC構成要素)は、リージョンマップ405へのアクセスを持っている;リージョンマップのコピー420も、それがJVMにコピーされた後、図4に示される。リージョンマップは、従って、所与のノード中のJVM及びデータベースの両方に利用可能である。このインスタスとなる実施例において、各メタデータオブジェクトは、0x0と0x3fffffff合計間の整数を産出するためにハッシュされる、つまり30ビットの値の属性(例えば名前)を持っている。これらの値は、オーバーフロー問題(例えば範囲の高域に1を加える時)にぶち当たる事なく、符号付き32ビット整数中で快適に表わす事ができる。30ビットは、大きなクラスタにさえ十分であるおよそ10億までのリージョンを考慮に入れる。リージョンは、1セットのハッシュ値を表わし、全てのリージョンのセットは、あらゆるハッシュ値をカバーする。各リージョンのための異なるビット位置があり、異なるビット位置は、好ましくは固定順になっている。従って、各リージョンは、ハッシュ値のRegionLevelMapビットの抽出により好ましくは引き出される数によって識別される。64リージョンを考慮に入れて、構成パラメータが6にセットされる場合、生じるハッシュ値は、0x0から0x3fの数である。

20

30

先述の通り、リージョンコピーは、3つの(3)段階、すなわち、「正式な」(A)、「バックアップ」(B)そして「不完全」(I)のうちの1つにある。リージョンコピーが正式な場合、リージョンへの全ての要求がこのコピーに行き、また、各リージョンにつき1つの正式なコピーがある。リージョンコピーがバックアップである場合、コピーは、バックアップ要求(正式なリージョンマネージャプロセスからの)を受信する。メタデータがロードされているが、コピーがまだ同期されない(典型的に他のバックアップコピーに関して)場合、リージョンコピーは、不完全である。同期が完了するまで、不完全なリージョンコピーは、別の段階への昇進の資格を有さない、すなわち、そのポイントではコピーは、バックアップコピーになる。各リージョンは、1つの正式なコピー、所与の数(MDPL構成パラメータによってセットされた)バックアップあるいは不完全なコピーを持っている。

40

【0035】

米国特許第7,657,581号明細書に記載通り、バックアップリージョンコピーは、正式なリージョンコピーとそのMDPLバックアップコピー間で所与のプロトコル(あるいは「契約」)を強化することにより、正式なリージョンコピーと同期され続ける。こ

50

のプロトコルは、今説明される。

【 0 0 3 6 】

説明されたように、リージョンマップは、各リージョンの各コピーのオーナー権について説明する。例えば、図5は、4つのノードクラスタ用のリージョンマップをmetadataMD PL=2で例証する。この例において、示されるように、ノード1はリージョン0は正式であり、ノード2及び3はバックアップとして指定され、ノード2はリージョン1は正式であり、ノード3及び4はバックアップとして指定されるなどである。ネームスペース分割スキームは、クラスタが増大するとともに、特定のリージョンのコントロール（オーナー権）を変更するために使用されても良い。動的成長を許可する1つの方法は、ハッシュ値番号を有するツツの数を決定するregionMapLevel構成パラメータをインクリメントすること
10
である。クラスタが増大するとともに、リージョンマップの1つ以上のパーティションが「分離した」工程を経る。分離は、ハッシュ値のもう1つのビットを使用し、その結果メタデータを再分散することを引き起こす。例えば、レベル6のマップ、及びハッシュ値0x1000002a及び0x1000006aを備えた2つのメタデータオブジェクトを考慮してほしい。これらのハッシュ値（2進法の「0010」である「2」、及び2進法の「0110」である「6」を備えた16進法0x2a）の最後の6ビットは、同じである：したがって、両方のオブジェクトは、リージョン0x2aに分類される。その後、マップレベルが7に増加される場合、リージョンは0から0x7fであり、それにより、異なるリージョン、すなわち0x2a、0x6aに入ること
20
を2つのオブジェクトに強いる。

【 0 0 3 7 】

このアプローチは使用されてもよいが、それは同時に分離していることをすべてのリージョンに要求する。よりよい技術は、リージョンを増加的に分離することである。これをするために、ネームスペース分割スキームは、リージョン0でスタートし、現在のレベルの最後のリージョンで終了する順番でリージョンを分離する。リージョンは、ハッシュ値のもう1つのビットの使用により分離される。図6はこのプロセスを示す。この例において、マップレベル1では、2つのリージョン602（ノード0）及び604（ノード1）があると仮定する。ノード番号は、2進法で示される。マップが増大する必要がある場合、分割スキームは、ハッシュ値のもう1つのビットの使用によりリージョン0を分離する。これは、3つのリージョン606、608及び610を作る。それらがリージョン606（ノード00）にあり、残りのオブジェクトが新しい最後のリージョン610（ノード
30
10）へ行く場合、新しいビットが0であるオブジェクトはとどまる。分裂により加えられるビットはイタリック体にされる、すなわち00と10。注目すべきは、最初のリージョン606及び最後のリージョン610が2ビットを使用し、その一方で中間（分割されていない）リージョンは1つだけ使用する、けれども、左から右まで見られた場合、番号付けをするスキームはそれでも、正確に機能する、すなわち、{0, 1, 2}。更なる増大については、リージョン1は4つのリージョン612（ノード00）、614（ノード01）、616（ノード10）及び618（ノード11）を作るために分割される。これは2レベルを完成させる。リージョンマップが再び増大する必要がある場合、スキームは、リージョン00~000（つまりハッシュ値のもう1つのビットを加えることによって）を分割し、最後に、新しいリージョン100（さらにハッシュ値のもう1つのビットを加えること
40
によって）を加える。その後、リージョンマップは、示されるように5つのリージョン620、622、624、626及び628を持つことになる。

【 0 0 3 8 】

リージョンの数がノードの数に相当するという必要はない。より一般に、リージョンの数は、独立したノードのレイ中のノードの数に関連しない。

【 0 0 3 9 】

従って、1つの実施例によれば、リージョンに対するコントロールは、リージョンにメタデータオブジェクトを割り当て、次にリージョンを増加的に分けることで果たされる。リージョンコピー（正式、バックアップ、あるいは不完全である）は、各ノード上のデータベースに格納される。記述されたように、メタデータオペレーションは、正式な R G M
50

によって実行される。しかしながら、ノードが失敗する場合、いくつかの数のリージョンコピーは失われる。記述されたように、有効性は、正式なリージョンのバックアップコピーのうちの1つを昇格させることにより復元される、すなわちそれは数秒で通常できる。バックアップが昇格される短い間隔中に、MMCによってリージョンに提出される要求は失敗する。この障害は、遅れの後に再試行を引き起こすMMCによって見つかる例外として現われる。要求が再試行される時には、しかしながら、MMCユーザに対してサービスが中断されないことをもたらしながら最新のマップが適所にあるようになる。記述されたように、このアプローチは、同期されたままであるリージョンのコピー（好ましくはそれらのすべて）に依存する。

【0040】

下記は、メタデータ管理システムの追加導入詳細を提供する。

【0041】

上記の通り、ノードがクラスタを残す場合、あるいはノードがクラスタを連結する場合、あるいは不完全なリージョンコピーがロードを終える場合、MMリーダーはリージョンマップを作成する。最初のケースでは、ノードがクラスタを残す場合、一時的にあるいは永久に、そのノード上のMMによって管理されるリージョンを再び割り当てなければならない。ノードがサービスに戻る場合、あるいはノードが初めてクラスタを連結する場合、第2のケースはその状況を含む：そのような場合、クラスタ中の他のMMのためのロードを軽くするために、リージョンはそれに割り当てられる。新しいノード上で作られたリージョンは、すべて不完全である。これらのリージョンは、一旦それらがデータをロードし終えたと、バックアップに昇格される。不完全なリージョンがそのデータをロードすることを完了すると、3番目の状況が起こる。この時に、リージョンは、バックアップになる。マップ生成アルゴリズムは、好ましくは、所与のノードが正式なリージョンがクラスタにわたって平衡を保たれる、また、全てのリージョンがクラスタにわたって平衡を保たれるいかなるリージョンの1つを超えるコピーを決して含まないことを、保証する。全てのRGMが全てのメタデータ更新を処理し、このようにクラスタにわたって広げられるので、後者の2つの制約は必要である。正式なRGMは、また、検索要求を処理する、したがって、それらもまたよく分散される。

【0042】

下記に、マップ生成アルゴリズムに関する追加の詳細を提供する。

【0043】

MMリーダーが新しいマップを作成する必要がある場合、それが最初にするのはリージョンセンサスである。これは、現在クラスタ中の各ノード上のMMに要求を送信しながら要求/応答メッセージパターンを使用して行われる。要求/応答パターンは、好ましくは、どのリージョンの全体像がアーカイブに存在するかを形成して全ての応答が組み合わせられる集合ステップを含む。リージョンセンサスによって提供される情報は、好ましくは、各リージョンコピーの用に下記を含む：リージョンコピーを所有するノード、リージョンマネージャ（もしあれば）によって処理される最後の更新、及びリージョンのデータベーススキーマに格納されるタイムスタンプ。リージョンタイムスタンプは、センサスから削除される無効のリージョンを識別するために使用される。これは、無効のリージョンが形成されているマップから外され、また、無効のリージョンスキーマが削除されることを保証する。ほとんどの場合、無効のリージョンコピーは、現在のリージョンコピーからのマップ番号より低いマップバージョン番号を持つ。しかしながら、これは必ずしもそうだとは限らないかもしれない。例えば、新しいマップがノードクラッシュにより作成されていると仮定する。リージョンセンサスは、残りのリージョンを発見し、新しいマップを形成する。機能不全のノードがリージョンセンサスにตอบสนองするのに間に合うように再開すれば、あたかもうまく行かないものもなかったかのように、ノードはそのリージョンを報告する。しかしながら、これらのリージョンはすべてノードが下がっていた間に更新を逃したため、無効になるかもしれない。この問題の解決策は、リージョンセンサスで含まれるリージョンタイムスタンプを検査することである。リージョンコピーは、最後の更新の

10

20

30

40

50

タイムスタンプが処理されたことを表すそれぞれのリージョンタイムスタンプを報告する。リージョンコピーが同期された状態になるので、有効なタイムスタンプは、マップバージョン変更及び最初のマップを考慮に入れなければならない。機能不全のリージョンが最新の又は無効のマップバージョン番号を持っているかどうか、これが無効のリージョンを識別する。ノードが機能しなくなり、サービスに速く戻り、次に、無効のリージョンに基づいた要求を処理し始める危険はない。この理由は、ノードがリブートでリージョンマップを持たない、また、マップが受信されるまで、RGMは存在しないからである。RGMが作成されるまで、MMCからの要求は処理できない。したがって、新しいマップを得るまで、素早く再開する機能不全のノードは、要求を処理することができない、そして新しいマップはノードにその古いリージョンを廃棄させる。

10

【0044】

リージョンセンサスの後、最初のリージョンマップは以下のように生成される。リージョンセンサスが全くリージョンを放棄しない場合、クラスタは初めてスタートしているに違いない。この場合、正式なリージョンの所有者が最初に割り当てられる。各割り当てについては、アルゴリズムは、最小の使用中のノードを選択する。最小の使用中のノードは、最も少ないリージョンコピーを備えたノードである。つながりは、所有された正式なコピーの数に基づいて解決される。正式なリージョンの所有者が割り当てられた後、バックアップリージョンの所有者は、バランスのとれた権限があり合計のリージョンオーナー権を追い求めながら割り当てられる。新しいマップはすべてのMMに送られる、すなわち、その後、マップによって説明されたリージョンを作る。

20

【0045】

一旦クラスタがスタートしたならば、マップ変更は、好ましくは、順番に次のマップ変更を行うことにより実施される：(1)リージョンに正式なコピー(ノード機能不全による)がない場合、バックアップを昇格させる；(2)、リージョンがMDPLバックアップ以上に持っている場合、超過バックアップを削除する；(3)リージョンがMDPLバックアップ(ノード機能不全、あるいは正式なことへの昇格による)より少なく持っている場合、新しい不完全なリージョンコピーを作成する；(4)オーナー権のバランスを再び取る、そして(5)正式なオーナー権のバランスを再び取る。ステップ(4)は、最も忙しいノードを見つけ、オーナー権計算が少なくとも2低いノードにそのリージョンのうちの1つを再び割り当てることを含む。(目標ノードのオーナー権計算が1低い場合、再配分は、仕こと量の平衡を保つのを助けない)好ましくは、これは、新しい不完全なリージョンを作ることにより行われる。これがいかなるノードによって所有されたリージョンの最大数を縮小し続ける限り、この操作が継続される。ステップ(5)は、権限のあるリージョンの最大数を所有するノードを見つけることと、正式なオーナー権計算が少なくとも2低いバックアップを見つけることを含む。このステップは、例えば、バックアップを昇格させる及び正式なものを降格させることで責任を交換する。この操作は、いかなるノードによって所有される正式なリージョンの最大数を縮小し続ける限り、継続される。

30

【0046】

ノードがクラスタを残す場合、その後、ステップ(1)及び(3)は、ノードの離脱によって残されたリージョンマップ中のどんなギャップも満たす。必要ならば、その後、ステップ(4)と(5)が仕事量を一定にするために使われる。

40

【0047】

ノードがクラスタを連結する場合、ステップ(1)-(3)は、何も変更しない。ステップ(4)は、対照的に、新しいノードに割り当てられながら不完全なリージョンをセットもたらず。不完全なリージョンがそのデータをロードすることを完了する場合、それはMMリーダーに通知する。マップは、バックアップに不完全なリージョンを昇格させる。その後、ステップ(5)は、新しいノードに正式なリージョンを割り当てる効果がある。

【0048】

不完全なリージョンがその同期を終了する場合、それはバックアップリージョンに変わり、MMリーダーに通知する。その後、MMリーダーは、少なくとも1つのリージョン用

50

のTPOFバックアップ以上に含んでいる新しいマップを発行する。ステップ(2)は、最も極度にロードしたMMに対する負担を軽くすることを選びながら、超過しているバックアップリージョンを削除する。

【0049】

MMが新しいマップを受信する場合、それは、新しいマップと現在のものとの比較し、MMによって管理される各リージョンのためにいかなる変更をも適用する必要がある。可能な変更は、以下の通りである：リージョンを削除する、リージョンを作る、バックアップリージョンを正式なものに昇格させる、バックアップに不完全なリージョンを昇格させ、バックアップに正式なリージョンを降格させる。最初のタイプの変更にに関して、ロードバランスは、コピーの削除をもたらせながら、あるノードから別のノードにリージョンコピーのコントロールを移動させることができる。そのような場合、ネットワークとデータベース資源がリージョンのデータを格納するスキーマの削除を含めながら返される。正式なそしてバックアップリージョンが作られるとともに、リージョンを作る第2のタイプの変更が典型的に新しいクラスタに生じる。その後、不完全なリージョンだけが作られる。リージョン生成は、各タイプのメタデータオブジェクト用のテーブルを含んでいるデータベーススキーマを作成することを含む。各リージョンのスキーマは、リージョン(正式、バックアップ、あるいは不完全)の役割を識別する情報を含んでいる。3番目のタイプの変更、すなわち、バックアップから正式への昇格は、リージョンの役割の修正を必要とする。他の変更タイプは、それらの名前が意味するように、不完全からバックアップへあるいは正式からバックアップへのリージョンの役割を変更することを含む。

10

20

【0050】

ノードのメタデータマネージャはそれぞれ、全部のクラスタ用のメタデータの所与の部分をコントロールする。したがって、所与のノードに格納されるメタデータは、クラスタ中のすべての(あるいは所与のサブセットの)ノード中に理論上平等に分散されているデータベースと共に、分散型データベース(メタデータの)の一部から成る。メタデータマネージャは、説明されてきたように、この機能を達成するために協力する。新しいノードがクラスタに加えられる場合、個々のノード責任は、新しいキャパシティに調節される；これは新メンバーが等しい割り当てを仮定するように、すべてのノードにわたってメタデータを再分散することを含む。反対に、ノードが機能しなくなるか、クラスタから取り除かれる場合、他のノードメタデータマネージャは、より大きな割り当てを仮定することにより縮小されたキャパシティを補う。データロスを防ぐために、ノードがそれぞれ直接すべてのクラスタメタデータのあるパーセンテージを管理する責任があり、他のノードのセット番号にこのデータをコピーする場合、メタデータ情報は、好ましくは、複数のノードにわたって再現される。

30

【0051】

新しいマップが生成される場合、MMリーダーは、他のノードへそのマップの分散を始め、全てのノードがそれを持つまで処理の保留を要求する。一旦ノードがすべて新しいマップを持っていることをシステムが確認すれば、通常の処理が再開される。

【0052】

増分リフレッシュ

40

【0053】

上記の通りシステムにおいて、メタデータは、システムでノードにわたって冗長格納されているリージョンへ分散される。メタデータマネージャは、これらのリージョンの位置を含んでいるリージョンマップを持っていて、システムはこれによって適切に要求を送ることができるようになる。リージョンの数は、メタデータロードがすべてのノードにわたって分割される粒度を決定する。マップは、以下の段階、すなわち、正式なリージョンコピー((A)リージョン)、バックアップリージョンコピー((B)リージョン)、及び、スクラッチからあるいは(A)または(B)のリージョンから復元される過程にある不完全なリージョンコピー((I)リージョン)、におけるリージョンのコピーを含んでいる。(A)または(B)のリージョンがクラスタを残す場合、マップがリフレッシュさ

50

れる。(B)リージョンがクラスタを残している場合、(I)リージョンは、(A)リージョンからメタデータをすべてコピーすることにより作られ追加される。(A)リージョンがクラスタを残す場合、対応する(B)リージョンは、(A)に昇格され、次に、(I)リージョンが作られる。(I)リージョンコピーが完全な場合、それは(B)リージョンへ昇格される。リージョンがすべて完成した場合、マップは再度MDPLにある。このプロセス中に、失われたリージョンが戻る場合、そのタイムスタンプが無効であるので典型的に捨てられる。

【0054】

この開示によれば、リージョンへの最後に適用された変更を格納する概念に基づくエンハンスされたリカバリスキームが説明される。現在のタイムスタンプだけでなくその値は戻るリージョンをどうするか決めるために使用される。

10

【0055】

したがって、(A)または(B)のリージョンがすぐに(例えば単純なノードリブートにより)戻る場合、ほとんどの場合、リージョンは、ちょうど1更新だけ異なる。このケースを確認するために、この開示によれば、好ましくは、最後の更新は各リージョンのリージョンタイムスタンプテーブルで続けられる。最初のシナリオでは、正に最後の更新で見つからないリージョンがある。これは、(A)リージョンから更新を受信する前に、(B)リージョンが消失する場合である。この場合、システムは、現在のマップでそのリージョンを最新にさせるために、最後の更新を適用する。更新を適用した後に、マップインストールの残りは通常通り継続する。2番目のシナリオは、システムの他のいかなる場所にも繁殖しなかったリージョンに適用された更新があるということである。これは、(A)リージョンの更新がその(B)リージョンに適用される前に、(A)リージョンが消失する場合である。この場合、その更新は、無効であると考えられ、そのリージョンがマップに戻される前に、そのリージョンから取り除かれる。この場合、更新の除去が行われる間に、下記の通り、システムは「P」リージョンと呼ばれるものを作成する。

20

【0056】

最後に、正式なコピー(例えば戻るノードが利用不可能だった間に、リージョンへ書き込みがなかった場合)と戻るリージョンコピーが全く異なるケースがある。このケースでは、上記のリージョンタイムスタンプの比較によって、メタデータマネージャリーダーは、戻るリージョンが完全に最新で、したがって、サービスに直ちに(バックアップ(B)リージョンとして)戻ることができることを決めることができる。リージョンタイムスタンプ比較は、MMリーダーが暫くの間見つからなかったかもしれない不正確に戻るリージョンコピーを検査することを可能にする1以上の失敗したマップを説明するのに十分に柔軟である。

30

【0057】

ここでは、部分的、あるいは、「P」段階は、あるリージョンタイムスタンプの時点で最新の、そして現在のマップで最新でないあらゆるリージョンに当てはまるが、別のリージョンコピー上に格納された保留の更新を当てはめることにより最新にできる。システムは、リージョンコピーが失われたことを検知する場合、(I)リージョンが直ちに作られるので(前に述べたように)及びその(A)リージョンが受信する(リージョンコピーが失われたマップからの最後の更新で始まる)全ての更新をそれらの(I)リージョンが受信するので、(I)リージョンコピーは、(P)リージョンを最新にするために必要とされる保留の更新を正確に含んでいる。

40

【0058】

したがって、この開示によれば、(P)リージョンは、先在する(A)か(B)リージョンから復元されているリージョンを表わす。(P)リージョンはそれぞれ、(I)リージョン((P)リージョンコピーを最新にさせるために必要とされる保留の更新を正確に含んでいる)に格納された保留の変更から更新される。(B)または(I)リージョンのように、(P)リージョンが作られるとすぐに、(A)リージョンは、すべての更新のためにバックアップリージョンとしてそれを扱う。(P)リージョン(それが作成されると

50

すぐに)は、バックアップ要求を受信し、pending_updateテーブル(あるいは他の便利な持続メカニズム中で)にそれらを格納する。(P)リージョンが更新し終えた場合、それは(B)リージョンに変換される。この変換に際して、(I)リージョンは撤去される。

【0059】

最後の更新を削除する必要がある場合、Pリージョンは、それが他のリージョンコピー上で適用されなかった場合に適用された最後の更新を最初に取り消す。

【0060】

図7は、第1の例のシナリオを説明する。ここで、(B)リージョン700は失われており、(A)リージョン702から更新(100,4)を一度も受信していない。この場合に、(I)リージョン704は作られる。(A)リージョン702からメタデータをコピーすることにより(I)リージョンは追加される(示されるように);(I)リージョン704は、また、(A)リージョン702へのどんなそれに続く更新も受信し始める。(B)リージョン700が戻る場合、それを含んでいた最後のマップにおいて最新だったため、それは(P)リージョン706に変換される。(P)リージョン706は、(I)リージョン704の保留の更新テーブルからその見つからない更新を追加する。(I)リージョン704が最新である前に、このプロセスが終わる場合、(P)リージョン706は、(B)リージョンへ昇格され、(I)リージョンが廃棄される。

【0061】

図8は、第2の例を説明する。ここで、更新(100,4)(B)リージョン802を適用する前に、(A)リージョン800はダウンする。リージョン(A)がクラスタを残す場合、(B)リージョンは新しい(A)リージョンへ昇格される。新しい(I)リージョン804は、リージョン(A)から新しい更新を受信するために作られる。前の(A)リージョンが戻る場合、それは(P)リージョン806に変換される。上記の通り、(P)リージョン806は、リージョン(I)804の保留の更新テーブルからのその見つからない更新を追加する。(P)リージョンは、また、有効でない最後の更新を取り消す。リージョン(I)が最新である前に、このプロセスが終わる場合、リージョン(P)806は、(B)リージョンへ昇格され、(I)リージョンが落とされる。

【0062】

したがって、図7は、どのように一例、(B)リージョンのリストアの次の損失中に、(P)リージョンコピー(Pリージョン)が作成され、次に、(B)リージョンへ昇格されるかを説明する。図8は、どのように一例、(A)リージョンのリストアの次の損失中に、(P)リージョンコピーが作成され、次に、(B)リージョンへ昇格されるかを説明する。

【0063】

マップが存在するとそのマップが示す場合、(P)リージョンが作られ、既存のリージョンは、(P)リージョンに変換される。その場合に、リージョンも、それがその最後の更新を後退させる必要があり、どの(I)リージョンからコピーするかをマップ中で伝えられる。

【0064】

クラスタが、既存の正式なリージョンからバックアップリージョンを移動又は作る必要がある場合、リフレッシュタスク(RefreshTaskと呼ばれる)が始められる。そのタスクの中心となるのは、ソース及びターゲットマシンの両方におけるコピーテーブルスクリプトの起動である。スクリプトはSQLに基づくデータベースと直接交信する。RefreshTaskは、(P)リージョンの機能性を実施する。その他のオペレーションの前に、(P)リージョンは、適用された最後の更新が他のリージョンコピー上で適用されなかった場合、それを最初に取り消さなければならない。これは、Undo Last Updateと呼ばれるサブタスクである。そのため、RefreshTaskは、更新又は(A)リージョンから同等の列のコピーによって影響されたかもしれないデータベーステーブル中のいかなる列のローカルコピーを削除する。単一の保留の更新中に複数の更新があるかもしれないので、このプロセスは、保留の更新中ですべての個々の更新を処理する。各更新については、タスクは、影響される

10

20

30

40

50

データベーステーブルと影響を受けた列のSQL（構造化照会言語）WHERE句を決定する。このWHERE句は、次に、その更新用の影響を受けたテーブル中のいかなるローカルの列を最初に削除するのに使われ、（A）リージョンソース（copytableスクリプトを使用して）からのすべての同様の列をコピーする。Undo Last Updateによって修正される必要のあるメタデータのサブセットを選択するために、その句は、メタデータに述語を適用する。

【0065】

その後、RefreshTaskは、Copy Pending Updatesと呼ばれるサブタスクを始める。これは、マップ中で指定される（I）リージョンからpending_updateテーブルをコピーする。このコピーは、正常なバックアップ要求から入って来るpending_updatesとの対立を避けるためにpending_updateテーブルではなくpending_update_copyテーブルをターゲットとする。コピーが完全な場合、RefreshTaskは、2つのpending_updateテーブルを統合させる。これは、pending_updatesをソートし、どんな重複も除外する。

10

【0066】

その後、RefreshTaskは、Apply Pending Updatesと呼ばれるサブタスクを始める。具体的には、一旦pending_updateテーブルがフル実装されると、（P）リージョンは、好ましくは、2パスに、pendings_updatesを適用する。このプロセスの終わりに、リージョンのregion_timestampは通常更新される。

【0067】

一旦、保留の更新が適用されると、RefreshTaskは、バックアップリージョンへConvertと呼ばれるサブタスクを行なう。ここで、（P）リージョンは、それ自体を（B）リージョンに変換する。

20

【0068】

（P）リージョンが、（I）リージョンをリフレッシュするコストを下げるので、（I）リージョンは、（P）（B）変換が完了するまで前進し続ける。これを遂行するために、マップは、変換が起こる場合にIリージョンを落すように指定する使い捨てとしてそれらをマークするフラグを含む。Aリージョンは、そのIリージョンにバックアップ更新を送るのをやめることを知る必要があるので、この変換は、Aリージョンによって統合される。（P）（B）変換が起こる場合、（A）リージョンは、ロックを解除し、なんらかの使い捨てのIリージョン自体を除外するためにそれらのリージョンのそれぞれに1つの新しいメッセージを送りながら、発表されるリージョン用のローカルマップからそれらを撤去することによって、このメッセージに反応するために1つの新しいリスナーを持つ。（P）又は（I）リージョンRefreshTaskのいずれかのステップが機能しなくなる場合、そのリージョンは、再度使用されることができないようになる。

30

【0069】

上記の通り、増分リフレッシュは、多数の利点をもたらす。主要な利点は、（I）リージョンの長く高価なリフレッシュを回避することである。具体的には、リージョンセンス（マップ生成プロセスの一部としての）がMDPLコピーより少なく持っているリージョンを識別する場合、（I）リージョンは、新しいコピーとして作られる。この（I）リージョンは、リージョンの（A）コピーから完全なスキーマをコピーするためにRefreshTaskを起動する。このRefreshTaskは、データをコピーし、インデックスの作成を含むいくつかの高価なステップを持っている。ここに記述される技術の使用によって、失われたデータベースの部分に生じる更新だけが適用され、このオペレーションは、はるかに少ない時間を要し、冗長性を通じてその部分を再構築するよりはるかに少数のシステムリソースを消費する。修復する時間の短縮は、また、クラスタがMDPL及びピーク性能以下である時間を短縮する。

40

【0070】

ここに記述されるようなアーカイブ管理ソリューションは、デジタル資産のキャプチャー、保存、管理及び検索を可能にする。設計は、多数の必要条件のアドレスを指定する：既存のアプリケーションを備えた統合の無制限のストレージ、高い信頼度、自主管理と各種規格との適合、ハードウェア・インディペンデンス及び軽減。

50

【 0 0 7 1 】

Linux（例えば）を実行するコモディティハードウェアのクラスタは、ロバストプラットフォーム及びこと実上無制限のアーカイブを提供する。システムは、計ることができる、例えば、少数のストレージノードサーバから何千ものテラバイトデータを格納する多くのノードまで。アーキテクチャは、ストレージ容量が、組織の増加するアーカイブ要求と歩調を常に合わせることを保証する。

【 0 0 7 2 】

システムは、ファイルを決して失わないことを指定される。アーカイブがデバイス障害から常に保護されるように、それは、クラスタにわたってデータを複製する。ディスク又はノードが機能しなくなると、クラスタは、同じデータの複製を維持するクラスタ中の他のノードにわたって自動的に機能しなくなる。

10

【 0 0 7 3 】

システムは、自律的処理を通じてアーカイブストレージのコストを下げる。例えば、ノードがクラスタ化されたアーカイブを連結するか離れる場合、システムは自動的にクラスタのロードバランスを調節し、メンバーノードにわたってファイルを再分散することによりパフォーマンスを最適化する。

【 0 0 7 4 】

システムは、ユーザ定義の保存政策への従順を促進する。

【 0 0 7 5 】

システムは、オープンプラットフォーム上で展開することによりハードウェア依存を排除する。コモディティプラットフォーム及び所有者の記憶デバイス間のコストギャップが大きくなると、ITバイヤは、もはや高いコストの電気器具ベンダーとの関係にはまりたくない。所与のノードがコモディティハードウェア及びオープンソース（例えばLinux）オペレーティングシステムソフトウェア上で典型的に動作するので、好ましくは、バイヤは、最良のソリューションのために多くのハードウェア選択肢の中で買い物をすることができる。

20

【 0 0 7 6 】

システムは、また、ファイルを格納し検索するNFS、HTTP、FTP及びCIFSのような業界基準インターフェースを提供する。これは、システムが、カスタマイズされたアーカイブアプリケーションと同様に、ほとんどの標準コンテンツ管理システム、検索システム、ストレージ管理ツール（HSMとバックアップシステムのような）に容易に接続することができることを保証する。

30

【 0 0 7 7 】

上記のものは、ある実施例によって行なわれたオペレーションの特定な順序について記述しているが、代替実施例が異なる順にオペレーションを行い、あるオペレーションを組み合わせ、あるオペレーションをオーバーラップさせる等を理由にそのような順序が例となるということを承知してほしい。所与の実施例へのスペック中の参照は、記述される実施例が特別な特徴、構造あるいは特性を含むかもしれないことを示し、しかし、すべての実施例が必ずしも特別な特徴、構造あるいは特性を含むとは限らないかもしれない。

【 0 0 7 8 】

公開された技術は、方法又はプロセスのコンテキストに述べられているが、ここでの内容は、また、ここでのオペレーションを行なうための装置に関する。この装置は、所要の目的のために特に構築されてもよい、もしくは、それは、コンピュータに格納されたコンピュータプログラムによって選択的に活性化されるか再構成される汎用計算機から成ってもよい。そのようなコンピュータプログラムは、例えば、光ディスクを含む任意のタイプのディスク、CD-ROM及び光磁気ディスク、読み出し専用メモリ（ROM）、ランダムアクセスメモリ（RAM）、磁気又は、光カード、あるいは電子マニュアルを格納するのにふさわしい任意のタイプの媒体、そしてそれぞれがコンピュータシステムバスに接続されているだがこれに限らないは、コンピュータ読取り可能な記憶媒体に格納されてもよい。

40

50

【 0 0 7 9 】

システムの所与のコンポーネントは、別々に説明されているが、通常のスクリルのうちの1つは、所与の指示、プログラムシーケンス、コード部分などにおいて機能のうちのいくつかを組み合わせられる又は共有されるかもしれないことを認識する。

【 0 0 8 0 】

「固定コンテンツ」のためのアーカイブのコンテキストで本発明を述べたが、これもまた限定ではない。本明細書に述べた技術は、コンテンツに対する付加タイプ及び置換タイプの修正を可能にするストレージシステムにも等しく適用することができる。

本発明について述べてきたが、次に以下のとおり特許請求する。

【 図 1 】

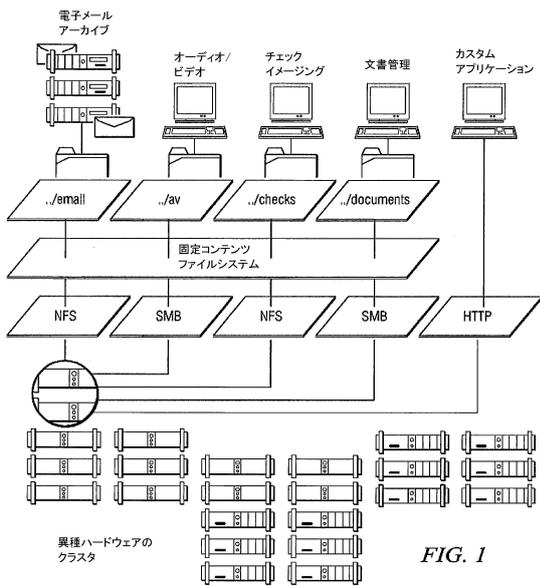


FIG. 1

【 図 2 】

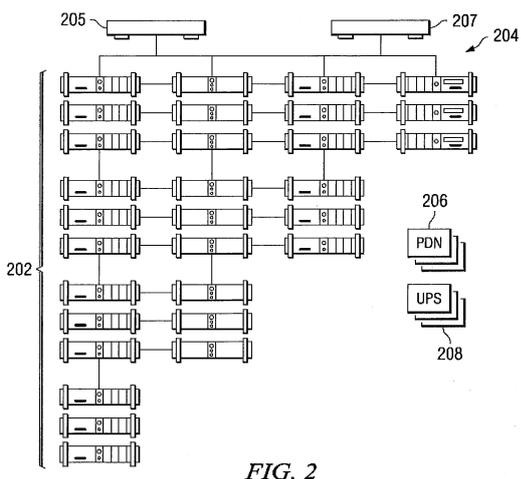
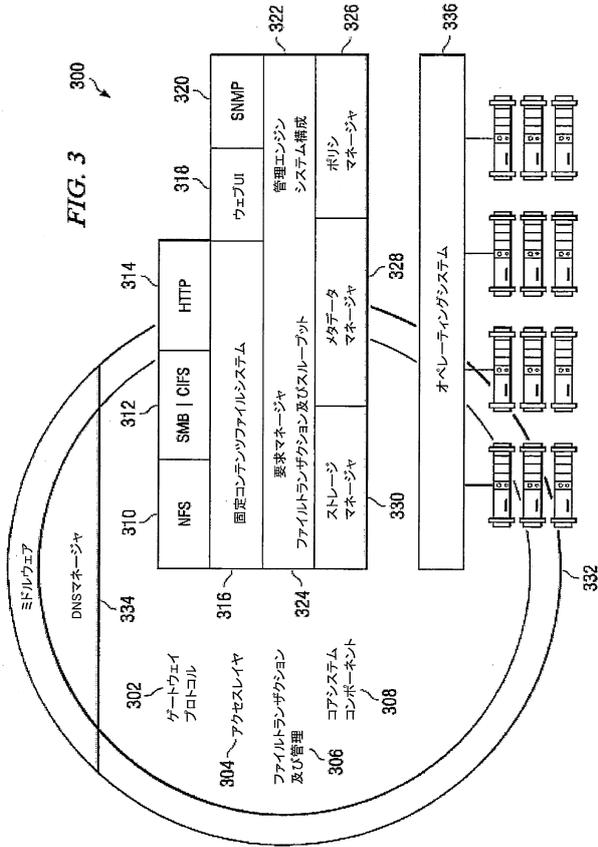
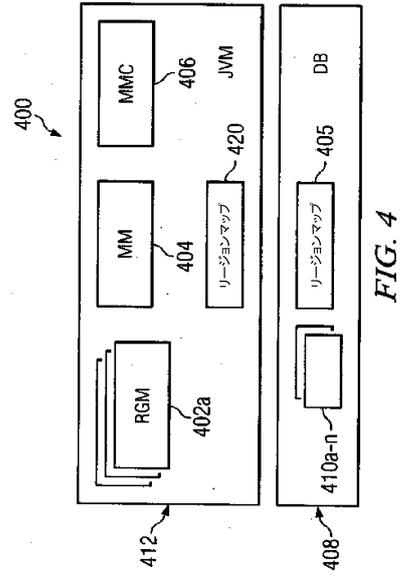


FIG. 2

【図3】



【図4】



【図5】

リージョン	正式	BACKUP 1	BACKUP 2
0	n1	n2	n3
1	n2	n3	n4
2	n3	n4	n1
3	n4	n1	n2
4	n1	n2	n3
...
63	n4	n1	n2

FIG. 5

【図6】

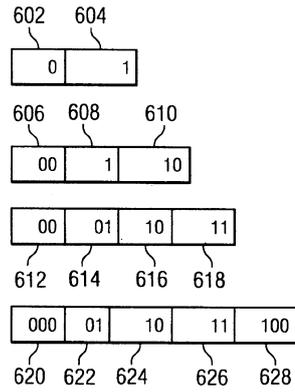


FIG. 6

【 図 7 】

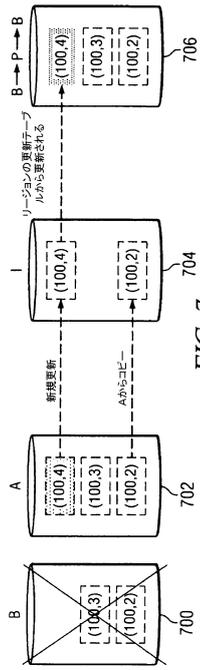


FIG. 7

【 図 8 】

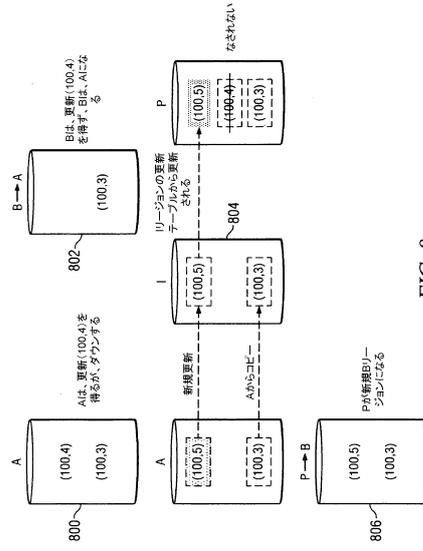


FIG. 8

フロントページの続き

- (72)発明者 グリマルディ, ケビン エス.
アメリカ合衆国 02144 マサチューセッツ州 サマービル, ベイステートアベニュー2番, 3
8
- (72)発明者 パーマー, トレック
アメリカ合衆国 02139 マサチューセッツ州 ケンブリッジ, パットナムアベニュー1番, 2
7
- (72)発明者 ピンクニー, デビッド
アメリカ合衆国 01810 マサチューセッツ州 アンドーバー, タングルウッドウェイ エヌ.
, 8

審査官 金沢 史明

- (56)参考文献 特表2008-518284(JP, A)
特開2001-290687(JP, A)
特開2006-072450(JP, A)
特開2005-141528(JP, A)

- (58)調査した分野(Int.Cl., DB名)
G06F 12/00