

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5798258号
(P5798258)

(45) 発行日 平成27年10月21日(2015.10.21)

(24) 登録日 平成27年8月28日(2015.8.28)

(51) Int.Cl.		F I			
G06F 3/06	(2006.01)	G06F 3/06	3 O 1 Z		
G06F 12/00	(2006.01)	G06F 3/06			
		G06F 3/06	3 O 1 X		
		G06F 3/06	3 O 4 N		
		G06F 3/06	3 O 4 F		
請求項の数 20 (全 25 頁) 最終頁に続く					

(21) 出願番号 特願2014-550282 (P2014-550282)
 (86) (22) 出願日 平成24年3月29日 (2012.3.29)
 (65) 公表番号 特表2015-509235 (P2015-509235A)
 (43) 公表日 平成27年3月26日 (2015.3.26)
 (86) 国際出願番号 PCT/US2012/031102
 (87) 国際公開番号 W02013/147783
 (87) 国際公開日 平成25年10月3日 (2013.10.3)
 審査請求日 平成26年6月30日 (2014.6.30)

(73) 特許権者 510071448
 ヒタチ データ システムズ コーポレー
 ション
 HITACHI DATA SYSTEM
 S CORPORATION
 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
 (74) 代理人 110000279
 特許業務法人ウィルフォート国際特許事務
 所
 (72) 発明者 ロジャース, リチャード
 アメリカ合衆国 98005 ワシントン
 州 ベルビュー, サウスイースト 4番地
 1002番, 11870
 最終頁に続く

(54) 【発明の名称】 記憶階層化のためのコンテンツ選択

(57) 【特許請求の範囲】

【請求項1】

1つ以上のノードを有し、各ノードが、プロセッサを有し、且つ、第1の階層としての複数の第1の論理記憶ユニットおよび第2の階層としての複数の第2の論理記憶ユニットの基になる複数の記憶媒体を有する記憶システムに接続され、

各ノードのプロセッサが、

前記複数の第1の論理記憶ユニットに複数の第1のコンテンツオブジェクトを記憶し、前記記憶システムに記憶すべき前記第1のコンテンツオブジェクトの数が前記第1のコンテンツオブジェクトのデータ保護レベルに基づいて規定されており、

前記複数の第1の論理記憶ユニットの使用された容量が閾値を上回る場合、前記第1のコンテンツオブジェクトのうちの1つの第1のコンテンツオブジェクト以外の全ての第1のコンテンツオブジェクトを前記複数の第2の論理記憶ユニットに移動し、且つ、前記複数の第1の論理記憶ユニットに前記1つの第1のコンテンツオブジェクトを維持する、ことを特徴とするシステム。

【請求項2】

前記複数の第1の論理記憶ユニットは、複数の稼働ユニットであり、前記複数の稼働ユニットは、それらの稼働ユニットに対応する1つ以上の第1のRAIDグループに対してスピンドウン機能の実行を可能にせず、

前記複数の第2の論理記憶ユニットは、複数のスピンドウンユニットであり、複数のスピンドウンユニットが、それらのスピンドウンユニットに対応する1つ以上の第2のRA

ＩＤグループに対してスピンドアウン機能の実行を可能にする、
ことを特徴とする請求項１に記載のシステム。

【請求項３】

前記プロセッサが、前記１つ以上のノードの論理パーティションであるネームスペースを管理し、

前記複数の第１の論理記憶ユニットにおけるどのコンテンツオブジェクトが前記複数の第２の論理記憶ユニットへ移動されるべき資格があるかを示す階層化ルールが、前記ネームスペースに設定され、

前記複数の第１の論理記憶ユニットの前記使用された容量が前記閾値を上回る場合、前記プロセッサが、前記階層化ルールに従って、前記複数の第２の論理記憶ユニットに前記ネームスペース内のコンテンツオブジェクトを移動する、
ことを特徴とする請求項１に記載のシステム。

10

【請求項４】

或るコンテンツオブジェクトのデータ保護レベルが１より大きい場合、前記階層化ルールが、前記ネームスペースにおける前記或るコンテンツオブジェクトのうちの１つの或るコンテンツオブジェクト以外の全ての或るコンテンツオブジェクトが前記複数の第２の論理記憶ユニットに移動されるべき資格があることを示す、
ことを特徴とする請求項３に記載のシステム。

【請求項５】

前記階層化ルールが、或る時間前に前記ネームスペースに取込まれたコンテンツオブジェクトが前記複数の第２の論理記憶ユニットへ移動されるべき資格があることを示す、
ことを特徴とする請求項３に記載のシステム。

20

【請求項６】

前記階層化ルールが、前記ネームスペース内のいずれのコンテンツオブジェクトも前記複数の第２の論理記憶ユニットへ移動されるべき資格がないことを示す、
ことを特徴とする請求項３に記載のシステム。

【請求項７】

階層化ルールを変更することによって前記第１のコンテンツオブジェクトのうちの１つの第１のコンテンツオブジェクト以外の全ての第１のコンテンツオブジェクトが、前記複数の第２の論理記憶ユニットに移動されるべき資格を有さないようになる場合、前記第１のコンテンツオブジェクトのうちの１つの第１のコンテンツオブジェクト以外の全ての第１のコンテンツオブジェクトが前記複数の第１の論理記憶ユニットへ戻される、
ことを特徴とする請求項１に記載のシステム。

30

【請求項８】

第１の階層としての複数の第１の論理記憶ユニットおよび第２の階層としての複数の第２の論理記憶ユニットの基になる複数の記憶媒体を有する記憶システムに接続された１つ以上のノードのための方法であって、

前記複数の第１の論理記憶ユニットに複数の第１のコンテンツオブジェクトを格納し、
前記記憶システムに記憶すべき前記第１のコンテンツオブジェクトの数が前記第１のコンテンツオブジェクトのデータ保護レベルに基づいて規定されており、

40

前記複数の第１の論理記憶ユニットの使用された容量が閾値を上回る場合、

前記第１のコンテンツオブジェクトのうちの１つの第１のコンテンツオブジェクト以外の全ての第１のコンテンツオブジェクトを前記複数の第２の論理記憶ユニットに移動し、
且つ、

前記複数の第１の論理記憶ユニットに前記１つの第１のコンテンツオブジェクトを維持する、
方法。

【請求項９】

前記複数の第１の論理記憶ユニットは、複数の稼働ユニットであり、前記複数の稼働ユニットは、それらの稼働ユニットに対応する１つ以上の第１のＲＡＩＤグループに対して

50

スピンドウン機能の実行を可能にせず、

前記複数の第2の論理記憶ユニットは、複数のスピンドウンユニットであり、複数のスピンドウンユニットが、それらのスピンドウンユニットに対応する1つ以上の第2のRAIDグループに対してスピンドウン機能の実行を可能にする、
ことを特徴とする請求項8に記載の方法。

【請求項10】

前記1つ以上のノードの論理パーティションであるネームスペースを管理し、

前記複数の第1の論理記憶ユニットにおけるどのコンテンツオブジェクトが前記複数の第2の論理記憶ユニットに移動されるべき資格があるかを示す階層化ルールを前記ネームスペースに設定する、
ことを更に実行し、

前記複数の第1の論理ユニットの前記使用された容量が前記閾値を上回る場合には、前記階層化ルールに従って、前記複数の第2の論理記憶ユニットに、前記ネームスペース内のコンテンツオブジェクトを移動する、
ことを特徴とする請求項8に記載の方法。

【請求項11】

或るコンテンツオブジェクトのデータ保護レベルが1より大きい場合、前記階層化ルールが、前記ネームスペースにおける前記或るコンテンツオブジェクトのうちの1つの或るコンテンツオブジェクト以外の全ての或るコンテンツオブジェクトが前記複数の第2の論理記憶ユニットに移動されるべき資格があることを示す、
ことを特徴とする請求項10に記載の方法。

【請求項12】

前記階層化ルールが、或る時間前に前記ネームスペースに取込まれたコンテンツオブジェクトが前記複数の第2の論理記憶ユニットへ移動されるべき資格があることを示す、
ことを特徴とする請求項10に記載の方法。

【請求項13】

前記階層化ルールが、前記ネームスペース内のいずれのコンテンツオブジェクトも前記複数の第2の論理記憶ユニットへ移動されるべき資格がないことを示す、
ことを特徴とする請求項10に記載の方法。

【請求項14】

前記階層化ルールを変更し、

階層化ルールを変更することによって前記第1のコンテンツオブジェクトのうちの1つの第1のコンテンツオブジェクト以外の全ての第1のコンテンツオブジェクトが、前記複数の第2の論理記憶ユニットに移動されるべき資格を有さないようになる場合、前記第1のコンテンツオブジェクトのうちの1つの第1のコンテンツオブジェクト以外の全ての第1のコンテンツオブジェクトを前記複数の第1の論理記憶ユニットへ戻す、
ことを特徴とする請求項8に記載の方法。

【請求項15】

第1の階層としての複数の第1の論理記憶ユニットおよび第2の階層としての複数の第2の論理記憶ユニットの基になる複数の記憶媒体を有する記憶システムに接続されたノードで実行されるコンピュータプログラムであって、

前記複数の第1の論理記憶ユニットに複数の第1のコンテンツオブジェクトを格納し、前記記憶システムに記憶すべき前記第1のコンテンツオブジェクトの数が前記第1のコンテンツオブジェクトのデータ保護レベルに基づいて規定されており、

前記複数の第1の論理ユニットの使用された容量が閾値を上回る場合、

前記第1のコンテンツオブジェクトのうちの1つの第1のコンテンツオブジェクト以外の全ての第1のコンテンツオブジェクトを前記複数の第2の論理記憶ユニットへ移動し、
且つ、

前記複数の第1の論理記憶ユニットに前記1つの第1のコンテンツオブジェクトを維持する

10

20

30

40

50

ことを前記ノードに実行させるコンピュータプログラム。

【請求項 16】

前記 1 つ以上のノードの論理パーティションである名前空間を管理し、

前記複数の第 1 の論理記憶ユニットにおけるどのコンテンツオブジェクトが前記複数の第 2 の論理記憶ユニットに移動されるべき資格があるかを示す階層化ルールを前記名前空間に設定する、

ことを更に前記ノードに実行させ、

前記複数の第 1 の論理ユニットの前記使用された容量が前記閾値を上回る場合には、前記階層化ルールに従って、前記複数の第 2 の論理記憶ユニットに、前記名前空間内のコンテンツオブジェクトを移動する、

10

ことを特徴とする請求項 15 に記載のコンピュータプログラム。

【請求項 17】

或るコンテンツオブジェクトのデータ保護レベルが 1 より大きい場合、前記階層化ルールが、前記名前空間における前記或るコンテンツオブジェクトのうちの 1 つの或るコンテンツオブジェクト以外の全ての或るコンテンツオブジェクトが前記複数の第 2 の論理記憶ユニットに移動されるべき資格があることを示す、

ことを特徴とする請求項 16 に記載のコンピュータプログラム。

【請求項 18】

前記階層化ルールが、或る時間前に前記名前空間に取込まれたコンテンツオブジェクトが前記複数の第 2 の論理記憶ユニットへ移動されるべき資格があることを示す、

20

ことを特徴とする請求項 16 に記載のコンピュータプログラム。

【請求項 19】

前記階層化ルールが、前記名前空間内のいずれのコンテンツオブジェクトも前記複数の第 2 の論理記憶ユニットへ移動されるべき資格がないことを示す、

ことを特徴とする請求項 16 に記載のコンピュータプログラム。

【請求項 20】

前記階層化ルールを変更し、

階層化ルールを変更することによって前記第 1 のコンテンツオブジェクトのうちの 1 つの第 1 のコンテンツオブジェクト以外の全ての第 1 のコンテンツオブジェクトが、前記複数の第 2 の論理記憶ユニットに移動されるべき資格を有さないようになる場合、前記第 1 のコンテンツオブジェクトのうちの 1 つの第 1 のコンテンツオブジェクト以外の全ての第 1 のコンテンツオブジェクトを前記複数の第 1 の論理記憶ユニットへ戻す、

30

ことを更に前記ノードに実行させることを特徴とする請求項 15 に記載のコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、一般に、記憶システムに、及びより詳しくは、記憶階層化のためのコンテンツ選択に基づいてコンテンツシステム内のコンテンツの記憶を管理するためのシステム及び方法に関する。

40

【背景技術】

【0002】

固定コンテンツシステム (FCS) は、変化しないデータを収容している。コンテンツの僅かな割合だけが、実際にその後アクセスされる；しかしながら、データはパワー資源を継続的に消費している記憶媒体上になおとどまる。たとえ媒体上のコンテンツの多くがアクセスされなくても、パワーは媒体をスピンさせ続けて無駄にされる。

【発明の概要】

【発明が解決しようとする課題】

【0003】

米国特許第 8,006,111 号明細書が、スピンドアウンの概念を開示する。このアプ

50

ローチの下で、ファイル移動及びディスクドライブのパワー状態が、ファイルの活動に基づいて決定される。活性/不活性ディスクドライブへのアクセスが活性/不活性ディスクドライブをパワーダウンするための不活性閾値に到達すると、各活性/不活性ディスクドライブが類似した確率のアクセスを有するファイルをロードされるように、能動記憶域内に不活性になっていたファイルの群が、それぞれの空にされた活性/不活性ディスクドライブに連続して移行される。焦点は、個々のディスクドライブにある。不活性閾値が到達された時だけ、ファイルが移動される。

【課題を解決するための手段】

【0004】

本発明の例示的な実施態様が、複製されたオブジェクト記憶システムまたはコンテンツプラットフォームであることができ、かつ、記憶媒体を2つのタイプ、すなわち、稼働(Run)ユニット(RU)及びスピンドウンユニット(SDU)に分類する、固定コンテンツシステム(FCS)を提供する。RUはおそらく使用されるコンテンツに対して指定され、及びSDUはおそらく使用されないコンテンツを収容する。定期的に、FCSはRU媒体とSDU媒体との間を移動させるための候補を識別する設定可能なサービスプラン及び記憶階層化ポリシーに基づいてコンテンツを評価する。記憶階層化は、オブジェクト移動及び装置操作を最適化するためにシステム状態に基づいてシステム挙動を決定し、それで、システム状態がニーズを決定する時だけ、データ移動が実行される。監視されるシステム状態は、オブジェクトが存在するRUが消費利用閾値(すなわちスペース/記憶消費閾値)に到達したかどうか、かつ、オブジェクトのためにSDU上で利用可能なスペースがあるかどうかである。消費利用閾値の一例は、消費された記憶域である(例えば、70%消費された記憶域がデフォルトとして用いられることができる)。記憶階層化ルール(STR)は、SDUまたはRU上に存在するオブジェクト立候補のためのポリシーを特定する。特定の実施態様において、設定は「決してない」「保護コピーだけ」または、「取込後N日」である。STRは、FCS内の全てのオブジェクトまたは一まとまりのオブジェクトに適用されることができる。

【0005】

特定の実施態様において、階層化のためのコンテンツ識別が、バックグラウンドで実行される。コンテンツが識別される時、SDU媒体が(必要に応じて)スピニアップされ、及びオブジェクト(複数オブジェクト)が移動される。不活性タイムアウトの後、SDU媒体はスピンドウンされる。コンテンツがFCSから要求され、及びそれがSDU上に存在する場合、SDUがスピニアップされ、及びコンテンツが読み出されて要求側に返される。不活性タイムアウトの後、SDU媒体は再びスピンドウンされる。

【0006】

本発明は、米国特許第8,006,111号明細書と種々の点において異なる。例えば、ファイルを移動させる判定は、個々のファイルの活動に基づかない。その代わりに、ファイル移動を確立するための基準が、ファイルの経過時間、冗長バックアップコピーの存在、などを含む。米国特許第8,006,111号明細書が個々のディスクドライブに焦点を合わせるとはいえ、本発明の特定の実施態様に従う解決策はディスク配列内に実現される機能を包含し、及び、焦点は一まとまりのディスクドライブの形のRAID群にある。さらに、それが存在する記憶域で、スペース/記憶消費閾値が到達されない限り、本発明の特定の実施態様に従う解決策は能動記憶域からコンテンツを移動させない。対照的に、米国特許第8,006,111号明細書は不活性閾値が到達された時にだけファイルを移動させる。

【0007】

固定コンテンツシステムは、スピンドウン機能を可能にしない稼働ユニット及びスピンドウン機能を可能にするスピンドウンユニットを含む複数の記憶ユニット内にコンテンツを記憶する複数の独立ノードを有する。本発明の一態様に従って、コンテンツの記憶を管理する方法が：どんなコンテンツが稼働ユニット上に記憶されるべき資格があるか、及びどんなコンテンツがスピンドウンユニット上に記憶されるべき資格があるかを示すポリシ

10

20

30

40

50

ーを設定する記憶階層化ルールであって、記憶ユニット内の記憶及び記憶ユニット間の移行のためにその記憶されたコンテンツの適格性を判定するためにコンテンツシステム内の少なくとも一群のコンテンツに適用可能である記憶階層化ルールを確立するステップ；コンテンツシステムの状態を監視するステップ；ならびに、記憶階層化ルール、コンテンツシステムの状態及び少なくとも一群のコンテンツの記憶されたコンテンツの適格性に基づいて、稼働ユニットとスピンドアウンユニットとの間の移行を含む記憶ユニット間を移行するコンテンツの候補を識別するステップを含む。

【0008】

いくつかの実施態様において、記憶階層化ルールが、期待される使用、ライフサイクル及びコンテンツの経過時間及びコンテンツの1つ以上の冗長バックアップコピーの存在を含む一組の基準に基づいて記憶されたコンテンツの適格性を判定するためにコンテンツを評価するように確立される。記憶階層化ルールが、決してスピンドアウンされない記憶ユニット上に記憶されるべきコンテンツにあてはまる「けっしてない」ルール、スピンドアウンユニット内の記憶の候補であるバックアップコピーのためのコンテンツにあてはまる「保護コピーだけ」ルール、及びスピンドアウンユニットに記憶されるべき候補になるのに十分である一定の時間Xの間コンテンツシステム内に存在したコンテンツにあてはまる「取込後X時間」ルールを含む。記憶階層化ルールは、資格のあるコンテンツが存在する稼働ユニットが設定可能な消費利用閾値に到達した時にだけ、資格のあるコンテンツをスピンドアウンユニットに移行するために確立される。コンテンツシステムの状態を監視するステップが、コンテンツが存在する稼働ユニットが、コンテンツをスピンドアウンユニットに移行するパーミッションを示す設定可能な消費利用閾値に到達したかどうかを判定するステップ、及び稼働ユニットからコンテンツを収容するためにスピンドアウンユニット上に利用可能なスペースがあるかどうかを判定するステップを含む。

【0009】

特定の実施態様において、この方法が識別に基づいてコンテンツを移行するステップ；及びスピンドアウンユニットの各々の状態を、コンテンツを記憶するためのシステムニーズに基づいてパワーアップまたはパワーダウンの適切な状態に管理するステップ、ならびに監視、識別及び移行の結果として、コンテンツを移行するステップを更に含む。この方法が、記憶サブシステム内の稼働ユニット及びスピンドアウンユニットを含む記憶ユニットの異なる記憶クラスの組合せを利用する記憶階層化戦略、及び記憶ユニットの異なる記憶クラス内のコンテンツの記憶のための記憶されたコンテンツの適格性を判定するポリシーを設定する記憶階層化ルールを各々特定する複数のサービスプランからサービスプランを選ぶステップを更に含む。

【0010】

本発明の別の態様が、スピンドアウン機能を可能にしない稼働ユニット及びスピンドアウン機能を可能にするスピンドアウンユニットを含む複数の記憶ユニット内にコンテンツを記憶する複数の独立ノードを有するコンテンツシステム内のコンテンツの記憶を管理するための器具を目的とする。器具は、プロセッサ、メモリ及び記憶階層化サービスモジュールを備える。記憶階層化サービスモジュールが：どんなコンテンツが稼働ユニット上に記憶されるべき資格があるか、及びどんなコンテンツがスピンドアウンユニット上に記憶されるべき資格があるかを示すポリシーを設定する記憶階層化ルールであって、記憶ユニット内の記憶及び記憶ユニット間の移行のためにその記憶されたコンテンツの適格性を判定するためにコンテンツシステム内の少なくとも一群のコンテンツに適用可能な記憶階層化ルールを確立し；コンテンツシステムの状態を監視し；かつ、記憶階層化ルール、コンテンツシステムの状態及び少なくとも一群のコンテンツの記憶されたコンテンツの適格性に基づいて、稼働ユニットとスピンドアウンユニットとの間の移行を含む記憶ユニット間を移行するコンテンツの候補を識別するように構成される。

【0011】

実施態様によっては、この器具が識別に基づいてコンテンツを移行するように構成される移行モジュール、及びコンテンツを記憶し、かつ監視、識別及び移行の結果として、コ

10

20

30

40

50

ンテンツを移行するためのシステムニーズに基づいて、スピンドウンユニットの各々の状態をパワーアップまたはパワーダウンの適切な状態に管理するように構成される記憶ユニット状態管理モジュールを更に備える。この器具は、記憶サブシステム内の稼働ユニット及びスピンドウンユニットを含む記憶ユニットの異なる記憶クラスの組合せを利用する記憶階層化戦略、ならびに記憶ユニットの異なる記憶クラス内のコンテンツの記憶のための記憶されたコンテンツの適格性を判定するポリシーを設定する記憶階層化ルールを各々特定する複数のサービスプランからサービスプランを選ぶためのユーザインタフェースを提供するように構成されるサービスプラン選択モジュールを更に備える。スピンドウン機能を可能にするように構成されないRAID群上に稼働ユニットが収容され、そして、スピンドウンを可能にするように構成されるRAID群上にスピンドウンユニットが収容されている。記憶ユニットは、異なる信頼性、性能またはコスト特性の少なくとも1つを有する異なるクラスの記憶ユニットを含む。この器具は、識別に基づいてコンテンツを移行し、かつコンテンツのメタデータ基準に基づいて異なるクラスの稼働ユニットとスピンドウンユニットとの間でコンテンツを移動させるように構成される移行モジュールを更に備え、このオブジェクトメタデータ基準が「データのタイプ」、「最後のアクセスからの時間」、「取込からの時間」及び「コンテンツのバージョン」の1つ以上を含む。

10

【0012】

本発明の別の態様が、スピンドウン機能を可能にしない稼働ユニット及びスピンドウン機能を可能にするスピンドウンユニットを含む複数の記憶ユニット内にコンテンツを記憶する複数の独立ノードを有するコンテンツシステム内のコンテンツの記憶を管理するデータプロセッサを制御するための複数の命令を記憶するコンピュータ可読の記憶媒体を目的とする。この複数の命令が：どんなコンテンツが稼働ユニット上に記憶されるべき資格があるか及びどんなコンテンツがスピンドウンユニット上に記憶されるべき資格があるかを示すポリシーを設定する記憶階層化ルールであって、記憶ユニット内の記憶及び記憶ユニット間の移行のためにその記憶されたコンテンツの適格性を判定するようにコンテンツシステム内の少なくとも一群のコンテンツに適用可能な記憶階層化ルールをデータプロセッサに確立させる命令；コンテンツシステムの状態をデータプロセッサに監視させる命令；ならびに記憶階層化ルール、コンテンツシステムの状態及び少なくとも一群のコンテンツの記憶されたコンテンツの適格性に基づいて、稼働ユニットとスピンドウンユニットとの間の移行を含む記憶ユニット間を移行するコンテンツの候補をデータプロセッサに識別させる命令を含む。

20

30

【0013】

本発明のこれらの、そしてまた他の、特徴及び効果が、特定の実施態様の以下の詳細な説明を考慮して当業者に明白になるであろう。

【図面の簡単な説明】**【0014】**

【図1】本発明の方法と器具が適用されることができる固定コンテンツ記憶アーカイブの簡略ブロック図である。

【図2】その各々が対称性でかつ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の冗長配列の簡単にされた表現である。

40

【図3】所定のノード上で実行するアーカイブクラスタアプリケーションの種々のコンポーネントのハイレベル表現である。

【図4】固定コンテンツを記憶するためにブロックベースの記憶サブシステムに連結される一群のノードを有する固定コンテンツシステムの簡略図である。

【図5】記憶階層化サービス(STS)モジュール、ノード内の記憶媒体ユニット及びそれらの状態をリストする記憶媒体状態テーブルならびに記憶マネージャの複数のインスタンスを有するノードに関して略図で例示する簡略図である。

【図6】記憶階層化のためのコンテンツ選択のプロセスを例示する一例を示す。

【図7】STSモジュールによって実行されるSTSを例示する流れ図の一例である。

【図8】(a) サービスプランを作り出してかつ階層化ポリシーを選ぶためのスクリーン

50

ショット、(b)作り出されるべきネームスペースのためのサービスプランを選ぶ時何が起こるかについて示すスクリーンショット及び(c)既存のネームスペース上のサービスプランをどのように変更するべきかについて示すスクリーンショットを含む記憶階層化戦略を各々特定する複数のサービスプランからサービスプランを選ぶためのユーザインタフェースの一例を示す。

【発明を実施するための形態】

【0015】

本発明の以下の詳細な説明では、参照が開示の一部を形成する添付の図面になされ、及び本発明が実施されることが出来る例示的な実施態様を例証としてかつ限定でなく示す。図において、同様な数字がいくつかの図の全体にわたって実質的に類似したコンポーネントを記述する。更に、詳細な説明が種々の例示的な実施態様を提供するとはいえ、後述するようにかつ図面内に図示するように、本発明は、本願明細書に図と共に説明される実施態様に限定されないが、しかし、公知であるように、または当業者にとって公知になるように他の実施態様に拡張することができる。「一実施態様」、「この実施態様」、または「これらの実施態様」への明細書における参照は、実施態様と関連して記述される特定の特徴、構造または特性が本発明の少なくとも1つの実施態様内に含まれることを意味し、及び、明細書内の種々の場所におけるこれらの句の出現が、必ずしも全て同じ実施態様を参照するというわけではない。加えて、以下の詳細な説明において、数多くの具体的な詳細が本発明の徹底的な理解を提供するために記載される。しかしながら、これらの具体的な詳細が、本発明を実践するために全てが必要であるというわけではないことは、当業者に明白である。他の状況では、本発明を不必要に不明瞭にしないために、周知の構造、材料、回路、プロセス及びインタフェースは詳述せず、及び/またはブロック図形式で例示されることが出来る。

【0016】

さらに、あとに続く詳細な説明のいくつかの部分が、コンピュータ内の操作のアルゴリズム及び記号表示に関して提示される。これらのアルゴリズム的記述及び記号表示は、データ処理技術の当業者によってそれらの革新の本質を最も効果的に他の当業者に伝えるために使用される手段である。アルゴリズムは、所望の終了状態または結果に至る一連の規定されたステップである。本発明において、実施されるステップは、具体的な結果を達成するための具体的な数量の物理操作を必要とする。通常、必然的にではないとはいえ、これらの数量は、記憶され、伝達され、組み合わせられ、比較され、かつさもなければ操作されることが可能な電気もしくは磁気信号または命令の形式をとる。これらの信号をビット、値、要素、シンボル、文字、用語、数、命令、などとして参照することは、主に一般的な使用の理由で折に触れて都合がいいと立証されている。しかしながら、これら及び類似の用語の全てが適切な物理量と関連しているべきであり、かつこれらの数量に適用される単に都合がいいラベルであるだけであることが念頭に置かれるべきである。とりわけ別の方法で述べられない限り、以下の説明から明白に認識されることは、記述の全体にわたって、「処理」、「計算」、「算出」、「判定」、「表示」などのような用語を利用する考察は、コンピュータシステムのレジスタ及びメモリ内で物理(電子)数量として表されるデータを、コンピュータシステムのメモリまたはレジスタまたは他の情報記憶、送信もしくはディスプレイ装置内で物理量として同じように表される他のデータに、操作してかつ変換するコンピュータシステムまたは他の情報処理装置のアクション及びプロセスを含むことができるということである。

【0017】

本発明は、さらにここで操作を実行するための器具に関する。この器具は必要とされた目的のために特別に構成されることが出来るかまたは、それが1つ以上の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または再構成される1台以上の汎用コンピュータを含むことが出来る。この種のコンピュータプログラムは、限定されないが、光ディスク、磁気ディスク、読取り専用メモリ、ランダムアクセスメモリ、固体素子装置及びドライブまたは電子情報を記憶することに適しているその他のタイプの媒体のようなコンピュ

10

20

30

40

50

ータ可読の記憶媒体内に記憶されることができる。本願明細書に提示されるアルゴリズム及び表示は、いかなる特定のコンピュータまたは他の器具にも本質的に関連がない。種々の汎用システムが、本願明細書における教示に従うプログラム及びモジュールとともに使用されることができるか、または、所望の方法ステップを実行するためにより専門の器具を構成することが、都合がいいと立証されることができる。加えて、本発明は任意の特定のプログラミング言語を参照して記述されない。理解されるであろうことは、本願明細書に記述されるように、種々のプログラミング言語が本発明の教示を実現するために使用されることができるということである。プログラミング言語（複数言語）の命令は、1台以上の処理装置、例えば中央処理ユニット（CPU）、プロセッサまたはコントローラによって実行されることができる。

10

【0018】

後で詳しく述べるように、本発明の例示的な実施態様は、記憶階層化のためのコンテンツ選択に基づいて固定コンテンツシステム内の固定コンテンツの記憶を管理するための器具、方法及びコンピュータプログラムを提供する。

【0019】

I. 固定コンテンツ分散データ記憶

【0020】

従来のテープ及び光記憶解法を置換するかまたは補充する高度に利用可能な、信頼性が高いかつ持続的な方法での、「固定コンテンツ」のアーカイブ記憶に対するニーズが出現した。用語「固定コンテンツ」は、一般的に参照または他の目的のために変更なしで保持されるのを期待される任意のタイプのデジタル情報を指す。この種の固定コンテンツの例は多くの他の中に、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、などを含む。従来の独立ノードの冗長配列（RAIN）記憶アプローチは、この種の固定コンテンツ情報資産の記憶用の大きなオンラインアーカイブを作り出すための選り抜きのアーキテクチャとして現れた。ノードを必要に応じてクラスタに接続してかつそれから出ることを可能にすることによって、RAINアーキテクチャは1つ以上のノードの故障から記憶クラスタを遮断する。複数ノード上にデータを複製することによって、RAINタイプアーカイブはノード故障または除去を自動的に補正することができる。一般的に、RAINシステムは主として閉システム内の同一のコンポーネントから設計されたハードウェア器具として供給される。

20

30

【0021】

図1は、そのような拡張可能なディスクベースのアーカイブ記憶管理システムを例示する。ノードは、異なるハードウェアを備えることができしたがって、「異種」とみなされることができる。ノードは一般的に、実際の物理記憶ディスクまたは記憶領域ネットワーク（SAN）内における仮想記憶ディスクであることができる1台以上の記憶ディスクにアクセスする。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、そのアプリケーションが実行する基本オペレーティングシステム）は同じであるかまたは実質的に同じであることができる。各ノード上のソフトウェアスタック（それはオペレーティングシステムを含むことができる）は、対称であるが、一方、ハードウェアは異質であることができる。このシステムを使用して、図1にて図示したように、とりわけ、企業は文書、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、などのような多くの異なるタイプの固定コンテンツ情報用のパーマネント記憶を作り出すことができる。もちろん、これらのタイプは、単に例証となるだけである。高レベルの信頼性が、独立サーバまたはいわゆる記憶ノード上にデータを複製することによって達成される。好ましくは、各ノードはそのピアと対称である。したがって、好ましくは、任意の所定のノードが全ての機能を実行することができるので、任意の1つのノードの故障はアーカイブの利用可能度にほとんど影響を及ぼさない。

40

【0022】

自己の米国特許第7,155,466号明細書にて説明したように、デジタル資産を収集して、保存して、管理してかつ読み出す各ノード上で実行される分散ソフトウェアアプ

50

リケーションを組み込むことは、RAINベースのアーカイブのシステムにおいて公知である。図2は、そのようなシステムを例示する。個々のアーカイブの物理境界は、クラスタ（またはシステム）と称する。一般的に、クラスタは単一装置でなく、しかし、むしろ一まとまりの装置である。装置は、均質であるかまたは異質であることができる。典型的装置は、Linuxのような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機械である。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、2、3の記憶ノードサーバから何千ものテラバイト単位のデータを記憶する多くのノードまでスケールアップされることができるアーカイブを提供する。このアーキテクチャは、記憶容量が組織の増大するアーカイブ要件と常に足並みをそろえることができることを確実にする。

10

【0023】

上記したような記憶システムにおいて、データは一般的にアーカイブが装置故障から常に保護されるようにクラスタ全体にランダムに分散される。ディスクまたはノードが故障する場合、クラスタは、同じデータの複製を維持するクラスタ内の他のノードまで自動的に故障する。このアプローチがデータ保護の見地からよく機能するとはいえ、クラスタのデータロスに対する算出された平均時間(MTDL)が要望するほど高くないかもしれない。特に、MTDLは一般的にアーカイブがデータを失う前に算出された時間を表す。デジタルアーカイブにおいて、いかなるデータロスも望ましくないが、しかし、ハードウェアコンポーネント及びソフトウェアコンポーネントの性質に起因して、(いかにみても)この種の発生の可能性が常にある。アーカイブクラスタ内のオブジェクト及びそれらのコピーのランダム分布の理由でたとえば、所定のノード内の所定のディスク(ここではミラーコピーが記憶される)が予想外に故障する場合、オブジェクトの必要とされたコピーが利用できないかもしれないので、MTDLは結局必要とされるより低いことになるかもしれない。

20

【0024】

図2に示すように、本発明が実現される例証となるクラスタが好ましくは、以下の全般的なカテゴリのコンポーネントを備える：ノード202、一对のネットワークスイッチ204、配電器(PDU)206及び無停電電源(UPS)208。ノード202は、一般的に1台以上のコモディティサーバを備えてかつ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台以上のハードディスク(例えば標準IDE/SATA、SCSI、など)及び2枚以上のネットワークインタフェースカード(NIC)を収容している。典型的ノードは、2.4GHzのチップ、512MB RAM及び6台の(6台の)200GBハードディスクを備えた2Uラックマウントユニットである。しかしながら、これは限定でない。ネットワークスイッチ204は、一般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クラスタアクセスを可能にする外部スイッチ207を備える。各スイッチは、クラスタ内の全ての可能性のあるノードを取り扱うのに十分なポートを必要とする。イーサネットまたはGigEスイッチが、このために使用されることができる。PDU 206は全てのノード及びスイッチを動かすために使用され、全てのノード及びスイッチを保護するUPS 208が使用される。限定するつもりではないとはいえ、一般的に、クラスタは公共インターネット、企業イントラネットまたは他の広域もしくはローカルエリアネットワークのようなネットワークに接続可能である。例示の実施態様では、クラスタは企業環境内に実現される。たとえばサイトの会社DNSサーバを通してナビゲートすることによって、それが到達されることができる。したがって、例えば、クラスタのドメインは既存のドメインの新規のサブドメインであることができる。代表実現において、サブドメインは会社DNSサーバにおいてクラスタ自体のサーバに委任される。エンドユーザは、任意の従来インタフェースまたはアクセスツールを使用してクラスタにアクセスする。したがって、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)を介して、API経由で、またはその他の周知のもしくは後で開発されたアクセス方式、サービス、プログラムもしくはツール

30

40

50

を通して実施されることができる。

【 0 0 2 5 】

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または HTTP API のような 1 つ以上のタイプの外部ゲートウェイを通してクラスタにアクセスする。アーカイブは好ましくは任意の標準 Unix ファイルプロトコル志向ファシリティの下に任意選択で位置することができる仮想ファイルシステムを通して露出される。これらは、NFS、FTP、SMB / CIFS、などを含む。

【 0 0 2 6 】

一実施態様において、アーカイブクラスタアプリケーションがクラスタとして共にネットワーク化される（例えばイーサネット経由で）独立ノードの冗長配列（H - R A I N）上で動作する。所定のノードのハードウェアは、異質であることができる。しかしながら最大信頼性のために好ましくは各ノードが、次に図 3 内に例示されるようにいくつかのランタイムコンポーネントを備える、分散アプリケーションのインスタンス 3 0 0（それは同じインスタンスまたは実質的に同じインスタンスであることができる）を実行する。したがって、ハードウェアが異質であることができるとはいえ、（少なくともそれが本発明に関する）ノード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 3 0 2、アクセスレイヤ 3 0 4、フィルタランザクション及び管理レイヤ 3 0 6 及びコアコンポーネントレイヤ 3 0 8 を備える。機能が他の意味がある方法で特徴づけられることができると当業者が認識するように、「レイヤ」名称は説明的な目的のために提供される。レイヤ（またはその中のコンポーネント）の 1 つ以上が、一体化されるかまたはその逆であることができる。いくつかのコンポーネントが、レイヤ全体に共有されることができる。

【 0 0 2 7 】

ゲートウェイプロトコルレイヤ 3 0 2 内のゲートウェイプロトコルが、既存のアプリケーションへの透明性を提供する。特に、ゲートウェイは NFS 3 1 0 及び SMB / CIFS 3 1 2 のようなネイティブファイルサービス、同じくカスタムアプリケーションを構築するウェブサービス API を提供する。HTTP サポート 3 1 4 が、さらに提供される。アクセスレイヤ 3 0 4 が、アーカイブへのアクセスを提供する。特に、本発明によれば、固定コンテンツファイルシステム（FCFS）3 1 6 がアーカイブオブジェクトへの完全なアクセスを提供するネイティブファイルシステムをエミュレートする。FCFS は、あたかもそれらが通常ファイルであるかのように、アーカイブコンテンツへのアプリケーションの直接アクセスを与える。好ましくは、アーカイブされたコンテンツがその本来の形態で表示され、一方メタデータがファイルとして露出される。彼らによく知られている方法で、管理者が固定コンテンツデータを供給することができるように、FCFS 3 1 6 はディレクトリ及びパーミッション及びルーチンファイルレベル呼出しの従来のビューを提供する。ファイルアクセス呼出しは好ましくはユーザスペースデーモンによってインターセプトされてかつ呼出しアプリケーションに対する適切なビューを動的に作り出す（レイヤ 3 0 8 内の）適切なコアコンポーネントに経由される。FCFS 呼出しは好ましくは、自律性のアーカイブ管理を容易にするためにアーカイブポリシーによって限定される。したがって、一例では、管理者またはアプリケーションはその保持期間（所定のポリシー）がなお有効であるアーカイブオブジェクトを削除することができない。

【 0 0 2 8 】

アクセスレイヤ 3 0 4 は、好ましくはさらにウェブユーザインタフェース（UI）3 1 8 及び SNMP ゲートウェイ 3 2 0 を含む。フィルタランザクション及び管理レイヤ 3 0 6 内の管理エンジン 3 2 2 への対話型アクセスを提供する管理者コンソールとして、ウェブユーザインタフェース 3 1 8 が好ましくは実現される。管理コンソール 3 1 8 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ビューを提供するパスワードで保護されたウェブベースの GUI である。SNMP ゲートウェイ 3 2 0 は管理エンジン 3 2 2 への記憶管理アプリケーションの簡単なアクセスを提供し、それらがクラスタ活動をしっかりと監視して制御することを可能にする。管理エンジンはシステ

10

20

30

40

50

ム及びポリシーイベントを含むクラスタ活動を監視する。ファイルランザクション及び管理レイヤ306は、さらにリクエストマネージャプロセス324を含む。リクエストマネージャ324は、(アクセスレイヤ304を通して)外界からの全てのリクエスト、同じくコアコンポーネントレイヤ308内のポリシーマネージャ326からの内部リクエストを組織化する。

【0029】

ポリシーマネージャ326に加えて、コアコンポーネントはさらにメタデータマネージャ328及び記憶マネージャ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、好ましくは各ノード上にインストールされる。集合的に、クラスタ内のメタデータマネージャが分散型データベースとして働き、全てのアーカイブオブジェクトを管理する。所定のノード上で、メタデータマネージャ328がアーカイブオブジェクトのサブセットを管理し、そこで好ましくは、各オブジェクトが、外部ファイル(「EF」、記憶用のアーカイブに入れられたデータ)とアーカイブデータが物理的に位置する一組の内部ファイル(各「IF」との間をマップする。同じメタデータマネージャ328が、さらに他のノードから複製された一組のアーカイブオブジェクトを管理する。したがって、すべての外部ファイルの現在の状態がいくつかのノード上の複数のメタデータマネージャに常に利用可能である。ノード故障の場合には、他のノード上のメタデータマネージャが、故障したノードによって以前に管理されていたデータへのアクセスを提供し続ける。記憶マネージャ330は、分散アプリケーション内の全ての他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、それがノードのローカルファイルシステム内にデータオブジェクトを記憶する。所定のノード内の各ドライブが、好ましくはそれ自体の記憶マネージャを有する。これは、ノードが個別ドライブを除去してスループットを最適化することを可能にする。記憶マネージャ330は、さらにシステム情報、データに関する完全性チェック及びローカル構造を直接横断する能力を提供する。

【0030】

さらに、図3内に例示されるように、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を通し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の間の通信を可能にする効率的なかつ信頼性が高いメッセージベースのミドルウェアレイヤである。例示の実施態様では、レイヤは、マルチキャスト及び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、全てのノードを企業サーバに接続する分散型のネームサービスを実行する。好ましくは、DNSマネージャ(単独でまたはDNSサービスと共に)負荷バランスが、全てのノードにわたって最大クラスタスループット及び利用可能度を確実にすることを要求する。

【0031】

例示の実施態様では、HCP(日立コンテンツプラットフォーム)アプリケーションインスタンスのようなアプリケーションが、Red Hat Linux 9.0、フェドラーコア6、などのような、基本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、任意の都合がいい分散通信機構である。他のコンポーネントがFUSE(ユーザスペース内のファイルシステム)を含むことができ、それが固定コンテンツファイルシステム(FCFS)316に対して使用されることができる。NFSゲートウェイ310は、標準nfsd LinuxカーネルNFSドライバによって実現されることができる。各ノード内のデータベースは、例えばPostgreSQL(また、本願明細書でPostgresと称される)を実現されることができる、それはオブジェクトリレーショナルデータベース管理システム(ORDBMS)である。ノードは、Java HTTPサーバ及びサブレットコンテナであるJettyのようなウェブサーバを含むことができる。もちろん、上記の機構は単に例証となるだけである。

【0032】

所定のノード上の記憶マネージャ330は、物理記憶装置を管理する役割を果たす。好ましくは、各記憶マネージャインスタンスが、全てのファイルがそのプレースメントアルゴリズムに従って配置される単一ルートディレクトリの原因となる。複数の記憶マネージャ

10

20

30

40

50

ャインスタンスが同時にノード上で実行することができ、各々が通常、システム内の異なる物理ディスクを代表する。記憶マネージャは、残りのシステムから使用されるドライブ及びインタフェース技術を分離する。記憶マネージャインスタンスがファイルを書き込むよう依頼される時、それはそれが原因である表現に対するフルパス及びファイルネームを生成する。代表実施態様において、記憶マネージャ上に記憶されるべき各オブジェクトは、次いで異なるタイプの情報を得続けるためにデータを記憶するにつれてファイルにそれ自体のメタデータを追加する記憶マネージャによって記憶されるべき生データとして収容される。例証として、このメタデータは以下を含む：EF長（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作成タイムスタンプ（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込みの時間（PUT）での内部ファイルのシグネチャ）及びEFファイルネーム（外部ファイルファイルネーム）。内部ファイルデータと共にこの追加的なメタデータを記憶することは、追加的なレベルの保護を提供する。特に、スカベンジングは内部ファイル内に記憶されるメタデータからデータベース内に外部ファイルレコードを作り出すことができる。他のポリシーは、内部ファイルが損なわれていないままであることを確認するために内部ファイルに対して内部ファイルハッシュを確認することができる。

【0033】

内部ファイルはアーカイブオブジェクト内の本来の「ファイル」の一部を表すデータの「チャンク」であることができ、及び、それらはストライピング及び保護ブロックを達成する異なるノード上に配置されることができる。より小さいチャンクユニットへの外部ファイルのこの分裂は、要件でないが、しかしながら；代替案では、内部ファイルが外部ファイルの完全コピーであることができる。一般的に、各外部ファイルエントリに対して多くの内部ファイルエントリがあることができる一方、1つの外部ファイルエントリが各アーカイブオブジェクトに対してメタデータマネージャ内にある。一般的に、内部ファイルレイアウトはシステムに依存する。所定の実現において、ディスク上のこのデータの実際の物理フォーマットが一連の可変長レコード内に記憶される。

【0034】

リクエストマネージャ324は、システム内の他のコンポーネントと対話することによってアーカイブアクションを実行するために必要な操作の組を実行する役割を果たす。リクエストマネージャは、異なるタイプの多くの同時アクションをサポートして、任意の失敗したトランザクションをロールバックすることが可能であり、かつ実行するのに長い間かかる可能性があるトランザクションをサポートする。リクエストマネージャは、更にアーカイブ内の読出し/書込み操作が適切に取り扱われることを確実にしてかつ全てのリクエストがいつでも既知の状態にあることを保証する。それはさらに、所定のクライアントリクエストを満足させるためにノードにわたって複数の読出し/書込み操作を調整するためのトランザクション制御を提供する。加えて、リクエストマネージャは最近使用されたファイルに対するメタデータマネージャエントリをキャッシュに置いてかつセッション、同じくデータブロックに対するバッファリングを提供する。

【0035】

クラスタの主要な応答性は、ディスク上に無制限の数のファイルを確実に記憶することである。所定のノードは、それがどんな理由にせよ到達できないかまたはさもなければ利用できなくなるかもしれないという意味で、「信頼できない」と考えられることができる。一まとまりのこの種の潜在的に信頼できないノードが、信頼性が高いかつ高度に利用可能な記憶域を作り出すために共同する。概ね、記憶される必要がある情報の2つのタイプ：ファイルそれ自体及びファイルについてのメタデータがある。固定コンテンツ分散データ記憶の追加的な詳細が、米国特許出願公開第2007/0189153号明細書及び米国特許第7,657,581号明細書内に見つけることができ、それらを本願明細書に引用したものとする。

【0036】

ここで使用しているように、ネームスペースはクラスタの論理パーティションであってかつ少なくとも1つの規定されたアプリケーションに特定の一まとまりのオブジェクトとして基本的に機能する。各ネームスペースは、他のネームスペースに対して私的ファイルシステムを有する。さらに、1つのネームスペースへのアクセスは別のネームスペースへのユーザアクセスを許可しない。ノードのクラスタ/システムは、物理アーカイブインスタンスである。テナントは、ネームスペース（複数スペース）及びおそらく他のサブテナントのグループ化である。クラスタ/システムは、物理アーカイブインスタンスである。同一譲受人の米国特許出願公開第2011/0106802号明細書を参照されたい、それを全体として本願明細書に引用したものとする。

【0037】

10

II . 記憶階層化のためのコンテンツ選択

【0038】

例示的な実施態様によれば、固定コンテンツシステム（FCS）400は、一般的に、図4に示すように固定コンテンツを記憶する複数の記憶媒体ユニット430を有するブロックベースの記憶サブシステム420にネットワーク経由で連結される一群のノード410を有する。好ましい実施態様において、記憶サブシステムは、一群のディスクがRAID群に共に構成されてかつ記憶ニーズを満足させるために個々の論理ディスクユニットに切り分けられることを可能にするバックエンドディスク配列である。より高度な記憶サブシステムが、どのディスクまたはディスクの群が、ディスクのパワーアップ及びダウンを可能にするためにRAID群に制御されるパワーを有することを可能にされるかを示す能力を有する。

20

【0039】

本発明の実施態様は、稼働ユニット（RU）及びスピンドウンユニット（SDU）と呼ばれる2つのクラスの論理ディスクユニットを作り出すためにこれらの記憶サブシステム特徴を利用する。パワーが常にディスクにあり、かつデータアクセスに利用可能なように、ディスクスピンドウン機能を可能にするように構成されないRAID群上に、RUが収容される。ディスクがパワーダウンされることができ、及び、それゆえに、ディスクが再起動されてかつスピンドウンされるまでそれらのディスク上のデータが直ちに利用可能でないように、ディスクスピンドウンを可能にするように構成されるRAID群上に、SDUが収容される。

30

【0040】

RU及びSDUがどのように使用されるかを規定するために、記憶階層化ルール（STR）が確立される。適切なSTRは、FCS上に記憶される予想される使用及びデータのライフサイクルから決定される。STRは、どんなコンテンツがRUまたはSDU上に記憶されるべき資格があるかを示す。一実施態様によれば、FCS上に記憶される全てのコンテンツまたは一群のコンテンツのどちらかに適用される規定される3つのSTRがある。第1のSTRは、「けっしてない」であり、RU上にまたはスピンドウンされないSDU指示された記憶域上に記憶されるべき全てのデータを指す。第2のSTRは、「保護コピーだけ」であってかつ保護コピーを指す。固定コンテンツシステムでは、コンテンツはコンテンツの複数コピーを記憶することによって保護されることができる。一般的に、バックアップコピーまたは複数コピーはまずアクセスされない。これらのバックアップコピーは、保護コピーとしてのSDUの候補である。第3のSTRは、「取込後X時間」であってかつSDU上に記憶されるべき候補になるのに十分な一定の時間の間FCS内に存在したコンテンツを指す。

40

【0041】

割り当てられるSTRに関係なく、RUが設定可能な使用されたスペース/記憶容量しきい値に到達した時にだけ、SDUに資格のあるコンテンツを移行するように、FCSがさらに構成されることができる。任意の資格のあるコンテンツが、追加的なコストを伴わずに資格のないコンテンツと共にRU上にとどまることができる；さらに、それをRU上に残すより、資格のあるコンテンツを移行するためにSDUをパワーアップすることはあ

50

まり有益でないかもしれない。

【 0 0 4 2 】

S T R 構成とともに R U 及び S D U を提供することによって、F C S は記憶階層化サービス (S T S) を使用して F C S のコンテンツ上で操作することが可能である。S T S は、定期的に行うか以下活動の原因となる：(1) システムニーズに基づいて S D U の状態をパワーアップまたはパワーダウンの適切な状態に管理する；及び(2) S T R 構成、F C S の状態及び記憶されたコンテンツの適格性に基づいて R U と S D U との間でコンテンツを移行する。記憶媒体ユニット (稼働ユニット及びスピンダウンユニット) の状態を含む、F C S の状態が監視される。特定の実施態様において、S T S が F C S 内の各ノード内に提供される。

10

【 0 0 4 3 】

図 5 は、S T S モジュール 3 4 0、ノードに対する記憶媒体ユニット及びそれらの状態 (例えば容量) をリストする記憶媒体状態テーブル 3 5 0 及び記憶マネージャ 3 3 0 の複数のインスタンス (図 3 を参照のこと) を有するノード (図 3 の 3 0 0 または図 4 の 4 1 0) に関して略図で例示する簡略図である。例えば、記憶マネージャ 3 3 0 の各インスタンスが、全てのファイルがそのプレースメントアルゴリズムに従って配置される単一ルートディレクトリの原因となる。各記憶マネージャインスタンスは、通常 F C S 内の異なる論理ボリュームを表す。図 5 は、記憶マネージャ 3 3 0 のインスタンス (一般的にネットワークを通しての接続) 経由でノード 3 0 0 と記憶ユニット 4 3 0 (R U 及び S D U) との間の接続を示す。好ましい一実施態様では、全てのノードが、全てのノードに対する記憶媒体ユニット及びそれらの状態をリストする同じ記憶媒体状態テーブル 3 5 0 を記憶する。記憶媒体ユニットの状態は別々に監視されることができ、及びその情報が次いで同じ記憶媒体状態テーブル 3 5 0 を形成するために相互接続されたノードの間で共有される。S T S モジュール 3 4 0 は、以下の論理アクションを有することができる：

20

【 0 0 4 4 】

S D U への移動に適格なコンテンツの M o v e T o S D U リストを構築する。

【 0 0 4 5 】

S D U に適格でない以外の S D U 上のコンテンツの M o v e F r o m S D U リストを構築する。

【 0 0 4 6 】

構成される容量しきい値に到達した R U 上に存在するコンテンツだけを収容するために M o v e T o S D U リストにフィルターをかける。

30

【 0 0 4 7 】

コンテンツが M o v e T o S D U 及び M o v e F r o m S D U リスト内に存在する全ての S D U の S p i n U p S D U リストを構築する。

【 0 0 4 8 】

S p i n U p S D U リスト内の S D U をスピンアップする。

【 0 0 4 9 】

M o v e T o S D U 及び M o v e F r o m S D U リストを処理するために、適切な宛先への移行処理を開始する。

40

【 0 0 5 0 】

図 6 は、記憶階層化のためのコンテンツ選択のプロセスを例示する一例を示す。図 6 では、F C S はクラスタにわたりノード全体に広がる R U 用の 3 台のディスク及び S D U 用の 2 台のディスクを有する。保護コピーが S D U 上に配置されるべきことを示すサービスプランによって、F C S が構成される。コンテンツは F C S 用の 2 のデータ保護レベル (D P L) を有し、及び F C S 上に 2 つのコピーを備えた 3 つのオブジェクト (A、B、C) がある (A 1、A 2、B 1、B 2、C 1、C 2)。第 1 の 2 台の R U (A 1 B 2 C 2 及び B 1 C 1) は、それらの消費閾値より上にある。S T S は、S D U への移動の候補を処理するために F C S を走査する。走査に基づいて、それらが二次保護コピーであってかつ容量しきい値より上にある第 1 の R U 上に存在するので、コピー B 2 及び C 2 が移動され

50

る。コンテンツの一次コピー（B1及びC1）だけがそこに存在するので、何のコンテンツも第2のRUから移動されない。それが存在する第3のRUが容量しきい値より上にならないので、第2のコピーのA2は移動されない。

【0051】

一般に、SDUへ移動されるオブジェクトは、サービスプランがオブジェクトに対して変更され、かつそれがサービスプランの下でもはや基準を満たさない限り、オブジェクトの寿命の間そこにとどまる。それから、オブジェクトがRUへ戻される。

【0052】

本発明の代替実施態様は、さらに種々の信頼性、性能及びコスト特性を有する種々のクラスの記憶域にコンテンツを移行するために使用されることができる。異なるクラスの記憶域は、限定されるものではないが固体ディスク（SSD）、Fibre Channelハードディスクドライブ（FC-HDD）または他のネットワーク記憶ユニットを含むことができる。

10

【0053】

コンテンツ移動に対する判断は、システムによって自動的に生成されるかまたはコンテンツの所有者によって提供されるオブジェクトメタデータに基づく記憶または位置のクラスの高機能な選択を提供するように拡張されることができる。以下は、いくつかの例基準である。第一は、「データのタイプ」である。例えば、怪我が回復したあと、X線画像はまず観察されない、したがって、この種のタイプのデータはスピンドアウンの良い候補であり、一方、全般的な医療記録は多用され、及びこの種のタイプのデータはスピンドアウンの良い候補でないかもしれない。第2の基準は、「最後のアクセスからの時間」である。そのアクセス履歴に基づいてデータを移動することが望ましくなることができる。しばらく（例えば6ヵ月以上）アクセスされなかったコンテンツは、スピンドアウンの良い候補である、一方よりしばしばアクセスされるコンテンツが装置を動かすために移動されることができる。第3の基準は、「取込からの時間」である。いくつかのタイプのデータを取込ですぐに（例えば、データがバックアップだけのためにある時）または取込の6ヵ月後にスピンドアウンに移すことなどは意味をなすことができる。第4の基準は、「オブジェクトのバージョン」である。例えば、最新版以外のオブジェクトの全てのバージョンはスピンドアウン上にあるべきである。

20

【0054】

長期目標は、クラスタ管理者がテナント管理者に利用可能にするかまたは販売することができる異なるサービスプランをFCSがサポートするためにある。サービスプランは、記憶クラス（例えば常にスピンドアアップされたディスク及びスピンドアウンされたディスク）ならびにいつデータが異なるクラスの記憶域上で記憶されるかの一組のルールの組合せである記憶階層化戦略を有する。例えば、クラスタ管理者はその戦略内のデータが常にスピンドアアップされたディスク上にあるプレミアムサービスプランをサポートしたいかもしれない。「標準」サービスプランは6ヵ月の間スピンドアアップされたディスク上にデータを含むことができるが、しかし、6ヵ月後にアクセスされない場合、データはスピンドアウンディスクへ移動される。「アーカイブ」サービスプランが、データアクセスがおそらくなく、かつ、コンテンツへのアクセスに対する遅延が受け入れられるという予想によってコンテンツをスピンドアウンディスクの候補にすぐにするすることができる。最終的に、目標は、クラスタ管理者が、サービスプランを規定してかつテナント管理者にそれらのプランを異なるレートで利用可能にするかまたは販売することが可能であるようにすることである。

30

40

【0055】

図7は、STSモジュール340によって実行される記憶階層化サービス（STS）を例示する流れ図の一例である。ステップ702において、STSモジュール340がどんな固定コンテンツが稼働ユニット上に記憶されるべき資格があるか、及び、どんな固定コンテンツがスピンドアウンユニット上に記憶されるべき資格があるかを示すポリシーを設定する記憶階層化ルールであって、記憶ユニット内の記憶及び記憶ユニット間の移行のためにその記憶されたコンテンツの適格性を判定するために固定コンテンツシステム内の少

50

なくとも一群の固定コンテンツに適用可能な記憶階層化ルールを確立する。これは、一般的に管理者からの入力によって実行される。特定の実施態様において、記憶階層化ルールが予想される使用、ライフサイクル及びコンテンツオブジェクトの経過時間及びコンテンツオブジェクトの1つ以上の冗長バックアップコピーの存在を含む一組の基準に基づいて記憶されたコンテンツの適格性を判定するためにコンテンツオブジェクトを評価するように確立される。さらに、記憶階層化ルールは、資格のあるコンテンツが存在する稼働ユニットが設定可能な消費利用閾値に到達した時にだけ資格のあるコンテンツをスピンドウンユニットに移行するために確立されることができる。

【 0 0 5 6 】

ステップ704では、STSモジュール340が固定コンテンツシステムの状態を監視する。例えば、オブジェクトが存在する稼働ユニットがスピンドウンユニットにオブジェクトを移行するパーミッションを示す設定可能な消費利用閾値に到達したかどうか、及び稼働ユニットからオブジェクトを収容するためのスピンドウンユニット上で利用可能なスペースがあるかどうかを判定することをこれが含むことができる。

10

【 0 0 5 7 】

ステップ706では、STSモジュール340が記憶階層化ルール、固定コンテンツシステムの状態及びこの少なくとも一群の固定コンテンツの記憶されたコンテンツの適格性に基づいて記憶媒体ユニット間を（例えば、稼働ユニットとスピンドウンユニットの間を）移行するコンテンツオブジェクトの候補を識別する。

20

【 0 0 5 8 】

ステップ708では、移行モジュール342（図5を参照のこと）が識別に基づいてコンテンツオブジェクトを移行する。ステップ710では、記憶ユニット状態管理モジュール344（図5を参照のこと）が、コンテンツオブジェクトを記憶してかつ監視、識別及び移行の結果として、コンテンツオブジェクトを移行するためのシステムニーズに基づいて、スピンドウンユニットの各々の状態をパワーアップまたはパワーダウンの適切な状態に管理する。

【 0 0 5 9 】

特定の実施態様では、記憶ユニットは異なる信頼性、性能またはコスト特性の少なくとも1つを有する異なるクラスの記憶ユニットを含む。移行モジュール342は、コンテンツオブジェクトのオブジェクトメタデータ基準であって、「データのタイプ」、「最後のアクセスからの時間」、「取込からの時間」及び「オブジェクトのバージョン」の1つ以上を含むオブジェクトメタデータ基準に基づいて異なるクラスの稼働ユニットとスピンドウンユニットとの間でコンテンツオブジェクトを移動するように構成される。

30

【 0 0 6 0 】

特定の実施態様に従って図7のSTSプロセスを開始するために、サービスプラン選択モジュール346（図5を参照のこと）が、記憶サブシステム内の稼働ユニット及びスピンドウンユニットを含む記憶ユニットの異なる記憶クラスの組合せを利用する記憶階層化戦略、及び記憶ユニットの異なる記憶クラス内のコンテンツオブジェクトの記憶のための記憶されたコンテンツの適格性を判定するポリシーを設定する記憶階層化ルールを各々特定する複数のサービスプランからサービスプランを選ぶためにグラフィカルユーザインタフェースを提供するように構成される。

40

【 0 0 6 1 】

図8は、ユーザ/管理者が、サービスプランを作り出す、及び/または割り当てることを可能にする表示の形の、この種のユーザインタフェース800の一例を示す。ここに示した例では、図8aはサービスプランを作り出して階層化ポリシーを選ぶためのスクリーンショットを示す。規定されたポリシーを閲覧して編集するオプションが、スクリーンの底部にある。図8bでは、スクリーンショットが、作り出されるべきネームスペースに対するサービスプランを選ぶ時何が起こるかを示す。図8cでは、スクリーンショットが既存のネームスペース上のサービスプランを変更する方法を示す。

【 0 0 6 2 】

50

もちろん、図1内に例示されるシステム構成は、本発明が実現されることができるコンテンツプラットフォームまたは複製オブジェクト記憶システムの純粹に例示的なものであり、及び、本発明は特定のハードウェア構成に限定されない。本発明を実現するコンピュータ及び記憶システムは、さらに上記の本発明を実現するために使用されるモジュール、プログラム及びデータ構造を記憶して読み出すことができる周知の入出力装置（例えばCD及びDVDドライブ、フロッピーディスクドライブ、ハードディスク、その他）を有することができる。これらのモジュール、プログラム及びデータ構造は、この種のコンピュータ可読の媒体上でコード化されることができる。例えば、本発明のデータ構造が、本発明内に使用されるプログラムがその上に存在する1つ以上のコンピュータ可読の媒体から独立にコンピュータ可読の媒体上に記憶されることができる。システムのコンポーネントは、デジタルデータ通信の任意の形式または媒体、例えば通信ネットワークによって相互接続されることができる。通信ネットワークの例は、ローカルエリアネットワーク、広域ネットワーク、例えばインターネット、無線ネットワーク、記憶領域ネットワーク、などを含む。

10

【0063】

記述では、数多くの詳細が本発明の徹底的な理解を提供するために説明のために記載されている。しかしながら、これらの具体的な詳細の全てが本発明を実践するために必要とはされないことが当業者にとって明らかである。さらに注意されることは、本発明が通常フローチャート、流れ図、構造線図またはブロック図として表されるプロセスとして記述されることができることである。フローチャートがシーケンシャルプロセスとして操作を記述することができるとはいえ、操作の多くは並列にまたは同時に実行されることができる。加えて、操作の順序は再配置されることができる。

20

【0064】

周知のように、上記した操作はハードウェア、ソフトウェアまたはソフトウェア及びハードウェアのいくつかの組合せによって実行されることができる。本発明の実施態様の種々の態様が、回路及び論理装置（ハードウェア）を使用して実現されることができ、一方他の態様が、プロセッサによって実行されると、本発明の実施態様を実施する方法をプロセッサに実行させる機械可読の媒体上に記憶される命令（ソフトウェア）を使用して実現されることができる。さらに、本発明のいくつかの実施態様はハードウェアでもっぱら実行されることができるが、一方、他の実施態様はソフトウェアでもっぱら実行されることができる。さらに、記述される種々の機能が、単一ユニット内に実行されることができるかまたは任意の数の方法で複数のコンポーネントにわたって広げられることができる。ソフトウェアによって実行される時、この方法は、コンピュータ可読媒体上に記憶される命令に基づいて汎用コンピュータのようなプロセッサによって実行されることができる。必要に応じて、命令は圧縮された及び/または暗号化された形態で媒体上に記憶されることができる。

30

【0065】

上記のことから、本発明が、記憶階層化のためのコンテンツ選択に基づいて固定コンテンツシステム内の固定コンテンツの記憶を管理するための方法、器具及びコンピュータ可読の媒体上に記憶されるプログラムを提供することが明白である。加えて、特定の実施態様がこの明細書内に例示されて記述されたとはいえ、当業者は、同じ目的を達成するために算出される任意の配置が開示される特定の実施態様と置換されることができると認識する。この開示は本発明のありとあらゆる適応または変形を包含することを意図され、及び、以下の請求項内に使用される用語が、本発明を明細書内に開示される特定の実施態様に限定するように解釈されるべきでないことが理解されるべきである。むしろ、本発明の範囲は以下の請求項によって完全に決定されるべきであり、この種の請求項が権利を与えられる等価物の完全な範囲とともに、それが請求項解釈の確立された原則に従って解釈されるべきである。

40

【0066】

「先行技術文献のリスト」

50

「特許文献」

「PTL1」米国特許第8,006,111号明細書

「PTL2」米国特許第7,155,466号明細書

「PTL3」米国特許出願公開第2007/0189153号明細書

「PTL4」米国特許第7,657,581号明細書

「PTL5」米国特許出願公開第2011/0106802号明細書

【図1】

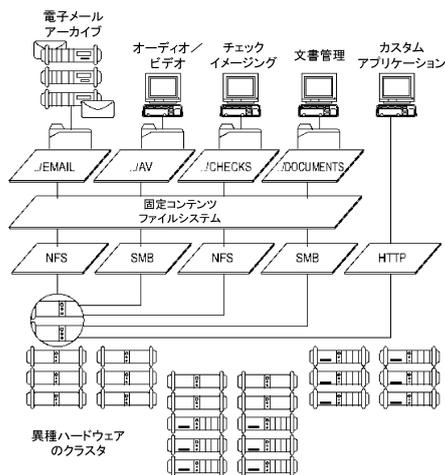


FIG. 1

【図2】

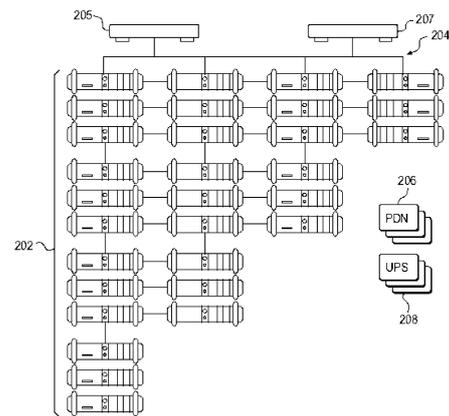


FIG. 2

【 図 3 】

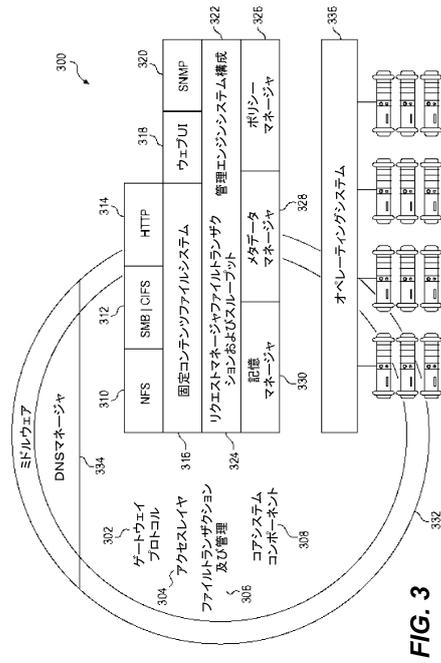


FIG. 3

【 図 4 】

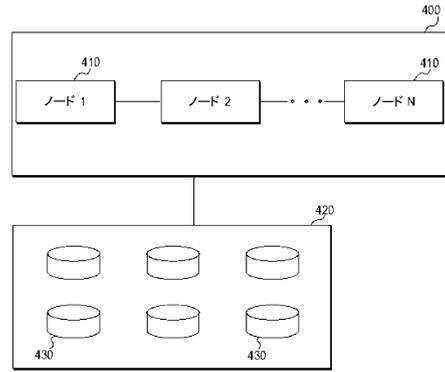


FIG. 4

【 図 5 】

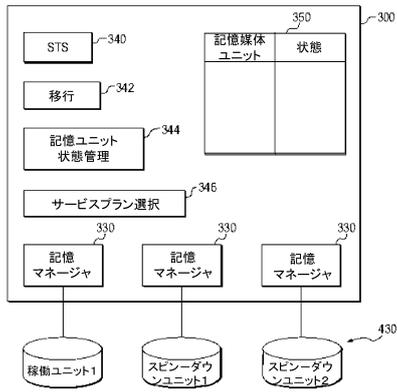


FIG. 5

【 図 6 】

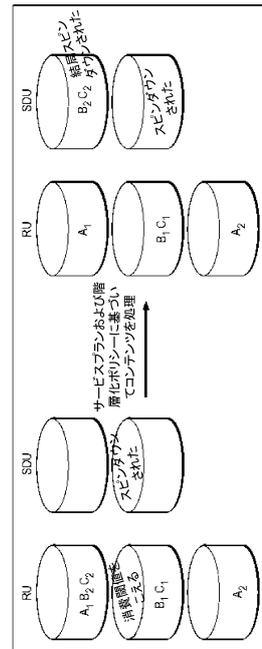


FIG. 6

【図7】

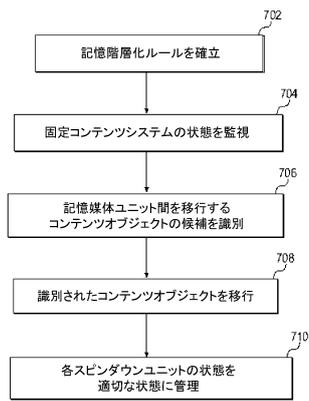


FIG. 7

【 図 8 A 】

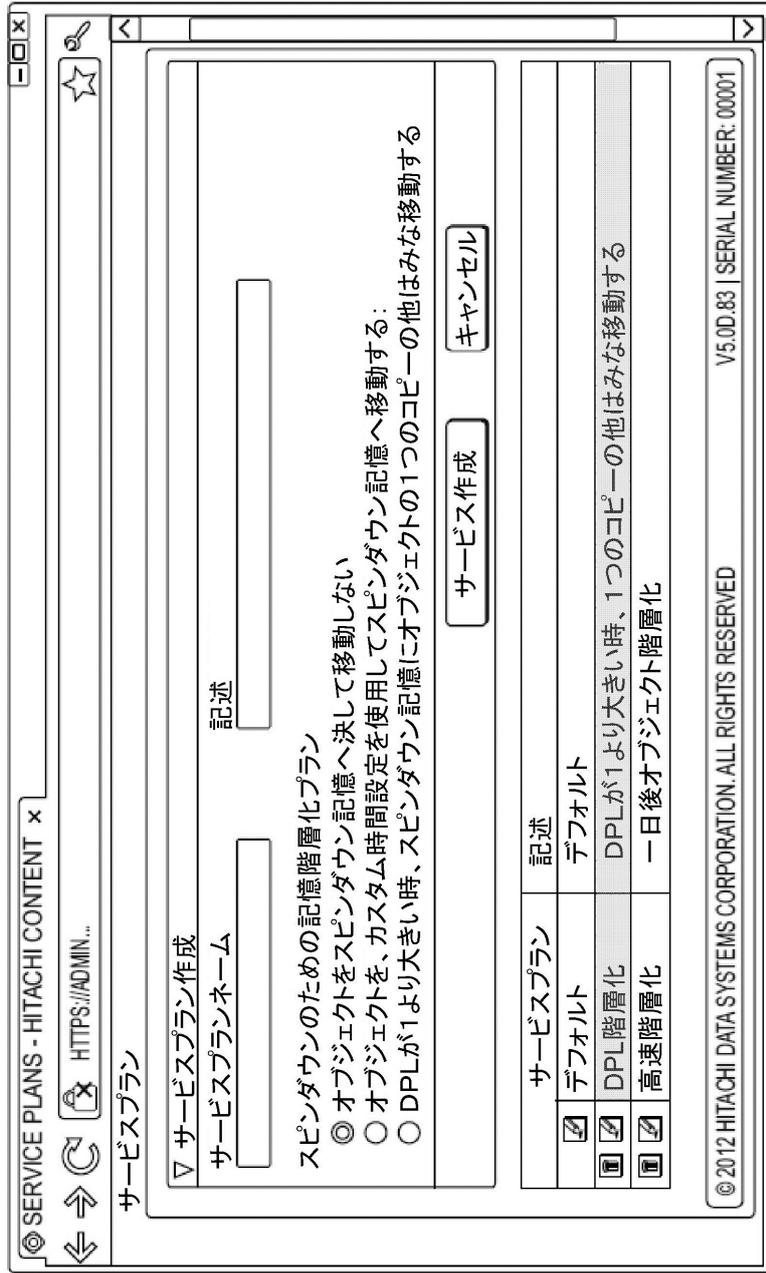


FIG. 8A

【 図 8 B 】

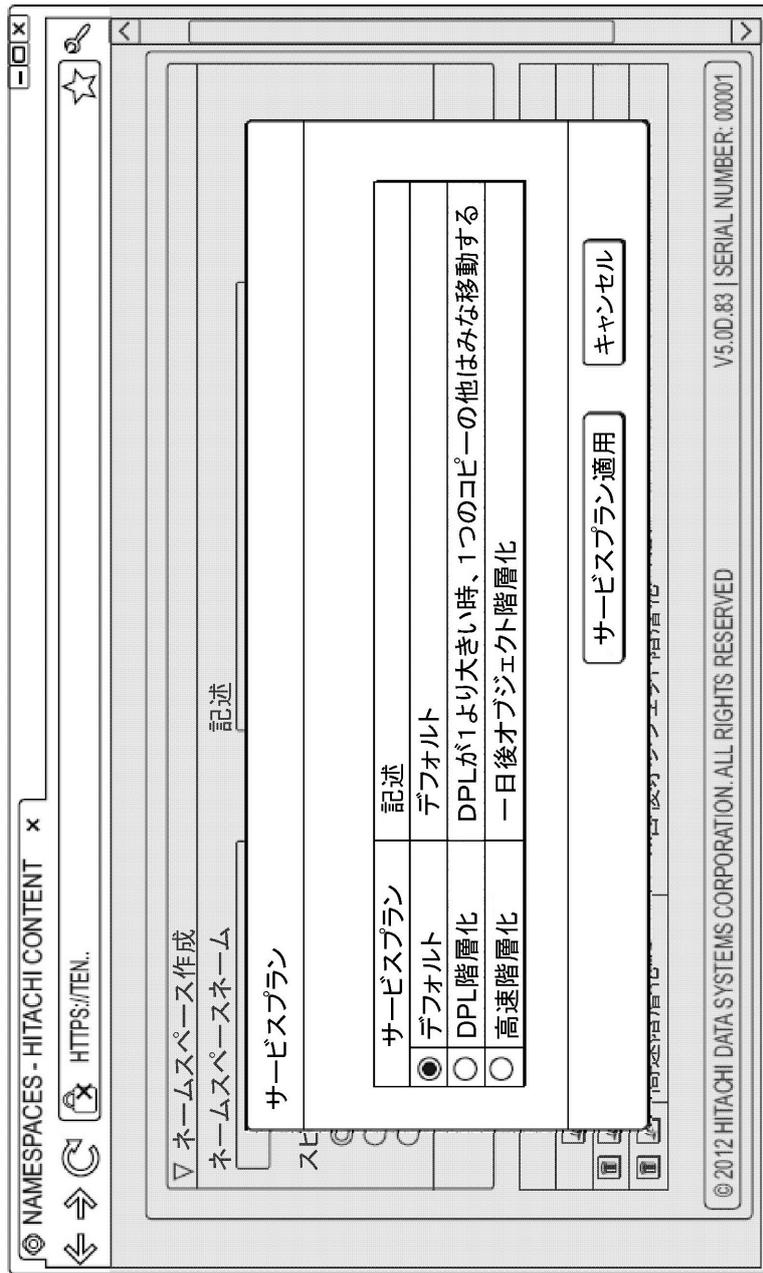


FIG. 8B

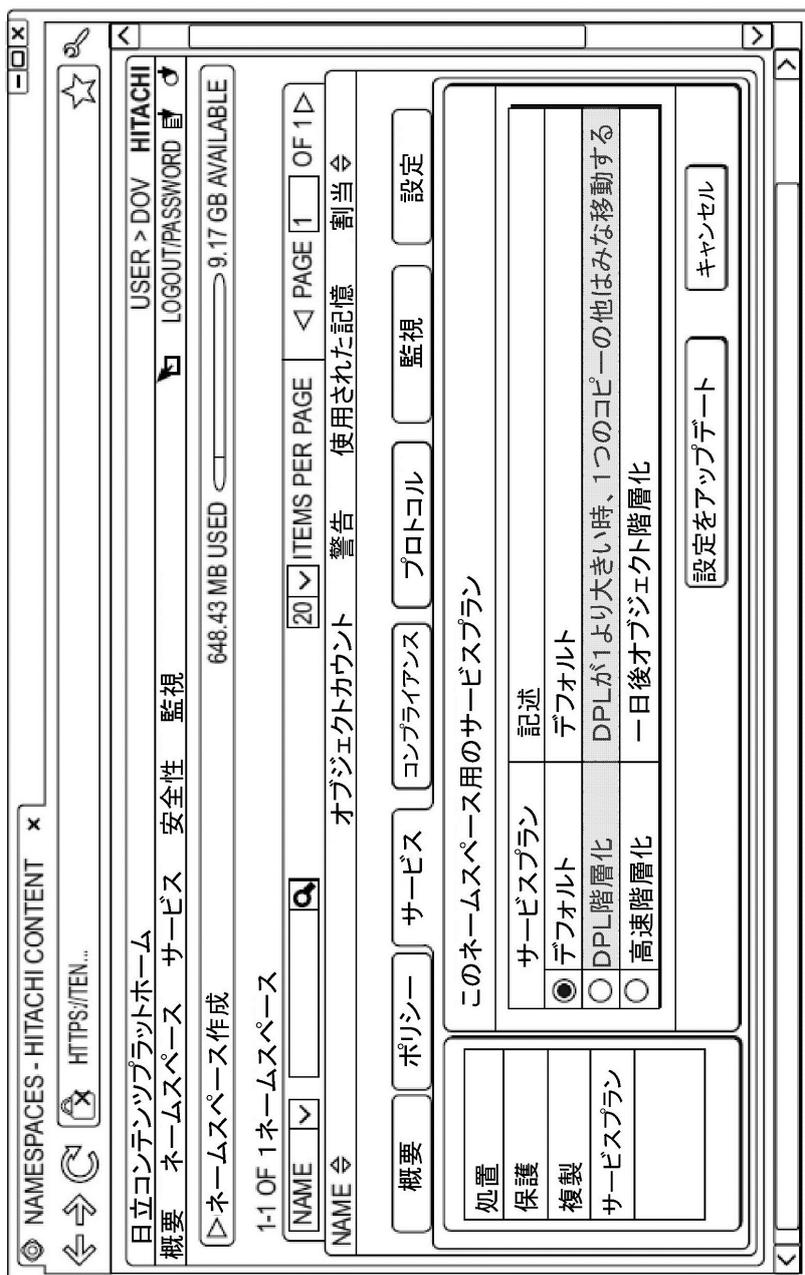


FIG. 8C

フロントページの続き

(51)Int.Cl. F I
G 0 6 F 12/00 5 0 1 B

(72)発明者 ゴロツスキー, ヴィタリー
アメリカ合衆国 0 2 4 7 4 マサチューセッツ州 アーリントン, カレッジ アベニュー 1 0 7

(72)発明者 ブライアント, アラン, ジー.
アメリカ合衆国 0 2 0 3 2 マサチューセッツ州 イースト ウォルポール, ユニオン ストリート 2 2 8

審査官 木村 雅也

(56)参考文献 特開平07 - 262058 (JP, A)
特開2012 - 027933 (JP, A)
特開2011 - 076294 (JP, A)
国際公開第2011 / 108021 (WO, A1)
国際公開第2010 / 131292 (WO, A1)
米国特許出願公開第2004 / 0243761 (US, A1)
米国特許第07567188 (US, B1)

(58)調査した分野(Int.Cl., DB名)
G 0 6 F 3 / 0 6
G 0 6 F 1 2 / 0 0