

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第5274772号
(P5274772)

(45) 発行日 平成25年8月28日 (2013. 8. 28)

(24) 登録日 平成25年5月24日 (2013. 5. 24)

(51) Int. Cl.

F I

G O 6 F 12/00 (2006. 01)

G O 6 F 12/00 5 O 1 A

G O 6 F 3/06 (2006. 01)

G O 6 F 12/00 5 1 7

G O 6 F 3/06 3 O 1 Z

請求項の数 15 (全 32 頁)

(21) 出願番号 特願2006-534180 (P2006-534180)
 (86) (22) 出願日 平成16年9月30日 (2004. 9. 30)
 (65) 公表番号 特表2007-507811 (P2007-507811A)
 (43) 公表日 平成19年3月29日 (2007. 3. 29)
 (86) 国際出願番号 PCT/US2004/032490
 (87) 国際公開番号 W02005/033938
 (87) 国際公開日 平成17年4月14日 (2005. 4. 14)
 審査請求日 平成19年9月19日 (2007. 9. 19)
 審判番号 不服2012-744 (P2012-744/J1)
 審判請求日 平成24年1月13日 (2012. 1. 13)
 (31) 優先権主張番号 10/675, 188
 (32) 優先日 平成15年9月30日 (2003. 9. 30)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 505014122
 シマンテック・オペレーティング・コーポ
 レーション
 アメリカ合衆国・94043・カリフォル
 ニア州・マウンテンビュー・エリス スト
 リート・350
 (74) 代理人 100147485
 弁理士 杉村 憲司
 (74) 代理人 100158148
 弁理士 荒木 淳
 (74) 代理人 100166213
 弁理士 永久保 宅哉

最終頁に続く

(54) 【発明の名称】 データ・ストレージ内の時相データを維持するためのシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

現在データと履歴データとを含む時相データに対する入出力要求を生成するように構成したアプリケーションと、

前記時相データを物理ストレージに格納するように構成され、前記物理ストレージにマップされた論理デバイスである時相ボリュームと、

該時相データに対して、前記入出力要求が示すオペレーションを時相ボリュームレベルで実行するための時相ボリューム・マネージャと

を備え、

前記時相ボリューム・マネージャは、さらに、

前記入出力要求によって指定されたタイムスタンプを前記時相データに割り当てるための複数のチェックポイント・オペレーションを行い、

生成されるべきスライス・イン・タイム・イメージの要求を指定し、且つ2つのタイムスタンプを指定するスライス・イン・タイム入出力要求を前記アプリケーションから受けとり、ここで、該スライス・イン・タイム・イメージは、前記時相データの一部であって、該2つのタイムスタンプの間の時間期間によって定められる部分に対応するものであり、

前記スライス・イン・タイム入出力要求に応答して、前記スライス・イン・タイム・イメージを生成する

ように構成されているシステム。

10

20

【請求項 2】

前記入出力要求の一つはタイムスタンプを指定する書込み要求であり、前記時相ボリューム・マネージャは、該書込み要求に応答して、その書込み要求によって指定された時相データを、前記時相ボリュームに書き込むとともにその書き込んだ時相データに前記タイムスタンプを付すようにさらに構成されてい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】

前記入出力要求の一つはタイムスタンプを指定する読取り要求であり、前記時相ボリューム・マネージャは、該読取り要求に応答して、その読取り要求のタイムスタンプによって指定された時相データを、前記アプリケーションに返却するようにさらに構成されてい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】

前記入出力要求の一つが、時相ボリュームのポイント・イン・タイム・イメージを作成するよう指定するとともにそのポイント・イン・タイム・イメージの作成時点を示すタイムスタンプを指定し、前記時相ボリューム・マネージャは、その作成時点のポイント・イン・タイム・イメージを生成するようにさらに構成されていることを特徴とする請求項 1 に記載のシステム。

【請求項 5】

前記入出力要求が、前記入出力要求によって指定されたタイムスタンプに対応する時点で前記時相ボリュームの履歴データを切り捨てるよう指定し、前記時相ボリューム・マネージャは、前記入出力要求によって指定された時点で前記時相ボリュームの履歴データを切り捨てるようにさらに構成されていることを特徴とする請求項 1 に記載のシステム。

20

【請求項 6】

前記入出力要求が、その入出力要求によって指定された時点まで前記時相ボリュームの特定の部分の履歴データをたどらせるよう指定し、前記時相ボリューム・マネージャは、前記時相ボリュームの当該特定の部分に対して、当該入出力要求が示すオペレーションを実施するようにさらに構成されていることを特徴とする請求項 1 に記載のシステム。

【請求項 7】

前記時相ボリュームは複数の領域に分割されており、

前記入出力要求が、該入出力要求が示すオペレーションが実施される 1 つまたは複数の時相ボリュームの領域を指示す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】

時相ボリューム・マネージャが、ハードウェア処理装置によって、

物理ストレージに格納された時相ボリュームに記憶された時相データにタイムスタンプを割り当てるための複数のチェックポインティング・オペレーションを行うステップであって、前記時相データは現在データと履歴データとを含み、前記時相ボリュームは前記物理ストレージにマップされた論理デバイスであるステップと、

前記チェックポインティング・オペレーションを行うステップ後に、スライス・イン・タイム入出力要求をアプリケーションから受けとるステップであって、前記スライス・イン・タイム入出力要求は、生成されるべきスライス・イン・タイム・イメージの要求を指定し、且つ 2 つのタイムスタンプを指定し、前記スライス・イン・タイム・イメージは、前記時相データの一部であって、該 2 つのタイムスタンプの間の時間期間によって定められる部分に対応するものであるステップと、

40

前記スライス・イン・タイム入出力要求を受けとるステップ後に、前記スライス・イン・タイム入出力要求に応答して、前記スライス・イン・タイム・イメージを生成するステップと

を含み、

前記時相ボリューム・マネージャは、前記ハードウェア処理装置によって、前記時相データに対するオペレーションを時相ボリュームレベルで実行し、

前記オペレーションは、前記時相ボリューム・マネージャが前記アプリケーションから

50

受け取った入出力要求によって指定され、

前記入出力要求は、前記チェックポイント・オペレーションのための前記タイムスタンプを指定すること

を特徴とする方法。

【請求項 9】

前記入出力要求の一つがタイムスタンプを指定する書込み要求であり、前記時相ボリューム・マネージャが、前記書込み要求に応答して、その書込み要求によって指定された時相データを前記時相ボリュームに書き込むとともにその書き込んだ時相データに前記タイムスタンプを付すことを特徴とする請求項 8 に記載の方法。

【請求項 10】

前記入出力要求の一つがタイムスタンプを指定する読取り要求であり、前記時相ボリューム・マネージャが、前記読取り要求に応答して、その読取り要求のタイムスタンプによって指定された時相データを前記アプリケーションに返却することを特徴とする請求項 8 に記載の方法。

【請求項 11】

前記入出力要求の一つは、ポイント・イン・タイム・イメージの作成時点を示すタイムスタンプを指定するとともにそのタイムスタンプに対応する時相ボリュームの部分のポイント・イン・タイム・イメージを作成するよう指定し、前記時相ボリューム・マネージャが該指定されたタイムスタンプに対応する時相ボリュームのポイント・イン・タイム・イメージを生成することを特徴とする請求項 8 に記載の方法。

【請求項 12】

前記入出力要求の一つは、該入出力要求によって指定されたタイムスタンプに対応する時点で前記時相ボリュームの履歴データを切り捨てるよう指定し、前記時相ボリューム・マネージャは、前記入出力要求によって指定された前記時点で前記時相ボリュームの履歴データを切り捨てることを特徴とする請求項 8 に記載の方法。

【請求項 13】

前記入出力要求の一つは、該入出力要求によって指定された時点まで前記時相ボリュームの特定の部分の履歴データをたどらせるよう指定し、前記時相ボリューム・マネージャは、前記時相ボリュームの当該特定の部分に対して、当該入出力要求が示すオペレーションを実施することを特徴とする請求項 8 に記載の方法。

【請求項 14】

前記時相ボリュームは複数の領域に分割されており、

前記入出力要求の一つが、該入出力要求が示すオペレーションが実施される 1 つまたは複数の時相ボリュームの領域を指示する

ことを特徴とする請求項 8 に記載の方法。

【請求項 15】

ハードウェア処理装置に、前記請求項 8 乃至 14 のいずれかの方法を実行させるためのプログラム命令を具備したコンピュータ可読記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、コンピュータ・システムの分野に関し、より詳細には、データ・ストレージ・システムに関する。

【背景技術】

【0002】

典型的な従来技術のデータ・ボリュームでは、データに変更があった場合には、その場で対応するブロックが変更される。これは、今日のリレーショナル・データベースのようなアプリケーションやファイル・システムに関する限りは適切である。しかし、従来のデータ・ボリューム管理技術は、時相情報すなわち径時変化情報を維持し処理することには不向きである。データ・マイニング、データ・ウェアハウジング、メディア・ライブラリ

10

20

30

40

50

、医療レコードなどのアプリケーションは、時相情報すなわち径時変化情報の維持／または処理が必要なこともある。従来技術では、これらの分野のアプリケーションが時相データを管理できる単一のインフラストラクチャが存在しない。ストレージのインフラストラクチャが欠けている故に、従来技術のアプリケーションは通常、時相パラダイムをサポートするのに非効率となりがちなブルート・フォース型の方法を使用する。したがって、様々なアプリケーションをまたいで時相データを一般的に管理する、論理デバイス・レベルのインフラストラクチャを提供することが望ましい。

【発明の開示】

【課題を解決するための手段】

【0003】

データ・ストレージ内の時相データを論理デバイス・レベルで維持するためのシステムと方法に関する諸実施態様を説明する。これらの諸実施態様は、コンテンツを保存するために、論理デバイス（ボリューム）レベルでデータにタイムスタンプを付けるための一般メカニズムを備えている。これらの諸実施態様は、時相ボリュームを管理しそれにアクセスするためのメカニズムを備えている。いくつかの実施態様をボリューム・マネージャと統合することができる。時相ボリュームを管理しそれにアクセスするためのメカニズムを時相ボリューム・マネージャと呼ぶ。時相ボリューム・マネージャの諸実施態様は、アプリケーションおよび／またはアプリケーション・エージェントが時相ボリューム・マネージャと通信して、1つまたは複数の時相ボリューム上の時相情報を管理し追跡するのを可能にするインターフェースを設けている。これらの諸実施態様は、時相データベース、バ

10

20

【0004】

時相ボリュームは、現在データに加えて非現在データも維持するボリュームである。時相ボリュームは、それに格納されているデータの履歴を維持し、過去の任意の時点におけるデータのコピーを検索するための方法をアプリケーションに提供することができる。時相ボリュームでは、データのブロックが変更されると常に、まず既存のブロックが保存され、次いで新しいデータが上書きされる。古いバージョンのブロックは、そのブロックがアプリケーションによってデータから削除された場合にも、維持される。これには、過去におけるデータの1つまたは複数の状態のコピーを維持する効果がある。時相ボリュームを、たとえば、ホスト・ベース環境、ネットワーク・ベース環境（スイッチまたはアプライアンス）、アレイ・ストレージ環境で 사용할 ことができる。さらに、時相ボリュームを、バンド内、バンド外仮想化において使用することができることに留意されたい。

30

【0005】

一実施態様では、アプリケーションは、時相ボリュームを（デバイスとして）直接使用するのを望まない場合は、アプリケーション・エージェント（以下で説明する）を使用して、それらがボリュームの時間特性を利用するのに役立てることができる。これにより、時相データの管理の負担がアプリケーションから取り除かれ、したがって、アプリケーションは、データをどのように格納するかについてではなく、データをどのように消費する

40

【0006】

これらの諸実施態様では、それだけに限らないが、入出力制御型チェックポイントイング、アプリケーション制御型チェックポイントイング、周期型チェックポイントイングを含めたいくつかの方法で、時相ボリューム内のデータの履歴を維持することができる。入出力制御型チェックポイントイングでは、時相ボリュームを使用するアプリケーションは、時相ボリュームへの書込みおよび／または時相ボリュームの読取りを行うときに、タイムスタンプを生成する。アプリケーションは、その必要がある場合またはそれが望まれる場合、書込みの度にその都度タイムスタンプを生成する代わりに、UNIX（登録商標）環境におけるIOCTL（I/O control command）など、ある領域（

50

または全体のボリューム)を対象とするタイムスタンプを指定する、入出力要求またはコマンドを発行することができる。この方法を、アプリケーション制御型チェックポイントイングと呼ぶこともある。別の代替態様は、時相ボリューム・マネージャが周期的に、たとえば10秒毎または10分毎に自動チェックポイントイングを行うためのものである。これを、周期型チェックポイントイングと呼ぶこともある。

【0007】

時相ボリュームは次元としての時間を有する。一実施態様では、時相ボリューム・マネージャは、ユーザが時相ボリューム内の時間次元(履歴)をトラバースするのを可能にするインターフェイスを設ける。一実施態様は、それ自体の独立した履歴を有する時相ボリュームの時相イメージを生成するためのメカニズムを備えていることができる。一実施態様は、2つのタイムスタンプ間での時相ボリュームのスライス・イン・タイム・イメージを生成するためのメカニズムを備えていることができる。一実施態様は、時相ボリュームのポイント・イン・タイム時相イメージを生成するためのメカニズムを備えることができる。ポイント・イン・タイム・イメージは、単一時間におけるボリュームのイメージであり、時相ボリュームの1次元イメージと見ることができる。スライス・イン・タイム・イメージも、ポイント・イン・タイム・イメージも、時相イメージのベースとして使用することができる。

【0008】

以下の詳細な説明は、後で簡潔に説明する添付の図面を参照して行う。

【0009】

本発明は、例示のため、いくつかの実施形態や例示的な図について説明してあるが、本明細書に記載の諸実施形態または図に限定されるものではないことを、当業者には理解されるであろう。添付の図面やその詳細な説明は、本発明を開示の特定の形態に限定するものではなく、本発明は、それとは対照的に、添付の特許請求の範囲で定義される本発明の趣旨と範囲内にあるすべての変更形態、均等物、代替形態を含むものであることを理解されたい。本明細書で使用する表題は、単なる構成上の目的にすぎず、本明細書の記載または添付の特許請求の範囲を限定するために使用されるものではない。本願の全体を通して使用する「することができる」という用語は、義務的な意味(すなわち、必須という意味)にではなく、任意的な意味(すなわち、そうする可能性があるという意味)で使用する。同様に、「含む」、「含めて」、「備える」という用語は、「それだけに限らない」ということを意味する。

【発明を実施するための最良の形態】

【0010】

データ・ストレージ内の時相データを論理デバイス・レベルで維持するためのシステムと方法の諸実施形態を説明する。これらの諸実施形態は、コンテンツを保存するために、論理デバイス(ボリューム)レベルでデータにタイムスタンプを付けるための一般メカニズムを提供する。これらの諸実施形態では、本明細書において時相ボリュームと呼ぶこともある時相データ・ボリュームは、管理、入出力オペレーション、レプリケーションのためのインターフェイスを設けていることができ、また、それだけに限らないが、バックアップ、復元、階層型ストレージ管理(Hierarchical Storage Management: HSM)を含めたオペレーションをサポートすることができる。これらの諸実施形態は、時相ボリュームを管理しそれにアクセスするためのメカニズムを提供する。いくつかの実施形態をベリタス社のVERITAS Volume Managerなどのボリューム・マネージャと統合することができる。本明細書では、時相ボリュームを管理しそれにアクセスするためのメカニズムを時相ボリューム・マネージャと呼ぶこともある。時相ボリューム・マネージャの諸実施形態は、アプリケーションが時相ボリューム・マネージャと通信して、1つまたは複数の時相ボリューム上の時相情報を管理し追跡するのを可能にするインターフェイスを設けている。

【0011】

本明細書内で使用されるいくつかの用語を定義する用語解説が、本明細書の最後に設け

10

20

30

40

50

られていることに留意されたい。

【0012】

これらの諸実施形態は、それだけに限らないが、時相データベース（それらが使用している時相モデルは問わない）、バージョン・ファイル・システム／レポジトリ、データ・アーカイブ、ストリーミング・メディアを含めたデータの履歴に働きかけるアプリケーションが、時相データを管理するためのインフラストラクチャを提供する。これらの諸実施形態は、時相データの管理が必要になる将来のアプリケーション用のビルディング・ブロックとしても働くことができる。

【0013】

時相ボリュームは、現在データに加えて非現在データも維持するボリュームである。時相ボリュームは、そこに格納されているデータの履歴を維持することができ、したがって、アプリケーションが過去の任意の時点におけるデータのコピーを検索するための方法を提供することができる。通常のボリュームでは、データに変更がある場合には、その場に対応するデータ・ブロックが変更される。時相ボリュームでは、データのブロックが変更される場合には、まず既存のブロックが保存され、次いで新しいデータが上書きされる。古いバージョンのブロックは、そのブロックがアプリケーションによってデータから削除された場合にも維持される。これには、過去におけるデータの1つまたは複数の状態のコピーを維持する効果がある。このプロセスは、ボリューム上のデータの途切れないバージョンングを行い、それに変更があった場合にはそのボリュームのスナップショットを作成するプロセスとも考えることができる。別の一実施形態では、新しいデータを離れた位置に書き込むことができ、時相ボリューム内のメタデータ（データ領域を指すポインタなど）を操作することができる。

【0014】

時相ボリュームを、それだけに限らないが、ホスト・ベース環境、ネットワーク・ベース環境（スイッチまたはアプライアンス）、アレイ・ストレージ環境を含めたストレージ環境で 사용할 ことができる。また、時相ボリュームを、バンド内やバンド外仮想化において使用することもできることに留意されたい。

【0015】

図1は、一実施形態による、時相ボリュームにおける時相オペレーションを管理する時相ボリューム・マネージャの図である。時相ボリューム・マネージャ100は、アプリケーション106および／またはオペレーティング／ファイル・システム104の代わりに、時相ボリューム102に対する時相データ・オペレーションを管理することができる。アプリケーション106は、時相データベース、バージョンング・ファイル・システム／レポジトリ、データ・アーカイブ、ストリーミング・メディアなど、時相データを管理するためにデータの履歴を維持することに関心のある、どんなタイプのアプリケーションでもよい。時相ボリューム・マネージャ100は、様々な読取り、書込み、または時相ボリューム102に対するその他のオペレーションを実施するための1つまたは複数のAPIを、アプリケーション106および／またはオペレーティング／ファイル・システム104に提供することができる。

【0016】

時相ボリュームは、時間次元に関心のある任意のアプリケーションが使用できる、一般的な時相データのインフラストラクチャと考えることができる。一実施形態では、アプリケーションは、時相ボリュームを（デバイスとして）直接使用するのを望まない場合は、アプリケーション・エージェント（以下で説明する）を使用して、それらがボリュームの時間特性を利用するのに役立てることができる。アプリケーション・エージェントは、アプリケーションがデータをどのように格納するかについてではなく、データをどのように消費するかについて集中することができるように、時相データの管理の負担をアプリケーションから取り除く。たとえば、データベースはテーブルのコピーをどのように格納するかを気にせずに、効率的な問合せ処理に集中ことができ、ファイル・システムはファイルをどのようにバージョン化するかを気にせずに、ファイル・オペレーションに集中

することができ、レポジトリはデータ履歴をアーカイブすることではなく、索引付けをより迅速に行うことに集中することができる。

【0017】

図2は、一実施形態による、アプリケーション・エージェントを使用して、時相ボリューム・マネージャとインターフェイスするアプリケーションの図である。時相ボリューム・マネージャ200は、アプリケーション206の代わりに、時相ボリューム202に対する時相データ・オペレーションを管理することができる。アプリケーション206は、時相データベース、バージョンング・ファイル・システム/レポジトリ、データ・アーカイブ、ストリーミング・メディアなど、時相データを管理するためにデータの履歴を維持することに関心のある、どんなタイプのアプリケーションでもよい。時相ボリューム・マネージャ200は、様々な読取り、書込み、または時相ボリューム202に対するその他のオペレーションを実施するための、1つまたは複数のAPIを提供することができる。アプリケーション・エージェント204は、アプリケーション206に代わって時相マネージャ200と相互作用するブローカであり、アプリケーション206の(一部または全部の)セマンティクスをサポートしている。アプリケーション・エージェント204は、アプリケーション206が時相ボリューム202と容易に相互作用できるようにすることが好ましい。さらに、アプリケーション・エージェント204は、アプリケーション206がデータを生データのままでなく、必要とするフォーマットで抽出することを可能にする。

10

【0018】

20

データの履歴をストレージ(論理デバイス)レベルで格納することによっても、管理が容易になる。これにより、アプリケーション・レベルでストレージを時間的に管理することの労力が軽減される。1つの例として、データベースは、Bツリーなどそれ自体のメタデータにおけるブランチのコピーを作成する必要がなく、ファイル・システムは、それ自体のiノードのコピーを維持する必要がなく、レポジトリは、それ自体の索引を複製する必要がない。このことは、システム管理者にとっては、異なる時相データベース、異なるバージョンング・ファイル・システム、異なるアーカイブ・レポジトリを管理するのではなく、1つの時相データソースだけを管理することを意味する。

【0019】

時相ボリュームは、時相情報を維持しそれにアクセスするためのインフラストラクチャを提供する。時相ボリュームは、ファイル・システムやデータベースを含む、すべてのレベルのアプリケーションが使用できる。さらに、時相ボリュームは、ファイル・システムやバックアップ製品(たとえば、VERITAS File System、およびVERITAS Net Backup製品)と統合することによって、データのアーカイブ、バージョンング、レプリケーション、バックアップ、HSMのためのビルディング・ブロックとしても使用することができる。時相ボリュームは、後のある時点でのスナップショット、増分バックアップ、レプリケーション、破損したボリュームまたは削除されたファイルの復元用にそれを使用できるように、時相コンテンツを保存しておく。

30

【0020】

時相ボリュームは、HSMなどのプロセスを使用して、データのオフライン・ストレージへの移行/オフライン・ストレージからの呼出しを自動化することにより、仮想的に無制限のストレージ容量をアプリケーションに提供することができる。増分バックアップを含めた自動バックアップを行うために、時相ボリュームを、VERITAS NetBackupなどのバックアップ・ユーティリティと統合することもできる。時相ウィンドウおよび/または周期型チェックポイントイングを使用して、時相ボリュームを対象とする周期的なレプリケーションを行うことができる。

40

【0021】

時相ボリュームは、それだけに限らないが、次のうちの1つまたは複数を含む分野のアプリケーションで使用することができる。

- ・ ビジネス・インテリジェンス/データ・マイニング：これは、顧客データのうしか

50

ら傾向と習性を発見することに関するものである。通常、かかる結果を発見するには、時相データが必要となる。時相データを使用すると、たとえば、小売業者は、自社の顧客の使用傾向を発見するためにデータ・マイニングを行うことができ、銀行は、クレジット履歴、詐欺などを発見するために顧客データを分析することができる。

- ・ データ・ウェアハウス：データ・ウェアハウスは、特定の企業に関するすべての情報を格納するために一般に使用される企業規模のデータベースである。データ・ウェアハウスは、データ・マイニング・ツールによって収集される有益な情報のリポジトリであり、知識発見に使用される。これらは、様々な企業のバックエンドとしても働くことができる。時間次元を格納し、それをデータの問合せの際に利用することができるのは、このデータ・ウェアハウス内であり、あらゆるデータ・ウェアハウスは、それ自体を時相データベースにする時間次元を有しており、したがって、本明細書に記載の時相ボリュームに適している。

10

- ・ マルチメディアとイメージング：ストリーミング・オーディオ/ビデオ、画像、電子書籍、さらにX線やMRIのような医療データは、ある種のバージョニングまたはストリーミング（経時変化コンポーネント）を関連付けることができるデータの例である。メディア・ライブラリは、メディア・ファイルを、時相フォーマットで格納し、接続速度などの要素に基づいてストリーミングし、衛星撮影や環境撮影などの画像データベースは、同じ位置の画像のバージョンを維持し、医療または臨床データベースは、X線、身体スキャンなどに伴う患者の履歴を維持する。これらは、経時変化メディアを何らかの形で使用でき、したがって、本明細書に記載の時相ボリュームに適しているアプリケーションのほんの一例にすぎない。

20

- ・ 固定コンテンツ：WORM (Write - Once - Read - Many：追記型) ストレージとも呼ばれる。本明細書に記載の時相ボリュームを使用すると、WORMストレージ上のデータを過去の時点までトレースすることが可能になる。

- ・ 科学技術計算と研究室：時相ボリュームは、遺伝子データベース、データ解析、パターン発見、予測、集中計算、簡易臨床検査、データ監視、信号処理、数学、バイオインフォーマティクスなど、経時変化データを有するこれらのすべてを対象とするストレージの問題を解決することができる。

- ・ データ・アーカイブと監査：時相ボリュームは、それが格納するデータの履歴を保存しているので、たとえばその履歴をオフラインで取得し、テープまたはその他のメディアに格納することによって、これをアーカイビングに使用することができる。アーカイブは、後に分析または監査に使用することができる。アーカイブの例としては、それだけに限らないが、旧従業員データベース、旧売上げデータ、センサス・データ、ログ、およびサーバ履歴がある。別の例としては、様々な法律および規制により必要とされる規制上のアーカイビングがある。

30

- ・ データ・バージョニング：文書のバージョニングを、時相ボリュームに実装することができる。

【0022】

一実施形態では、時相ボリューム・マネージャ・レベルで、索引付けシステムまたは構造（たとえば、キャッシュ・オブジェクト、アレイなど）を使用して、時相ボリュームのコンテンツを保護することができる。キャッシュ・オブジェクトは、制限されたストレージ・スペースを使用しながらも、そのユーザにストレージが無限にあるかのような錯覚を投影するオブジェクトとして定義することができる。キャッシュ・オブジェクトにストレージを割り付けることにより、キャッシュ・オブジェクトを、スペース最適化スナップショットを作成するのに使用することができる。キャッシュ・オブジェクトは、そのベースとなるキャッシュ・ボリュームと呼ばれるボリュームから、それ自体のストレージを導出する。ポリシーに基づく無限容量の取決めを守る必要がある場合には、キャッシュ・ボリュームを増大させることもできる。

40

【0023】

一実施形態では、ボリュームを1つまたは複数の領域に分割することができる。領域は

50

、ディスクの1つの物理ブロックから、数キロバイト、数メガバイト、数ギガバイトなどまでのどの大きさの領域でもよい。ボリュームは、複数の領域に分割することができ、各領域は、それぞれに関連するタイムスタンプを有する。一実施形態では、アプリケーション（ファイル・システム、データベースなど）は、その領域にどのようなタイムスタンプを関連付けるべきかを指定することができる。一実施形態では、データが時相ボリュームに書き込まれるとき、アプリケーションがタイムスタンプを指定する。従来技術のシステムでは、書込みに関するタイムスタンプをアプリケーションが指定することが許可されていないことに留意されたい。

【0024】

これらの諸実施形態では、それだけに限らないが、入出力制御型チェックポイントイン
10 グ、アプリケーション制御型チェックポイントインテグ、周期型チェックポイントインテグを含めたいくつかの方法で、時相ボリューム内のデータの履歴を維持することができる。

【0025】

入出力制御型チェックポイントインテグでは、時相ボリュームを使用するアプリケーションは、時相ボリュームへの書込みおよび/または時相ボリュームの読取りを行ったときに、タイムスタンプを生成する。入出力制御型チェックポイントインテグでは、アプリケーションは、入出力要求がある度にその都度、タイムスタンプを時相ボリューム・マネージャに供給する、あるいは時相データと何らかの関係がある入出力要求のときにだけ、タイムスタンプを時相ボリューム・マネージャに供給することもできる。本明細書では、入出力要求という用語は、任意の入出力要求、コマンド、入出力制御要求、I O C T L、あるいはアプリケーションまたはその他のエンティティがボリュームと相互作用するために、たとえばボリュームにデータを読み書きするために使用するその他のメカニズムを指すために使用する。一実施形態では、アプリケーションは、何らかの入出力要求（たとえば、読取りまたは書込み）があれば、供給したタイムスタンプを用いて書込みを行う領域にタイムスタンプを付けるよう指定することができる。一実施形態では、時相ボリューム・マネージャは、入出力制御型チェックポイントインテグにおいて、タイムスタンプを指定する入出力要求（U N I X の I O C T L などの入出力制御コマンドまたは要求）をアプリケーションが発行するのを可能にする入出力要求インターフェースを設けていることができる。一実施形態では、時相ボリューム・マネージャとの時相読取りと時相書込みインターフェイスを使用して、入出力制御型チェックポイントインテグを実現することができる。たと
20 30 40 50

【0026】

図3は、一実施形態による、入出力制御型チェックポイントインテグを使用して、時相ボリュームを論理デバイス・レベルで管理するための方法の流れ図である。300で示されるように、まず、あるアプリケーションを対象とする時相データを格納するために、時相ボリュームを生成する。302で示されるように、時相ボリューム・マネージャは、アプリケーション（またはエージェントが使用されている場合はアプリケーションのエージェント）から入出力要求を受信する。入出力要求が時相要求である場合は、次いで、入出力要求は、時相ボリューム上の時相データを対象とする1つまたは複数のタイムスタンプを指定する。時相ボリューム・マネージャは、時相ボリュームにアクセスし、時相ボリュームに対する時相オペレーションを要求する際に使用するためのA P I を、アプリケーション（またはアプリケーション・エージェント）に提示する。入出力要求が時相読取り要求である場合は、306で示されるように、時相ボリューム・マネージャは、その入出力要求の指定する1つまたは複数のタイムスタンプで指示された時相データを、アプリケーションに返却する。入出力要求が時相書込み要求である場合は、308で示されるように、時相ボリューム・マネージャは、その入出力要求の指定するタイムスタンプに従って、時相ボリュームにおける1つ（または複数）の領域のチェックポイントを生成し、310で示されるように、その入出力要求の指定するデータを時相ボリュームに書き込む。

【0027】

10

20

30

40

50

アプリケーションは、その必要がある場合またはそれが望まれる場合、書込みの度にその都度タイムスタンプを生成する代わりに、ある領域または全体のボリュームを対象とするタイムスタンプを指定する入出力要求を発行する。この方法を、アプリケーション制御型チェックポインティングと呼ぶこともある。アプリケーション制御型チェックポインティングでは、書込みの度にその都度タイムスタンプを供給するのではなく、時相ボリューム内の新しいチェックポイント（ポイント・イン・タイム・コピー）を指定する入出力要求を発行する。アプリケーション制御型チェックポインティングでは、アプリケーションは、チェックポイント/バージョン（時相ボリュームのポイント・イン・タイム・コピー）をいつ作成するかを時相ボリュームに告げる。一実施形態では、これを、入出力制御（I/OCTL）オペレーションなどの入出力要求を使用して行うことができる。アプリケーション制御型の入出力要求がある度にその都度、1つまたは複数の領域にタイムスタンプが付けられる。一実施形態では、アプリケーション制御型チェックポインティングを、2つ以上の時相ボリュームにまたがって自動的に行うことができる。

10

【0028】

図4は、一実施形態による、アプリケーション制御型チェックポインティングを使用して、時相ボリュームを論理デバイス・レベルで管理するための方法の流れ図である。400で示されるように、まず、あるアプリケーションを対象とする時相データを格納するために、時相ボリュームを生成する。次いで、アプリケーションは、タイムスタンプを指定できず、したがって時相ボリュームのチェックポイントが生成されない1つまたは複数の「通常の」読み取りおよび/または書き込みを、時相ボリュームに対して行う。402で示されるように、時相ボリューム・マネージャは、時相ボリュームを対象とした新しいチェックポイントを生成するよう指定する入出力要求を、アプリケーションから受信する。次いで、404で示されるように、時相ボリューム・マネージャは、その入出力要求の指定するタイムスタンプに従って、時相ボリュームのチェックポイントを生成する。

20

【0029】

別の方法としては、周期的に、たとえば10秒毎または10分毎に自動チェックポインティングを行う方法がある。これを、周期型チェックポインティングと呼ぶこともある。周期型チェックポインティングでは、時相ボリューム・マネージャは、時相ボリューム・マネージャによる論理デバイス（ボリューム）レベルでの周期型チェックポインティングの構成を可能にするインターフェースを設けている。周期型チェックポインティングでは、時相ボリュームのインフラストラクチャ（たとえば、時相ボリューム・マネージャ）が周期的に、データの新しいチェックポイントを作成する。これにより、好ましくは、データへのそれぞれの変更を格納することによって浪費されていたかもしれないストレージ・スペースが節約される。従来技術では通常、あらゆる書き込みが保存されることに留意されたい。

30

【0030】

図5は、一実施形態による、周期型チェックポインティングを使用して、時相ボリュームを論理デバイス・レベルで管理するための方法の流れ図である。500で示されるように、まず、あるアプリケーションを対象とする時相データを格納するために、時相ボリュームを生成する。502で示されるように、次いで、時相ボリューム・マネージャは、タイムスタンプを指定できず、したがって時相ボリュームのチェックポイントが生成されない入出力要求をアプリケーションから受信する。一実施形態では、入出力要求は、チェックポイントを取得すべき時間間隔（または期間）を指定する。504で示されるように、時相ボリューム・マネージャは周期的に（たとえば、n分、n時間、またはn日毎に）、時相ボリュームのチェックポイントを生成する。一実施形態では、時相ボリューム・マネージャは、タイムスタンプを指定し、周期型チェックポインティングを使用する場合はチェックポイントを生成するよう指定する入出力要求を、アプリケーションから受信することに留意されたい。

40

【0031】

アプリケーションが時相書き込みインターフェースを使用している場合は、そのアプリケ

50

ーションは、書込みデータと共にタイムスタンプを指定することができ、ボリューム・マネージャは時相書込みを実行する。一実施形態では、タイムスタンプは、特定の書込みに適用され、ボリューム上に格納されているその他のデータは、それに時相情報を格納することも、格納しないこともできる。言い換えれば、一実施形態では、ボリューム上には時相データも存在すれば、時相データ以外のデータも存在し得る。データが時相データでない場合は、タイムスタンプが付されない。一実施形態では、アプリケーションがデータにタイムスタンプを付けるための入出力要求を発行した場合、あるいは周期型チェックポイントニングが実施された場合にだけ、データにタイムスタンプが付される。これを、1つまたは複数の領域あるいは全体のボリュームに対して適用する。

【0032】

時相ボリュームは、ある次元としての時間を有する。一実施形態では、時相ボリューム・マネージャは、ユーザが時相ボリューム内の時間次元（履歴）をトラバースすることを可能にするインターフェースを設けている。ユーザは、たとえば、時相ボリュームの履歴における以前の各時点での時相ボリュームのデータに対して、1つまたは複数のオペレーションを実施するためのインターフェースを使用して、時相ボリュームの時間次元をトラバースする。かかるオペレーションの例として、それだけに限らないが、時相ボリューム上の特定のバージョンまたは値のデータの有無を時間次元において探索すること、時相ボリュームをその時相ボリュームの履歴における特定の時点の状態に復元すること、時相ボリュームのスライス・イン・タイム・イメージまたはポイント・イン・タイム・イメージを生成することなどのうちの、1つまたは複数のオペレーションを挙げることができる。一実施形態は、それ自体の独立した履歴を有する時相ボリュームの時相イメージを生成するためのメカニズムを用意する。

【0033】

一実施形態は、時相ボリュームのスライス・イン・タイム・イメージを生成するためのメカニズムを用意する。一実施形態では、時相ボリューム・マネージャは、ユーザ/アプリケーションがある領域の、たとえばタイムスタンプT1とT2の間のコンテンツ（スライス・イン・タイム・イメージ）を要求するのを可能にするインターフェースを設けている。スライス・イン・タイム・イメージはそれ自体、時相ボリュームであるが、そのエクステンツは当初、タイムスタンプT1とT2の間に制限される。スライス・イン・タイム・イメージは、2つのタイムスタンプの間の時相ボリュームのスライスである。時相ボリュームのスライス・イン・タイム・イメージは、その時相ボリューム上にあるデータの履歴のスライスを提示する。スライス・イン・タイム・イメージを、それ自体の履歴も保存する別個の時相ボリュームとして閲覧することができる。一実施形態では、スライス・イン・タイム・イメージに対して読取りと書込みを実行する。スライス・イン・タイム・イメージに対しては、その他のオペレーション、たとえばデータ・マイニングを実施することもできる。たとえば、データのある期間にわたって収集することができ、その期間内のある月（たとえば、8月）を対象とするスライス・イン・タイム・イメージを取得することができ、そのスライス・イン・タイム・イメージに対して、データ・マイニングまたはその他のオペレーションを実施することができる。

【0034】

図6は、一実施形態による、時相ボリュームのスライス・イン・タイム・イメージを生成するための方法の流れ図である。600で示されるように、まず、あるアプリケーションを対象とする時相データを格納するために、時相ボリュームを生成する。602で示されるように、次いで、時相ボリューム・マネージャは、アプリケーションから、少なくともそのうちのいくつかは時相ボリュームのチェックポイントを生成できる時相および/または非時相入出力要求を受信する。周期型チェックポイントニングが使用されている場合は、次いで、時相ボリューム・マネージャが周期的に、時相ボリューム内のチェックポイントを生成していたことに留意されたい。604で示されるように、時相ボリューム・マネージャは、時相ボリュームを対象とした、時相ボリュームのスライス・イン・タイム・イメージを作成するよう要求する入出力要求をアプリケーションから受信する。606

で示されるように、時相ボリューム・マネージャは、その入出力要求に応答して、それが指定する2つのタイムスタンプに従い、時相ボリュームのスライス・イン・タイム・イメージを生成する。

【0035】

一実施形態は、時相ボリュームのポイント・イン・タイムの時相イメージを生成するためのメカニズムを備えている。ポイント・イン・タイム・イメージは、単一時間におけるボリュームのイメージであり、時相ボリュームの1次元イメージと見ることができる。ポイント・イン・タイム・イメージを時相イメージのベースとして使用する。その場合、ポイント・イン・タイム・イメージは、そのイメージのベースとなる時点に初期化された時相ボリュームと見ることができる。入出力（書込み）がある度にその都度タイムスタンプが供給される場合は、ボリュームのポイント・イン・タイム・イメージの数は有限であるはずである。時相ボリュームは、2つ以上のポイント・イン・タイム・イメージを含むデータ・ボリュームと考えられる。2つ以上のポイント・イン・タイム・イメージを組み合わせることもできる。時相ボリュームのスライス・イン・タイム・イメージは、2つ以上のポイント・イン・タイム・イメージを含む。

【0036】

図7は、一実施形態による、時相ボリュームのポイント・イン・タイム・イメージを生成するための方法の流れ図である。700で示されるように、まず、あるアプリケーションを対象とする時相データを格納するために、時相ボリュームを生成する。702で示されるように、次いで、時相ボリューム・マネージャは、アプリケーションから、少なくともそのうちのいくつかは時相ボリュームのチェックポイントを生成できる時相および/または非時相入出力要求を受信する。周期型チェックポイントイングが使用されている場合には、次いで、時相ボリューム・マネージャが周期的に、時相ボリューム内のチェックポイントを生成することに留意されたい。704で示されるように、時相ボリューム・マネージャは、時相ボリュームを対象とした、時相ボリュームのポイント・イン・タイム・イメージを作成するよう要求する入出力要求をアプリケーションから受信する。706で示されるように、時相ボリューム・マネージャは、その入出力要求に応答して、それが指定するタイムスタンプに従い、時相ボリュームのポイント・イン・タイム・イメージを生成する。

【0037】

一実施形態は、時相ボリューム上にある2つのバージョンのデータ間における増分変更を判定するためのメカニズムを備えている。判定された増分変更を、増分バックアップやレプリケーションなどのアプリケーションで使用する。一実施形態では、時相ボリューム・マネージャは、任意の2つの時点の間に時相ボリュームにどのような変化があったのかを判定するためのインターフェースを設けている。たとえば、タイムスタンプ付けを、ある期間にわたって実施する。ユーザは、2つの時間の間でデータのバックアップまたはアーカイブをとることを望むかもしれない。増分イメージを作成する。増分イメージは、差異を示すイメージである。増分イメージを使用して、時間T1からT2までの（たとえば、ある時点から以前のバックアップ時までの）バックアップまたはレプリケーションを生成する。時相ボリュームのコンテンツを使用して、それらの差異を識別する。このメカニズムは、好ましくは、増分バックアップまたはレプリケーション用アプリケーションまたはユーティリティが容易に実装できるようにし、また、アプリケーションまたはユーティリティが、変更部分を算定し、スナップショットを管理することを強いられるのではなく、時相ボリューム・マネージャがスナップショットを管理する。バックアップまたはレプリケーション用アプリケーションは単に、時相ボリューム・マネージャに、時間T1とT2の間における時相ボリュームのコンテンツ（あるいは、コンテンツの差異または差分（delta））を要求するだけである。

【0038】

一実施形態では、時相ボリュームの内部で維持されているデータの履歴を、たとえばIOCTLなどの入出力要求を使用して切り捨てる。スペース不足の故に切捨てが必要なこ

10

20

30

40

50

とも、履歴がもはや必要とされなくなった故に切捨てが望まれることもある。一実施形態では、このインフラストラクチャは、時相ボリュームのユーザがどれほどの量の履歴を維持するのかを決定する時相ウィンドウを設定できるようにする。時相ウィンドウの範囲を超える変更は自動的に切り捨てる。一実施形態では、時相ボリューム・マネージャは、スペースを解放するために時相ボリュームを切り捨てるためのインターフェースを設けている。たとえば、このインターフェースは、時相ボリュームが3ヶ月間保持すべき情報をアプリケーションによって指定できるようにする。3ヶ月よりも古い情報を、時相ボリューム・マネージャによって削除する。3ヶ月よりも古いデータを、周期的に、たとえば毎日または毎週削除する。

【0039】

10

一実施形態は、それ自体のポイント・イン・タイム・イメージの1つ、および/または2つの特定のタイムスタンプ間のスライス・イン・タイム・イメージから、ボリュームを迅速に復元するためのメカニズムを備えている。たとえば、あるアプリケーションがそのデータに10分毎にタイムスタンプを付けており、そのボリュームが何らかの形で破損している場合は、時相データを使用してそのボリュームの(たとえば、最近の破損していないタイムスタンプ付きの時相データへの)復元を実施する。

【0040】

時相ボリュームと時相ボリューム・マネージャ

以下では、ボリューム・マネージャとキャッシュ・オブジェクトを使用して、データ・ストレージ内の時相データを論理デバイス・レベルで維持するためのシステムと方法に関する例示的な諸実施形態を説明する。これらの諸実施形態は、例示的なものであり、時相ボリュームをイネーブルするためのその他のメカニズムを使用することもできることに留意されたい。これらの諸実施形態は、時相ボリュームを対象とする入出力要求をアプリケーションから受信する手段と、入出力要求にตอบสนองして時相ボリュームに対する論理デバイス・レベルの時相オペレーションを実施する手段とを備えている。

20

【0041】

これらの諸実施形態は、たとえば、VERITAS Volume Manager (VxVM)などのボリューム・マネージャを使用して、時相ボリュームを作成し使用するためのメカニズムを備えている。一実施形態は、データの履歴のアーカイブをとるために、キャッシュ・オブジェクトを使用する。一実施形態は、時相ボリュームのアプリケーションまたはユーザが、通常のボリューム・オペレーションに加えて時間ベースのオペレーションも実施できるようにする、1つまたは複数のインターフェースを設けている。

30

【0042】

索引付けシステムまたは構造(たとえば、キャッシュ・オブジェクト、アレイ、B+ツリーなど)を使用して、ボリューム・マネージャ内で、たとえばVxVM内で時相ボリュームを実装する。一実施形態では、キャッシュ・オブジェクトを使用して、時相ボリュームを実装する。図8は、一実施形態による例示的なキャッシュ・オブジェクトを示している。図8に示されるように、キャッシュ・オブジェクトは、制限されたストレージ・スペースを使用しながらも、そのユーザに「無限の」ストレージを投影するオブジェクトである。キャッシュ・オブジェクトは永続的なスペース最適化ストアと考えることもできる。一実施形態では、キャッシュ・オブジェクトは、ユーザがスペース最適化スナップショットを作成できるようにする。スペース最適化スナップショットは、キャッシュ・オブジェクトを介して作成されることにより、好ましくは、その元のボリュームほど多くのストレージを使用しない。実際のデータは、それが元のボリューム上で変更されたときにだけ、スペース最適化スナップショット上に書き込まれる(これをコピー・オン・ライトと呼ぶ)。スペース最適化スナップショットが必要とするどんなストレージ・スペースも、そのベースとなるキャッシュ・オブジェクトから提供される。次いで、キャッシュ・オブジェクトは、それ自体のストレージをキャッシュ・ボリュームから導出する。キャッシュ・オブジェクトは、ユーザが定義したいいくつかのポリシーに従い、必要に応じてキャッシュ・ボリュームを増大させることによって、無限のストレージという錯覚を維持する。

40

50

【 0 0 4 3 】

一実施形態では、キャッシュ・オブジェクトのスペース最適化を可能にする基本的な技術は、入出力ブロックの再ベクトル化(re-vectoring)である。キャッシュ・オブジェクトは、この再ベクトル化を、永続的な変換マップを使用して実現する。スペース最適化スナップショットに対してコピー・オン・ライトで書き込まれたデータを格納することになる同じ物理ストア(キャッシュ・ボリューム)にこれらの変換マップを格納する。さらに、キャッシュ・オブジェクトを用いると、その上に複数のストレージ・ユニットを切り出すことが可能となり、それによって、キャッシュ・オブジェクト上で複数のスペース最適化スナップショットを作成することが可能になる。

【 0 0 4 4 】

キャッシュ・オブジェクトは、複数のストレージ・ユニットをその上で作り出すことも可能にする。これは、キャッシュ・オブジェクトが異なるボリュームを対象とする複数の仮想アドレス可能範囲を提供することを意味する。したがって、各変換は、ストレージ・ユニットの識別子と、キャッシュ・ボリューム上のデータの物理オフセットにマッピングされるボリューム・オフセットとを備える探索キーから成る。探索可能なアドレス空間を、探索キーのサイズに従ってかなり広くすることができ、(無限ではないにせよ)膨大なストレージ・スペースが与えられている場合には、さらに広くすることもできる。これは、実際のデータ・ブロックへの再ベクトル化または変換に使用される構造が、探索可能性が高く、非常に効率的なものになることを意味する。B+ツリーは、この要件を満たす1つの構造であり、例示のために本明細書の諸実施形態において使用されているが、その他の実施形態では、他の構造も使用することに留意されたい。

【 0 0 4 5 】

一実施形態では、キャッシュ・オブジェクト用のB+ツリーを、スーパー・ブロック、復元ログ・エリア、フリー・ブロック・リストなど、その他のメタデータと共にキャッシュ・ボリュームに格納する。メタデータ以降のエリアは、データ領域を格納するのに使用される。キャッシュ・オブジェクト内のB+ツリーの各ノードが、(たとえば、ボリューム・マネージャによって決定される)1つのページ・サイズである。このツリーのリーフ・ノードが、キャッシュ・ボリューム上の実際の物理データ・ブロックを指す。

【 0 0 4 6 】

一実施形態では、キャッシュ・オブジェクト上のあらゆる入出力が、再ベクトル化される。再ベクトル化は、B+ツリーを使用して行うことができ、入出力ではまず、キャッシュ・ボリューム上のデータの物理オフセットを割り出すために、B+ツリー・ウォークを実施する。「有効」領域へのすべての読取りと書込みは、再ベクトル化する必要があることもある。領域は、アプリケーションによって1回または複数回そこへの書込みが行われた場合に「有効」と見なされ、そうでない場合は「無効」と見なされる。ボリュームの無効領域への書込みの場合は、キャッシュ・オブジェクトは、B+ツリー内の新しいエントリを新しいオフセットに割り付けることができ、これを、割付け書込み(allocating write)と呼ぶこともある。

【 0 0 4 7 】

以後の書込みで、それ以前のコンテンツを上書きする。一実施形態では、書込み以前のコンテンツが保護されるように、書込みはすべて、独立した割付け書込みとして扱われる。

【 0 0 4 8 】

図9は、一実施形態による時相ボリュームTVOLの構成を示している。TVOLは、通常データ・プレックス(data plex)P1、P2などを有する。プレックスは、以下の用語解説で定義されている。時相データを保持するPtと呼ばれる特別なプレックスも存在する。Ptは、キャッシュ・オブジェクトを介して作成できる1つのストレージ・ユニットSDtを有している。次いで、このキャッシュ・オブジェクトは、それ自体のストレージをキャッシュ・ボリュームから検索する。一実施形態では、キャッシュ・オブジェクトは、先に説明したB+ツリーを使用して、書込みを再ベクトル化する。別の一実施形態

では、時相ボリュームに、通常のデータ・ブックスを所在させず、 P_t だけを所在させる。

【0049】

通常の手込みを使用してボリュームに変更が加えられた場合には、その変更を、 P_t を含めたすべてのブックスに反映させる。 P_1 や P_2 のような通常のブックスは、単に新しいデータを用いて既存のデータを上書きすることしかできないが、 P_t はこれに対し、まだタイムスタンプの付けられていない最新のデータを上書きする。最新のデータにタイムスタンプが既に付けられている場合は、どんなタイムスタンプも割り当てずに、それを新しい位置に書き込む。

【0050】

一実施形態では、時相オペレーションを対象とする手込みは、 P_t を含めたすべてのブックスに及ぶことになる。時相オペレーションでは、タイムスタンプを指定することが必要となる。一実施形態では、時相手込みを、図10に示したように処理する。

【0051】

図10は、一実施形態による、時相ボリューム上のデータ・ブロックの変更を処理する様子を示し、さらに、所与の時点でのキャッシュ・オブジェクト下の構造を示している。キャッシュ・オブジェクトのベースとなるキャッシュ・ボリュームは、ディスク上の実際のデータ・ブロックを指すポイントをそれ自体のリーフが保持するツリー型の構造を含む。キャッシュ・オブジェクト内にまだ所在していない領域 B への手込み（割付け手込み）では、タイムスタンプ T_0 をエントリ B_0 に割り付ける。 B のコンテンツに変更があった場合は、 B_0 を上書きするのではなく、キャッシュ・ボリューム（ B_1 ）上の新しい領域にタイムスタンプ T_1 を割り付ける。キャッシュ・オブジェクトは、 B_1 をリストに追加する。 B_1 は、そのブロックの最近のコピーを含んでおり、 B_0 は、それ以前のコピーである。この連鎖は、ブロック B にさらなる変更が加えられる限り、引き続き行う。 n 個の変更が加えられた場合は、タイムスタンプ $T_n \sim T_0$ を用いて、ブロックを $B_n \dots B_2, B_1, B_0$ の順に連鎖する。この連鎖が、非常に長い連鎖となることもある。一実施形態では、かかる場合には、より高速な索引付けのために、連鎖自体をツリーの形に転換する。一実施形態では、後で説明するように、タイムスタンプを、そのコンテンツを識別するためのコピーのリスト中にある各ノードを用いて維持する。

【0052】

一実施形態では、ブロックの変更に従って、より古いバージョンをキャッシュ・オブジェクトの下に連鎖する。実際、キャッシュ・オブジェクトは、そのボリュームに属するすべてのデータ・ブロックの履歴を保持している。一実施形態では、最新のコピーは、通常のブックス P_1 、 P_2 内などで発見するが、 P_t （図9参照）の下には、より古いバージョンしか格納することができない。このように、現在データと履歴データの間には、明確な区画が存在する。このため、新しい時相オペレーションを、すべてのブックスまたはそのサブセットに対して行うのではなく、直接 P_t に対して行うことも可能になる。

【0053】

一実施形態では、時相ボリュームは、バージョンングを実現するためにタイムスタンプを使用する。時相ボリュームに対する通常の手込みは時相を実現しない。一実施形態では、ユーザまたはアプリケーションは、変更データの履歴を保持するために、時相手込みインターフェイス（以下でさらに説明する）を使用する。通常の手込みが使用された場合は、新しいデータが単に、古いデータを上書きするだけである。時相手込みに関しては、新しいデータは、古いデータに連鎖され、所与のタイムスタンプが付けられる。一実施形態では、タイムスタンプをアプリケーションから供給する。

【0054】

一実施形態では、各タイムスタンプは、相互に比較可能な状態にすることができ、時間的順序に配列できるように構成する。これは、時相読取りオペレーションを行うために、またデータの履歴に関する情報を得るために望まれ、あるいは必要とされることもある。タイムスタンプに使用できるデータ型としては、それだけに限らないが、整数データ型を

10

20

30

40

50

挙げることができる。これらの諸実施形態では、タイムスタンプは、実際のシステム時間、バージョン番号、増分カウンタ、または他の任意の適当なフォーマットでよい。一実施形態では、タイムスタンプがボリューム・マネージャによって解釈可能であることを必要としなくてよい。ボリューム・マネージャがタイムスタンプを解釈する必要がないことにより、ユーザがデータについて、将来のある時間とし得る有効時間を有することが可能になる。しかし、ボリューム・マネージャがタイムスタンプを解釈する必要がないということは、時相ボリュームがタイムスタンプの示す時間の一貫性を保証できないことも意味する。したがって、一実施形態では、一貫性のあるタイムスタンプが維持されるかどうかはアプリケーションまたはユーザ次第である。一実施形態では、ユーザは、後で説明するアプリケーション・エージェントのフレームワークを使用して、タイムスタンプの維持の負担を軽減する。

10

【0055】

一実施形態では、時相ボリュームは、(キャッシュ・オブジェクトによって提供される「無限の」ストレージのおかげで)「無限の」履歴を維持するが、ユーザは、履歴の一部または1ウィンドウを保持することしか望まないこともある。たとえば、現在とそれ以前のバージョンのデータにしか関心がなく、そのため、1つの変更/バージョンの時相ウィンドウしか必要としないアプリケーションがあるかもしれない。あるいは、ある期間だけに、たとえば最近のn分、n時間、n日などの期間だけに関する履歴を維持するのを望み、時相ウィンドウが、n分、n時間、n日などとなるアプリケーションもあるかもしれない。

20

【0056】

一実施形態では、時相ウィンドウを、時相ボリュームの作成中に設定することができ、それが望まれまたは必要とされる場合は、後にそれを変更することもできる。時相ウィンドウの範囲を超える時相データへのどんな変更も、格納されず、したがって失われる。したがって、たとえば時相ウィンドウが5分の場合には、5分よりも古いどんな変更も、使用不可能となる。一実施形態では、時相ウィンドウはデフォルトで、「無限」とされる。時相ウィンドウを増加させると、履歴の量が増加するが、時相ウィンドウを減少させると、時相ボリュームが履歴を切り捨てるが必要になることもある。

【0057】

時相ウィンドウの用途は、上記で示した例から明らかである。1つの変更の時相ウィンドウに関しては、現在バージョンと最終バージョンのデータだけが維持される。これは、たとえば、2つのバージョン間の差異を使用する増分バックアップまたはレプリケーションに有益なこともある。バージョン間の差異だけが使用されるので、好ましくは、バンド幅の使用を最適化する。第2の例では、時相ウィンドウは5分である。これは、たとえば周期的な増分バックアップまたは周期的なレプリケーションに有益なこともある。ユーザは、n分毎、n時間毎、n日毎など増分的にデータのバックアップをとることを望む場合は、nよりも若干広い時相ウィンドウを使用することができ、バージョン間の差異を増分バックアップに使用する。

30

【0058】

ボリュームの時相ウィンドウを、1対のタイムスタンプ $\langle T_a, T_b \rangle$ として表す。ただし、 T_a は、最も古いタイムスタンプであり、 T_b は、最新のタイムスタンプである。時相ボリュームが作成された場合でも、そのボリュームに対する最初の時相書込みがあるまで、それ自体の時相ウィンドウを定義せずにおくこともできる。タイムスタンプ T_x を伴う最初の書込みがあると、時相ウィンドウは、 $\langle T_x, T_x \rangle$ に初期化される。以後の(たとえばタイムスタンプ T_i を伴う)時相書込みが、時相ウィンドウを $\langle T_x, T_i \rangle$ に変更する。履歴が T_x から T_y に切り捨てられた場合(以下で詳細に説明する)は、時相ウィンドウは、 $\langle T_y, T_i \rangle$ に変わる。

40

【0059】

一実施形態は、時相ボリュームに対する時相入出力オペレーションを可能にする1つまたは複数のインターフェースを設けている。一実施形態では、通常のボリュームで許され

50

ているオペレーションを、時相ボリューム内で、それと同様にまたは同じ形で処理する。一実施形態では、時相ボリュームに対する通常の（非時相）読取りを、通常のボリューム内でそうするのと同様にまたは同じ形で進める。（時相プレックス P t 以外には）通常のプレックスが存在しない場合は、いかなるタイムスタンプも付いていないデータを、返却する。タイムスタンプの付いていないデータが何ら存在しない場合は、最新のタイムスタンプを伴うデータを返却する。一実施形態では、時相ボリュームに対する通常の書込みを、データレス・スペース最適化スナップショット（本明細書の後段で説明する）を作成することを除けば、通常のボリューム上でそうするのと同様にまたは同じ形で進める。一実施形態では、通常の読取りを、通常のボリュームにおいてそうするのと全く同様に時相ボリューム上で進めることができ、過去のデータのコピーを検索するための 1 つまたは複数の入出力要求を定義する。

10

【 0 0 6 0 】

一実施形態は、ユーザが過去における特定の時点のエクステンツのコンテンツを読み取れることを可能にするインターフェースを設けている。一実施形態では、エクステンツを、ボリューム上のオフセットとそのレングス (length) によって指定する。一実施形態では、タイムスタンプを使用して、そのコンテンツが必要とされる時間を指定することができ、そのタイムスタンプの解釈をユーザに委ねる。一実施形態では、タイムスタンプのタイプは、それに対して整数比較を実施できるように、ボリューム・マネージャ（たとえば V x V M ）によって指定されたタイプであることが好ましい。タイムスタンプの例としては、それだけに限らないが、エポックからの秒数として表されるトランザクション時間を表現するデータベース用のタイムスタンプ、バージョン番号を表現するバージョニング・システム用のタイムスタンプ、システム・クロックによって表される時間とするファイル・システム用のタイムスタンプ、整数カウンタである他の何らかのアプリケーション用のタイムスタンプを挙げることができる。

20

【 0 0 6 1 】

一実施形態では、タイムスタンプが 0 の場合、現在のデータのコピーを返却する。これは、通常の読取りと等価なものである。タイムスタンプが負である、たとえば - n である場合は、そのデータの n 番目に新しいコピーが返却される。これは、たとえば実際の時間には関心がなく、変更またはバージョンにしか関心のないアプリケーションに有益なこともある。

30

【 0 0 6 2 】

以下は、諸実施形態が提供できる、時相ボリューム上の様々な時相オペレーションを論理デバイス・レベルで行うための例示的な A P I である。これらの A P I は例示的なものであり、限定的なものではないことに留意されたい。

【 0 0 6 3 】

一実施形態は、エクステンツのコンテンツを非解釈バイト・バッファとして返却するのに使用する A P I を提供する。ボリュームは、それ自体のデバイス番号 (d e v i c e) を使用して指定する。以下は、U N I X 環境での A P I の例示的なフォーマットである。

```
void*  vol_temporal_read(voldevno_t    device,
                           timestamp_t   timestamp,
                           voff_t        offset,
                           size_t        length);
```

40

【 0 0 6 4 】

一実施形態では、v o l _ t e m p o r a l _ r e a d により、I O C T L（たとえば V O L _ T E M P O R A L _ R E A D）をデバイス上で使用する。V O L _ T E M P O R A L _ R E A D は、[o f f s e t , l e n g t h] の対で指定されたデータを、ユーザが用意したバッファに読み込むための例示的な I O C T L である。読み取るべきコピーは、タイムスタンプによって指定する。V O L _ T E M P O R A L _ R E A D は、指定されたタイムスタンプ（または最近のタイムスタンプ）にマッチする 1 つ（または複数）の領域を返却する。以下は、返却される構造体の例示的なフォーマットである。

50

```

struct vol_temporal_io {
    int            flags;           /* フラグ */
    timestamp_t    timestamp;       /* タイムスタンプ */
    caddr_t        buf;            /* バッファ */
    size_t         buf_len;        /* バッファ長 */
    voff_t         offset;         /* オフセット */
}

```

【 0 0 6 5 】

一実施形態は、期間を読み取るのに使用される A P I を提供する。このインターフェイスは、ユーザが 2 つの時点におけるエクステンツのすべてのバージョンを読み取ることができる。このエクステンツは、offset と length を使用して指定する。このインターフェイスでは、2 つのタイムスタンプを、すなわち、期間の開始に対応する第 1 のタイムスタンプと、その期間の終了に対応する第 2 のタイムスタンプとを提供する。一実施形態では、どちらのタイムスタンプも包括的である。一実施形態では、2 つのタイムスタンプ（たとえば period_start と period_end）が等しい場合は、期間の開始におけるコピーが、返却される。2 つのタイムスタンプがどちらも 0 である場合は、現在のコピーが返却される。タイムスタンプがどちらも負である、たとえば、それぞれ - m と - n (m > n) である場合は、m 番目に新しいコピーと n 番目に新しいコピーの間のすべてのコピーが返却されることになる（一実施形態ではどちらも包括的）。以下は、A P I の例示的なフォーマットである。

```

void* vol_temporal_period_read (voldevno_t    device,
                                timestamp_t     period_start,
                                timestamp_t     period_end,
                                voff_t         offset,
                                size_t         length);

```

【 0 0 6 6 】

一実施形態では、vol__temporal__period__read により、I O C T L（たとえば V O L __ T E M P O R A L __ P E R I O D __ R E A D）をデバイス上で使用する。この I O C T L は、2 つの時点（ある期間）におけるエクステンツ [offset , length] のすべてのコピーを返却する。この期間は、開始と終了のタイムスタンプを使用して指定される。実際のデータを返却する前に、この I O C T L を使用して、バッファのサイズを見出す。このために、空のバッファがサブミットされることもあり、I O C T L は、それが発見し得る領域の数 (reg__num) だけを返却する。ユーザは、領域の数を見出すと、メモリをバッファに割り付けることができ、2 回目の I O C T L の呼出しを行う。ユーザは、構造体（たとえば period__reg__t）がバッファ内での領域のコピーのレイアウトを記述するのに十分なメモリを割り付けることが好ましい。この 2 回目では、I O C T L は、2 つのポインタを返却する。たとえば、第 1 のポインタは、その領域のコンテンツのコピーを次々とそこに詰め込む、ユーザ・バッファ自体 (buf) である。その領域からはブロックを削除できるので、それぞれのコピーは、異なるサイズになることもある。第 2 のポインタ (reg__ptr) は、ユーザ・バッファ内での領域とそのコピーのレイアウトを記述する構造体を指す。この構造体は、各行がエクステンツの個々のコピーを記述し、各列が 1 つのコピー内の領域を記述するマトリックスと考える。以下は、返却される構造体 vol__temporal__region__period__io の例示的なフォーマットである。

```

typedef struct period_reg {
    voff_t         vtbr_offset;      /* 領域のオフセット */
    timestamp_t    vtbr_timestamp;  /* タイムスタンプのコピー */
    struct period_reg *vtbr_next_region; /* 次の領域 */
    struct period_reg *vtbr_prev_region; /* 前の領域 */
    struct period_reg *vtbr_next_copy; /* 次のコピー */
}

```

```

    struct period_reg *vtbr_prev_copy;      /* 前のコピー */
} *period_reg_t;

struct vol_temporal_region_period_io {
    int          flags;          /* フラグ */
    timestamp_t  start_timestamp; /* 期間の開始 */
    timestamp_t  end_timestamp;  /* 期間の終了 */
    caddr_t      buf;           /* バッファ */
    period_reg_t reg_ptr;       /* データのレイアウト */
    size_t       buf_len;       /* バッファの長さ */
    caddr_t      reg_num;       /* レジスタの数 */
    voff_t       offset;        /* エクステントのオフセット */
    size_t       len;           /* エクステントの長さ */
}

```

【 0 0 6 7 】

一実施形態では、時相書込みは、通常の手込みと全く同様であるが、ユーザが定義したタイムスタンプを含む。時相手込みの例示的なインターフェイスは、次の通りである。

```

size_t vol_temporal_write(voldevno_t device,
                           timestamp_t timestamp,
                           void *buffer,
                           voff_t offset,
                           size_t length);

```

【 0 0 6 8 】

このボリュームは、それ自体のデバイス番号 (d e v i c e) を使用して指定することができ、書き込むべきデータは、バッファに供給する。書き込むべき時間は、エクステントと共にタイムスタンプの形で供給され、(何らかのエクステントが存在すれば) そのエクステントの現在のタイムスタンプ以上であることが好ましい。エクステントは通常通り、[o f f s e t , l e n g t h] の対を使用して指定する。この関数に使用されるカーネル I O C T L (たとえば V O L _ T E M P O R A L _ W R I T E) は、所与のタイムスタンプを使用して時相ボリューム上の領域への手込みを行う。

```

struct vol_temporal_region_io {
    int          flags;          /* フラグ */
    timestamp_t  timestamp;      /* タイムスタンプ */
    caddr_t      buf;           /* バッファ */
    size_t       buf_len;       /* バッファの長さ */
    voff_t       offset;        /* オフセット */
}

```

【 0 0 6 9 】

以下は、タイムスタンプの指定した時間より前に時相ボリューム (デバイス) の履歴を切り捨てるのに使用できる例示的なインターフェイスである。

```

int vol_temporal_truncate_history (voldevno_t device, timestamp_t timestamp);

```

【 0 0 7 0 】

この返却値がステータス・インジケータである。この関数では、I O C T L (たとえば V O L _ T E M P O R A L _ T R U N C A T E) を使用し、その I O C T L によって時相ボリュームの履歴を切り捨てる。ユーザは、タイムスタンプを使用して、履歴をそれ以前に削除すべき時間を指定する。

```

struct vol_temporal_truncate{
    int          flags;          /* フラグ */
    timestamp_t  timestamp;      /* タイムスタンプ */
}

```

```

}

```

【 0 0 7 1 】

以下は、時相ウィンドウを変更するための例示的なインターフェイスである。

```

int  vol_temporal_change_window(voldevno_t    device,
                                timestamp_t    new_window,
                                twindow_t      window_type);

```

【 0 0 7 2 】

新しい時相ウィンドウ (new__window) は、タイムスタンプのタイプを使用して指定する。タイムスタンプの解釈は、ウィンドウのタイプがどのようなタイプであるかに依存する。このタイプは、たとえば、時間に基づくものでも、変更の数に基づくものでも、(タイムスタンプが解釈されない) 現在のタイムスタンプと最も古いタイムスタンプとの間の差異に基づくものでもよい。この関数は、I O C T L (たとえば V O L__T E M P O R A L__C H A N G E__W I N D O W) を使用して、時相ボリュームの時相ウィンドウを変更する。新しいウィンドウが古いウィンドウよりも狭い場合は、この I O C T L は、必要な履歴切捨て・アクティビティを実施する。以下は、この I O C T L の例示的な構造体フォーマットである。

```

struct  vol_temporal_change_window{
        int                flags;                /* フラグ */
        timestamp_t        window;               /* 新しいウィンドウ */
        twindow_t          type;                 /* ウィンドウのタイプ */
}

```

【 0 0 7 3 】

時相ボリューム (デバイス) 内の変更またはバージョニングの維持を中断するために、以下の例示的なインターフェイスを使用する。

```

int  vol_temporal_pause_history(voldevno_t  device);
int  vol_temporal_resume_history(voldevno_t  device);

```

【 0 0 7 4 】

これらの 2 つの関数は、対応するオペレーションのステータスを返却する。これらの 2 つの関数は、あるカーネル I O C T L (たとえば V O L__T E M P O R A L__P A U S E、V O L__T E M P O R A L__R E S U M E) を使用し、それらの I O C T L によりそれぞれ、履歴の保持が中断し再開する。このオペレーションは、たとえば、アーカイビングのために全体の履歴をオフラインで取得するのに有益なこともある。中断期間中にあたるタイムスタンプを伴ってボリュームに加えられたどの変更も、ボリュームが中断されている間は失われる可能性がある。あるいは、それらの入出力は、受け付けられない。

【 0 0 7 5 】

時相ボリューム内のデータの履歴に関する情報を取得するために、以下の例示的なインターフェイスを使用する。

```

int  vol_temporal_info(  voldevno_t    device,
                        voff_t          offset,
                        unsigned int    *copy_num,
                        timestamp_t     *min_timestamp,
                        timestamp_t     *max_timestamp,
                        timestamp_t     **change_timestamps,
                        voff_t          **changed_regions);

```

【 0 0 7 6 】

時相ボリュームは、それ自体のデバイス番号 (d e v i c e) を使用して指定する。特定の領域についての履歴が望まれている場合は、その履歴を、それ自体のオフセットを使用して指定する。このオフセットが負である場合は、全体のボリュームについての情報を返却する。返却される最初の情報は、その領域のコピーの数 (c o p y__n u m) である。領域上の最小タイムスタンプと最大タイムスタンプはそれぞれ、m i n__t i m e s t

ampとmax_timestampである。change_timestampsは、領域上のすべてのタイムスタンプのリストである。全体のボリュームについての情報が望まれている場合は、すべての変更された領域(changed_regions)のリストを返却する。上記のインターフェイスは、IOCTL(たとえばVOL_TEMPORAL_INFO)を使用して、コピーの数、タイムスタンプ、変更されたブロックのリストなど、履歴に関する情報を取得する。以下は、VOL_TEMPORAL_INFOの例示的な構造体である。

```
struct vol_temporal_info{
    voff_t    offset;           /* 領域のオフセット */
    caddr_t   copy_num;        /* コピーの数 */
    caddr_t   start_timestamp; /* 最初の変更 */
    caddr_t   end_timestamp;   /* 最後の変更 */
    caddr_t   change_timestamps; /* タイムスタンプのリスト */
    caddr_t   changed_regions; /* 変更された領域のオフセットのリスト */
}
```

10

【0077】

時相ボリュームが過去におけるボリュームのすべてのコピーを保持しているので、一実施形態は、特定の時点のボリュームの完全なイメージを取得するためのメカニズムを備えている。このイメージを、その特定の瞬間におけるボリュームのコンテンツを表す別個のボリュームである。このイメージは、データで一杯の状態(data-full)であることがあり、その場合は、すべてのブロックが新しいストレージ上にコピーされ、あるいは、データレスであることもあり、その場合は、実際のデータのコピーは存在せず、元の時相ボリュームを指す何らかの追跡構造体(tracking structure)だけが存在する。

20

【0078】

図11は、一実施形態による非時相スナップショットを示している。データのコピー・オン・ライトは、スナップショットの元のボリューム上にあるデータが上書きされたときに起こるが、その元のボリュームが時相ボリュームである場合には、これを回避する。図11に示されるように、TVOLは、通常の入出力が発生する時相ボリュームである。この時相ボリュームに時相アプリケーションが存在しない場合は、履歴は、プレックスP_tの下で維持されず、その領域には、タイムスタンプが付けられない。あるスナップショットがT₀において取得されたときは、P_tの下でのすべての領域が、タイムスタンプT₀を用いてマークされ、したがって、これらの領域に対する以後のどんな書込みも回避される。TVOLへの新しい書込みは、たとえばT₁まで達した場合、その領域にT₀のタイムスタンプが付いているのを知り、上書きを行わない。その代わりに、この新しい書込みには、(T₀にある)古いデータが連鎖されることになる。

30

【0079】

TVOLが通常のボリュームであれば、T₁の書込みが発生した場合は、コピー・オン・ライト・ポリシーの故に、その領域の古いコンテンツ(T₀のコンテンツ)をSTVOL(スナップショットの時相ボリューム、または時相スナップショット)にプッシュしているはずである。これは、STVOLがT₀のTVOLのイメージを表す故に必要とされることがある。しかし、TVOLは、時相ボリュームであるので、コピー・オン・ライトは必要とされない。TVOLは、T₁の書込みが発生した場合でも、その領域にはT₀のスタンプが付いているので、それらのバージョニングを自動的に行うことができ、したがって、T₀の領域のコンテンツは依然として、P_tの下で使用可能なままである。通常の下では、TVOLの読取りは、その領域が無効である場合にしか満足されることはない。したがってSTVOLの場合は、すべての領域は当初、無効としてマークされる。STVOLからある領域を読み取る場合には、その読取りは、タイムスタンプT₀と共にTVOLへとリダイレクトされる。この読取りでは、(T₁の)現在のコンテンツではなく、T₀の領域のコンテンツを読み取る。

40

【0080】

50

STVOLに対する書込みが発生した場合は、まず、それに対応するブロックが、(タイムスタンプ T_0 を伴う)TVOLからフェッチされ、次いで、STVOLに書き込まれる。次いで、この領域は、それ自体のコピーがSTVOL内にあり、TVOLにリダイレクトされる必要はないので、有効なものとしてマークされる。全体の領域を書き込む場合は、これをさらに最適化する。かかる場合には、フェッチングの部分を省略する。

【0081】

図12は、一実施形態による、スペース最適化時相スナップショットを時相ボリュームから導出する様子を示している。この元の時相ボリュームは、時相ウィンドウ $\langle T_a, T_b \rangle$ を有するTVOLである。ある時相スナップショットSTVOLが、 T_n において作成された場合は、その時相ウィンドウは、 $\langle T_n, T_n \rangle$ となる。ただし、 $T_a \leq T_n \leq T_b$ である。スナップショットの作成オペレーションではまず、プレックス(P_2 、 P_x など)を作成し、次いで、キャッシュ・オブジェクトとスナップショット・ボリュームSTVOLを作成する。次いで、STVOLのすべての領域を「無効」にセットする。一実施形態では、STVOLに対する入出力を、次のように進める。

- ・ 時間 T_x でのSTVOLに対する時相読取りは、

その領域が、無効である場合には、TVOLに対する T_n の時相読取りを実施し、データを返却する。 T_n は、スナップショットが作成された時間である。

その領域が、有効であり、 $T_x > T_n$ の場合には、STVOLに対する T_x の時相読取りを実施する。ただし、 T_x は、STVOL上で領域が変更された時間であり、その領域における最新のコピーとなる。そうでない場合には、TVOLの時相読取りを実施する。

- ・ 時間 T_x でのSTVOLに対する時相書込みは、

その領域が、無効である場合には、 T_n の時相読取りを使用して、その領域についての古いデータをTVOLから取得し、新しいデータを古いデータ上にオーバーレイし、時相書込みおよびタイムスタンプ T_x ($T_x > T_n$)を使用して、STVOLへの書込みを行う。マップにおいてその領域を有効にセットし、STVOLの時相ウィンドウを、 $\langle T_n, T_x \rangle$ に変更する。

その領域が、有効である場合には、タイムスタンプ T_x を用いてSTVOL自体に対する時相書込みを行い、STVOLの時相ウィンドウを、 $\langle T_n, T_x \rangle$ に変更する。

【0082】

ポイント・イン・タイム・イメージであるスナップショットとは異なり、スライスは、元の時相ボリュームをある期間だけ切り取ったものである。したがって、時相ウィンドウ $\langle T_a, T_b \rangle$ を伴う時相ボリュームのスライスは、 T_i から T_j の履歴のサブセットを保持する。ただし、 $T_i \geq T_a$ かつ $T_j \leq T_b$ である。

【0083】

一実施形態では、時相ボリュームのスライスを作成することは、スナップショットを作成するのと同様に行う。まず、割付けが行われるが、これには、プレックスとキャッシュ・オブジェクトを作成し、最終的に、ボリュームSTVOL(スナップショット時相ボリューム、または時相スナップショット)を作成することを要する。次いで、マップ内のすべての領域が無効としてマークされ、スライスの時相ウィンドウが、 $\langle T_i, T_j \rangle$ にセットされる。一実施形態では、STVOLに対する入出力を、次のように進める。

- ・ 時間 T_x でのSTVOLに対する時相読取りは、

その領域が、無効としてマークされている場合には、タイムスタンプ T_x を用いてTVOLに対する時相読取りを実施する。

その領域が、変更されたものとしてマークされており、 $T_x > T_i$ である場合には、STVOLに対する T_x の時相読取りを実施する。ただし、 T_x は、STVOL上で領域が変更された時間であり、その領域における最新のコピーとなる。そうでない場合には、TVOLの時相読取りを実施する。

- ・ 時間 T_x でのSTVOLに対する時相書込みは、

その領域が、無効としてマークされている場合には、 T_i の時相読取りを使用して、その領域についての古いデータをTVOLから取得し、新しいデータを古いデータ上にオー

10

20

30

40

50

バーレイし、時相書込みとタイムスタンプ T_x ($T_x > T_j$) を使用して、S T V O L への書込みを行う。マップにおいてその領域を有効にセットし、S T V O L の時相ウィンドウを、 $\langle T_j, T_x \rangle$ に変更する。

その領域が、有効なものとしてマークされている場合には、タイムスタンプ T_x を用いて S T V O L 自体に対する時相書込みを行い、S T V O L の時相ウィンドウを、 $\langle T_j, T_x \rangle$ に変更する。

【 0 0 8 4 】

一実施形態では、時相ボリューム上のデータの履歴がボリューム自体に格納されているので、ボリュームを、その過去のある状態に戻す。時相ボリューム (T V O L) が、時相ウィンドウ $\langle T_a, T_b \rangle$ を有し、ユーザが、その時相ボリュームの状態を時間 T_k における状態に復元することを望んでおり、 $T_a \leq T_k \leq T_b$ である場合には、一実施形態において、このプロセスを次の方法で実施する。

- ・ ユーザが履歴を削除することを望む場合には、

T V O L と有効とされているあらゆる領域を対象とするマップを参照し、 T_k より大きいタイムスタンプを有するすべての領域のコピーを削除する（これは基本的に、データ・ブロックをフリー・プールに戻すことを意味する）。

- ・ ユーザが履歴を保持することを望む場合には、

T V O L と有効とされているあらゆる領域を対象とするマップを参照し、（キャッシュ・オブジェクト上での）時間 T_k における領域のコピーのオフセットを発見し、たとえばそれが B_x であれば、時間 T_r における領域のコピーのオフセットを B_x にセットする。ただし、 T_r はその復元が開始された時間である。

【 0 0 8 5 】

一実施形態は、時相ボリュームをそのポイント・イン・タイム・イメージの 1 つから復元するためのメカニズムを備えている。このプロセスは、前段で説明したプロセスと同様に進める。T V O L 用のマップを使用する代わりに、スナップショット用のマップを使用する。しかし、この種の復元には、スナップショットが独立して変更されている可能性もあるので、データ移動が必要なこともある。一実施形態は、時相ボリュームをそのスライスの 1 つから復元するためのメカニズムを備えている。スライスは、変更されていることも、変更されていないこともある。この復元では、履歴の一部または全部を変更することができ、複雑なデータ移動が必要なこともある。

【 0 0 8 6 】

以上、入出力と関連するオペレーションを時相ボリュームに対して行うためのインターフェイスに関する例示的な諸実施形態を、説明してきた。以下では、先に説明したような一般的なインターフェイスではなく、限定的なものでもない、アプリケーションから見て有益である例示的なユーザ・レベルのインターフェイスを説明する。

【 0 0 8 7 】

アプリケーションは、ボリュームの領域の変更を途切れなく追跡することを望む場合、タイムスタンプを供給する必要がある。ボリュームに書込みを行うための時相書込みインターフェイスを使用する必要がある。このため、履歴の維持にしか関心のないアプリケーションについては、大きなオーバーヘッドが生まれる可能性がある。以下で説明する例示的なインターフェイスは、ユーザが時相ボリューム (d e v n a m e) への書込み中にタイムスタンプを指定することを必要としない。以下は、タイムスタンプを必要としない時相ボリュームへの書込みのための例示的な A P I である。

```
int temporal_sys_write (char          *dev_name,
                        void           *buffer,
                        unsigned int    offset,
                        unsigned int    length);
```

【 0 0 8 8 】

読取りを行うときは、ユーザは、時相ボリュームを読み取るための以下の例示的な A P I に示すように、時間列（たとえば「2003年1月20日 午後2時10分」）を使用

10

20

30

40

50

してタイムスタンプを指定する。

```
void temporal_sys_read (char          *dev_name,
                        char          *time_string,
                        unsigned int   offset,
                        unsigned int   length);
```

【 0 0 8 9 】

データベースに関しては、ブロックに対する入出力を行うことが望まれることもある。一実施形態は、データベース（またはその他のアプリケーション）が使用できる、時相ボリュームとのブロック・レベルのAPIを提供する。一実施形態では、以下の例示的なインターフェイスを、ブロック・レベルのAPIとして提供する。

```
void temporal_block_read (char          *dev_name,
                          unsigned int   block_offset);
```

```
void temporal_block_period_read (char    *dev_name,
                                 timestamp_t period_start,
                                 timestamp_t period_end,
                                 unsigned int block_offset,
                                 unsigned int *copy_num);
```

```
int temporal_block_write (char*    dev_name,
                          void      *buffer,
                          timestamp_t timestamp,
                          unsigned int block_offset);
```

```
int temporal_block_info (char    *dev_name,
                        unsigned int block_offset,
                        timestamp_t *min_timestamp,
                        timestamp_t *max_timestamp,
                        timestamp_t **change_timestamps);
```

【 0 0 9 0 】

アプリケーションは、先に説明したインターフェイスを使用して、また一実施形態では、1つまたは複数のボリューム・マネージャ・ライブラリ（たとえばVxVMライブラリ）を使用して、時間特性にアクセスするが、提供されたインターフェイスまたはAPIが適切でなく、インターフェイスをそれ自体のニーズに合わせることを望むこともある。たとえば、データベースは問合せ処理のためのブロック抽出により関心があることもあり、ファイル・システムはエクステンツを用いてその仕事を行うこともあり、レポジトリはすべてのオブジェクトを抽出することもある。一実施形態は、アプリケーションがよりその個々のニーズに特有のエージェントを書くことを可能にし得る、エージェント・ベースのフレームワークを提供する。

【 0 0 9 1 】

本明細書では、エージェントは、アプリケーションの代わりに時相ボリュームと相互作用するブローカであり、アプリケーションの（一部または全部の）セマンティクスをサポートしている。エージェントは、アプリケーションが時相ボリュームと容易に相互作用できるようにすることが好ましい。さらに、エージェントは、アプリケーションがデータを生データのままでなく、必要とするフォーマットで抽出することを可能にする。これは、時相ボリュームには何らの影響もないが、アプリケーションには有利なことがある。

【 0 0 9 2 】

一実施形態では、エージェントは、そのアプリケーションのニーズに適した任意のインターフェイスを、アプリケーション側にエクスポートする。エージェントは、ボリューム側では、時相ボリュームのインフラストラクチャによって提供されたインターフェイス、

10

20

30

40

50

たとえば先に説明したインターフェイスを使用する。エージェントは、ライブラリ、ユーティリティ、あるいはスクリプトの形さえとる。エージェントは、アプリケーションへのフック（特定の呼出しのトラップングなど）として書くことができ、あるいは、アプリケーションと密接に統合することもできる。一実施形態では、時相ボリュームは、標準ライブラリとコマンド・セットを提供することができ、それらを使用してエージェントを書くことができる。

【 0 0 9 3 】

エージェントは、たとえば時相オペレーションをサポートしていないが、本明細書に記載の時相ボリュームを使用した時相オペレーションをサポートすることを望む既存のシステム内で使用することができる。以下は、エージェントを使用するいくつかの例示的なシナリオである。

10

【 0 0 9 4 】

第1の例示的なシナリオでは、いくつかのデータベースは、時相オペレーションまたは時相データをサポートすることができないが、時相ボリュームを使用することを望む可能性がある。これらのデータベースを対象として、データベースを作成し、1つまたは複数のデータベース・エージェントを使用して、時間次元がこれらのデータベースに関連付けられるように、オペレーションを再ベクトル化する。このようにして、時相でない既存のデータベースは、インターフェイスの働きをするデータベース・エージェント層を1つ（または複数）の時相ボリュームに追加して、時相オペレーションをサポートする。

【 0 0 9 5 】

20

第2の例示的なシナリオでは、レポジトリは、エージェントを時相ボリュームとのインターフェイスとして使用して、それ自体のデータ・オブジェクトを時間的に格納し検索する。オブジェクトに変更がある場合には、その変更要求を、古いオブジェクトがあったのと全く同じ位置にそのオブジェクトを書き込む、オブジェクト・エージェントへと経路変更する。このレポジトリは、時相ボリューム上で作成されるので、変更のあったブロックだけが、バージョン化される。オブジェクト・エージェントは、特定のタイムスタンプ中にオブジェクトのバージョンを抽出している間には、その適切なタイムスタンプを、ボリュームの読取りに追加し、したがって、索引付けまたは抽出のパフォーマンスが低下することなく、オブジェクトの過去のバージョンが返却される。

【 0 0 9 6 】

30

一実施形態では、先に説明した時相ウィンドウのパラメータおよび/または時相ボリュームの周期型チェックポイントング機能を、周期型レプリケーションを実施する際に使用する。周期型レプリケーションでは、変更の発生に伴ってそれを送り届けるのではなく、周期的に1組のボリュームの複製を行う。これは、それだけに限らないが、（同期レプリケーションと比較して）入出力の待ち時間が短縮されること、および（非同期レプリケーションと異なり）ログが不要であることを含めて、複数の利点を有する。時相ウィンドウは、レプリケーションの期間よりも若干広く設定することができ、したがって、その期間内に発生した変更だけが格納される。全体の領域のコピーではなく、変更だけをこの領域に送ることによって、レプリケーションをさらに最適化することができる。これらの変更は、データの最新のコピーと最後のレプリケーションの時点（最後の期間の終了時点）におけるデータのコピーを使用して計算する。

40

【 0 0 9 7 】

時相アプリケーションは、極めて膨大な量のデータを生成することができ、これらのデータは、長期間にわたり時相ボリューム内で保持することを求められることもある。より古い履歴をオフラインで取得し、最新の履歴をオンラインで維持することは、たとえばユーザが、変更は失いたくないが、最近の履歴を使用したいという場合に有益である。かかる場合は、時相ウィンドウは、「無限」とされる。これにより、過去におけるデータのすべての状態を時相ボリューム内にもつことになり、そのため、大容量のストレージが必要になる。一実施形態では、アクセス頻度の低いデータおよび/または履歴を、アプリケーションがオンラインの状態にある間でも透過的に他のメディア（たとえば、テープまたは

50

光デバイス)に移動させるために、H S M (階層型ストレージ管理)を使用する。

【0098】

一実施形態では、時相ボリュームをバックアップ・ユーティリティと統合する。バックアップ・ユーティリティは、差分バックアップ技術を使用して、ファイルに加えられた変更だけを格納することができ、それにより、バックアップを作成するのに必要なバンド幅が最適化される。ネットワーク・バックアップ・ユーティリティは、バックアップを中央のバックアップ・サーバ(必ずしも単一である必要はない)に格納することができ、それにより、多数のアプリケーション・クライアントが、1つのバックアップ・インスタンスだけを共用することが可能になる。時相ボリュームを使用すると、この差分保存を論理デバイス・レベルで自動的に行うことができ、したがって、ファイルの差分を計算する必要はなく、明示的にそれを格納する必要もない。時相ボリュームは、複数のクライアントを対象とする差分のバージョニングに使用することもできる。複数のユーザが独立してファイルを変更した場合でも、それらのバージョンを同じ時相ボリューム内に格納することができ、それにより、後のファイルの各バージョンの抽出が、円滑に進められる。

10

【0099】

一実施形態では、時相ウィンドウを用いると、バックアップ・ユーティリティがクライアント側のファイルへの変更を追跡することを可能にする。1つの変更の時相ウィンドウは、途切れなく変更を追跡する際に使用する。元のコンテンツは、(必ずしも)サーバからフェッチする必要はなく、その代わりに、データの以前の状態を対象とする時相読取りを使用して検索するので、差分の取得を、最適化する。一実施形態では、この差分をクライアント側で作成する。

20

【0100】

用語解説

階層型ストレージ管理(H S M)：ポリシー・ベースの自動データ移行、自動バックアップ、圧縮、セキュリティ、アーカイビングなど、企業規模のストレージ管理サービスを提供する1組のソフトウェアおよび/またはハードウェアである。すべてのH S Mソリューションにおいて最も一般的な機能は、データ移行、すなわちデータをそのアクセス・パターンに基づいて、テープ、C D、W O R Mドライブなど、より廉価なストレージまたはニアライン・ストレージにオフロードする能力である。H S Mは、そのアクセスがあったときは、かかるデータを透過的に呼び出す。

30

【0101】

ミラー、ブレックス：ボリュームのデータの論理マッピングまたはコピーである。ボリューム内には、1つまたは複数のブレックスが存在することがある。個々のミラーはそれぞれ、完全なコピーを有している必要はないが、物理メディア上の何らかのエラーが原因で、またはそのミラーのジオメトリの特性の故に、データがマッピングされていない穴が存在することもある。これを、スパス・ブレックスと呼ぶこともある。しかし、完全なブレックスが、ボリューム内に存在している必要がある。

【0102】

元のボリューム：1つまたは複数のミラーが、スナップショットまたはデタッチのオペレーションの結果としてそこから分離される(またはされた)ボリュームである。

40

【0103】

レプリケーション：1次サイト上のデータ・ボリュームと、地理的に離れた2次サイト上のそれ自体のイメージとの間の一貫性を維持するプロセスである。3つのモードまたはレプリケーションがある。すなわち、同期型、非同期型、周期型がある。同期型レプリケーションは、2次サイトも更新されたときにだけ書込みを返却する。非同期型レプリケーションは、変更をログの形で記録し、書込みを直ちに返却し、その後(バックグラウンドで)2次サイトの更新を行う。周期型レプリケーションは、一貫性を維持するために、増分変更部分を2次サイトに転送する。

【0104】

復元：ある論理的な破損の後に、元のボリュームをその以前のポイント・イン・タイム

50

・イメージの1つに復帰させる、ボリューム・マネージャのオペレーションである。

【0105】

再同期化：コンテンツの古いプレックス(stale plex)をボリュームのコンテンツと同期させ、あるいはボリュームを（インスタント・スナップショットを用いて）別のボリュームと同期させる、ボリューム・マネージャのオペレーションである。

【0106】

スナップショット・ボリューム：元のボリュームのポイント・イン・タイム・イメージである。

【0107】

スナップショット：スナップショット・ボリュームを作成するボリューム・マネージャのオペレーションである。

【0108】

ボリューム：物理ブロック・ストレージ・デバイス（ディスク）のように振る舞うが、任意複雑形状の、ことによると冗長な内部ジオメトリを有する、ボリューム・マネージャの疑似デバイスである。ボリュームは最終的に、その論理スペースを実際のディスク・ドライブ上の物理スペースにマッピングする（ボリューム・マネージャから物理ディスクとして見えるものは、別のドライバまたはストレージ・アレイによって作り出された疑似デバイスであるかもしれないが）。

【0109】

VxFS：VERITAS File Systemである。これは、ロールバックとリドゥを可能にすることによって、ファイルに関するインテントまたはファイルに対して行われたオペレーションが保護される、ジャーナリング・ファイル・システムである。VxFSは、クローニングも可能にするが、これは、ファイル・システムのポイント・イン・タイム・イメージを取得するのと等価である。

【0110】

チェックポイントイング：データ・ボリュームの論理的なポイント・イン・タイム・イメージを取得するプロセスである。

【0111】

無効領域：アプリケーションによって書込みが依然として行われておらず、したがっていかなる有効なアプリケーション・データも含んでいないデータ・ボリュームの領域である。

【0112】

有効領域：アプリケーションにとって有効なデータを含んでいるデータ・ボリュームの領域である。

【0113】

結論

これらの様々な実施形態はさらに、前段の記載に従って、搬送媒体に実装された命令および/またはデータを受信し、送信し、または格納することを含む。一般的に言えば、搬送媒体としては、磁気または光メディアなどのストレージ・メディアまたはメモリ・メディア、たとえばディスクまたはCD-ROM、RAM（たとえば、SDRAM、DDR SDRAM、RDRAM、SRAMなど）やROMなどの揮発性または不揮発性メディア、ならびにネットワークおよび/または無線リンクなどの通信メディアを介して運搬される電気信号、電磁気信号、デジタル信号などの伝送メディアまたは信号を、挙げることができる。

【0114】

添付の図面に示し本明細書において説明した様々な方法は、それらの例示的な諸実施形態である。これらの方法は、ソフトウェア、ハードウェア、またはそれらの組合せで実装する。方法の順序は、変更してよく、また、様々な要素を追加し、順序を入れ替え、組合せ、省略し、修正するなどしてもよい。

【0115】

10

20

30

40

50

本開示の利益を有する当業者には自明であるように、様々な修正および変更を加える。本発明は、かかる修正および変更をすべて包含するものであり、したがって、上記の記載は、限定的な意味にではなく、例示的なものと見なされるべきである。

【図面の簡単な説明】

【 0 1 1 6 】

【図 1】一実施形態による、時相ボリュームにおける時相オペレーションを管理する時相ボリューム・マネージャの図である。

【図 2】一実施形態による、アプリケーション・エージェントを使用して、時相ボリューム・マネージャとインターフェイスするアプリケーションの図である。

【図 3】一実施形態による、入出力制御型チェックポイントイングを使用して、時相ボリュームを論理デバイス・レベルで管理するための方法の流れ図である。

10

【図 4】一実施形態による、アプリケーション制御型チェックポイントイングを使用して、時相ボリュームを論理デバイス・レベルで管理するための方法の流れ図である。

【図 5】一実施形態による、周期型チェックポイントイングを使用して、時相ボリュームを論理デバイス・レベルで管理するための方法の流れ図である。

【図 6】一実施形態による、時相ボリュームのスライス・イン・タイム・イメージを生成するための方法の流れ図である。

【図 7】一実施形態による、時相ボリュームのポイント・イン・タイム・イメージを生成するための方法の流れ図である。

【図 8】一実施形態による例示的なキャッシュ・オブジェクトの図である。

20

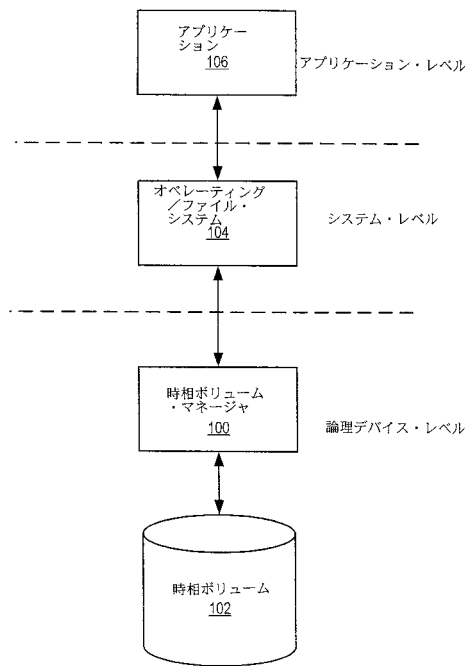
【図 9】一実施形態による時相ボリューム T V O L の構成図である。

【図 10】一実施形態による、時相ボリューム上のデータ・ブロックの変更を処理する様子を示し、さらに、所与の時点でのキャッシュ・オブジェクト下の構造を示す図である。

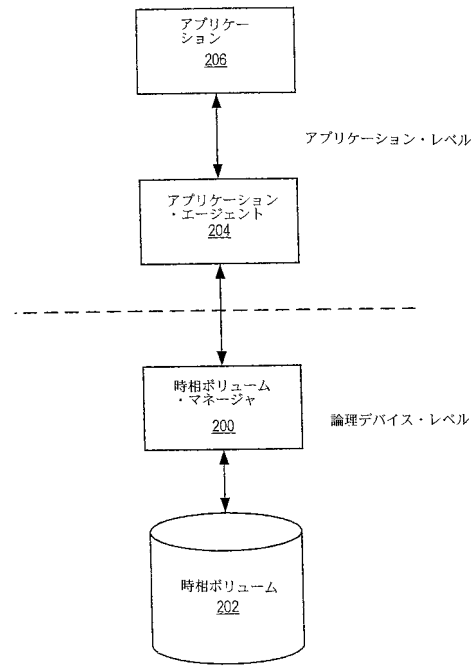
【図 11】一実施形態による、非時相スナップショットの図である。

【図 12】一実施形態による、スペース最適化時相スナップショットを時相ボリュームから検索する様子を示す図である。

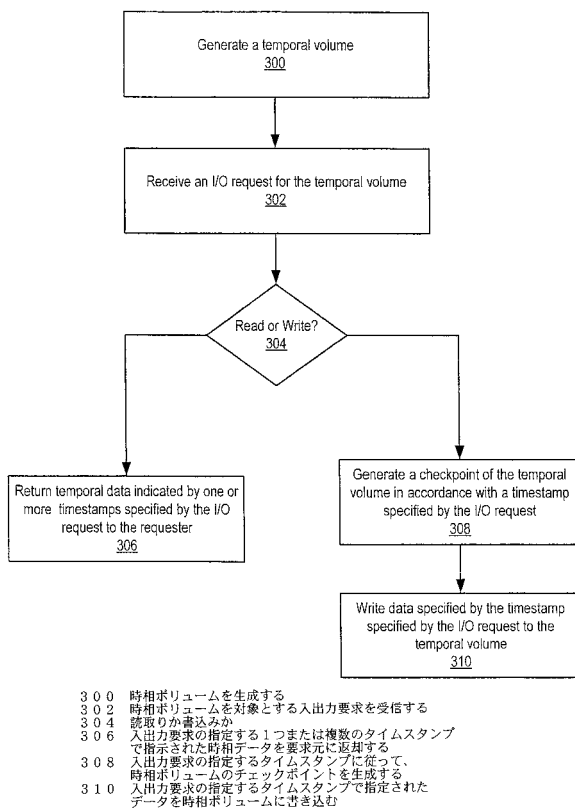
【図 1】



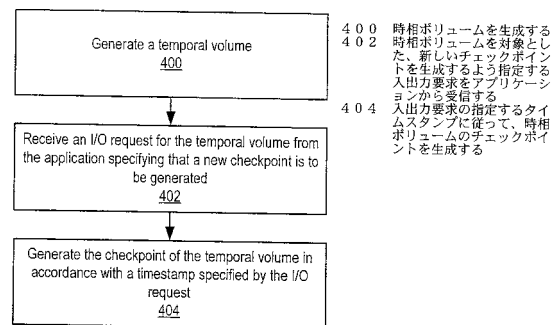
【図 2】



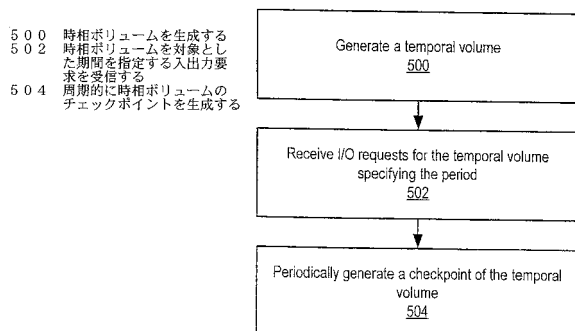
【図 3】



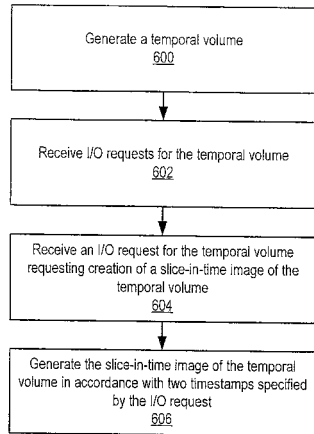
【図 4】



【図 5】

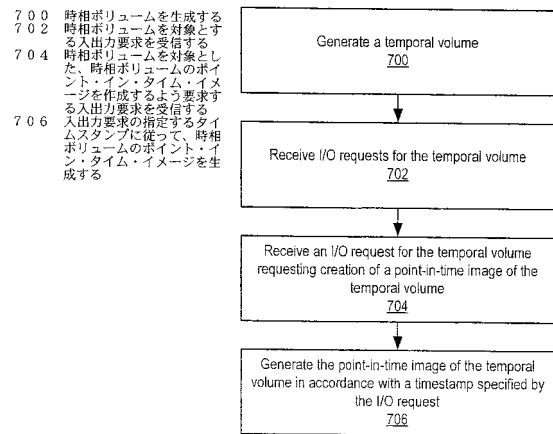


【図 6】



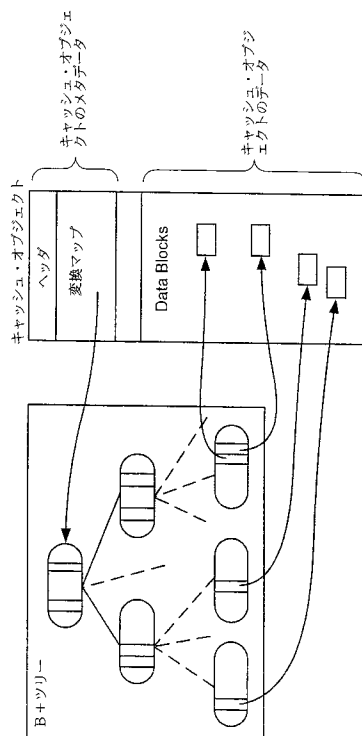
6 0 0 時相ボリュームを生成する
 6 0 2 時相ボリュームを対象とした入出力要求を受信する
 6 0 4 時相ボリュームを対象とした時相ボリュームのスライス・イン・タイム・イメージを作成するよう要求する入出力要求を受信する2つのタイムスタンプに従って、時相ボリュームのスライス・イン・タイム・イメージを生成する

【図 7】

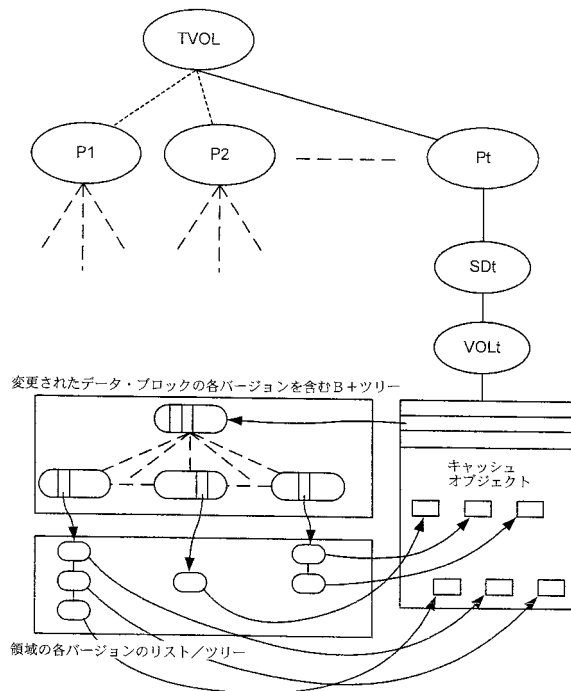


7 0 0 時相ボリュームを生成する
 7 0 2 時相ボリュームを対象とする入出力要求を受信する
 7 0 4 時相ボリュームを対象とした、時相ボリュームのポイント・イン・タイム・イメージを作成するよう要求する入出力要求を受信するタイムスタンプに従って、時相ボリュームのポイント・イン・タイム・イメージを生成する

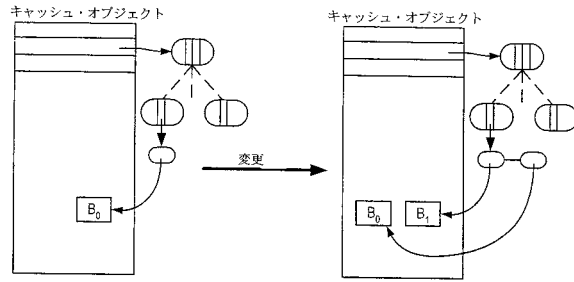
【図 8】



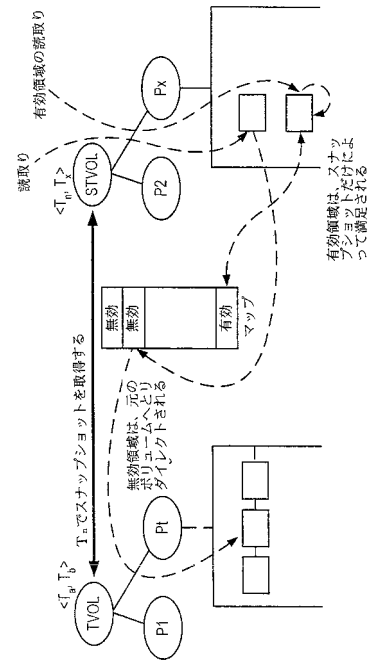
【図 9】



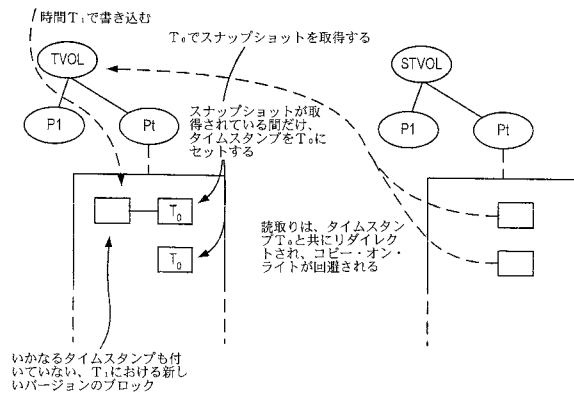
【図 10】



【図 12】



【図 11】



フロントページの続き

- (72)発明者 ケクレ, アナンド・エイ
インド国・ブーン 4 1 1 0 4 5 ・バナー・エス ナンバー 7 4 ・コスモス ガーデنز・1 0
- (72)発明者 パンシュブッデ, アンカール
インド国・ナグプール 4 4 0 0 0 9 ・ハヌマン ナガー・ヴィーナ・1 9 9

合議体

審判長 水野 恵雄

審判官 稲葉 和生

審判官 衣川 裕史

- (56)参考文献 特開2 0 0 1 - 1 5 9 9 9 3 (J P , A)
特開2 0 0 2 - 1 5 7 1 5 6 (J P , A)

- (58)調査した分野(Int.Cl. , D B 名)

G06F3/06

G06F12/00

G06F12/16