

(11) 特許出願公開番号

特開2012-123793

(P2012-123793A)

(43) 公開日 平成24年6月28日(2012.6.28)

(51) Int.Cl.

G O 6 F 12/00 (2006.01)

F 1

G O 6 F 12/00 5 3 1 M

G06F 12/00 517

テーマコード (参考)

審査請求 未請求 請求項の数 18 O L (全 23 頁)

(21) 出願番号 特願2011-247465 (P2011-247465)

(22) 出願日 平成23年11月11日 (2011.11.11)

(31) 優先權主張番号 12/963146

(32) 優先日 平成22年12月8日 (2010.12.8)

(33) 優先權主張国 米国 (US)

(特許庁注：以下のものは登録商標)

1. RRAM

(71) 出願人 390009531

インターナショナル・ビジネス・マシーン

ズ・コーポレーション

INTERNATIONAL BUSIN

ESS MACHINES CORPO

RAT ION

アメリカ合衆国10504 ニューヨーク

州 アーモンク ニュー オーチャード

ロード

(74) 代理人 100108501

弁理士 上野 剛史

(74) 代理人 100112690

一 種 佐 太 士 理 弁

(74) 代理人 100091568

弃理士 市位 嘉宏

[最終頁に続く](#)

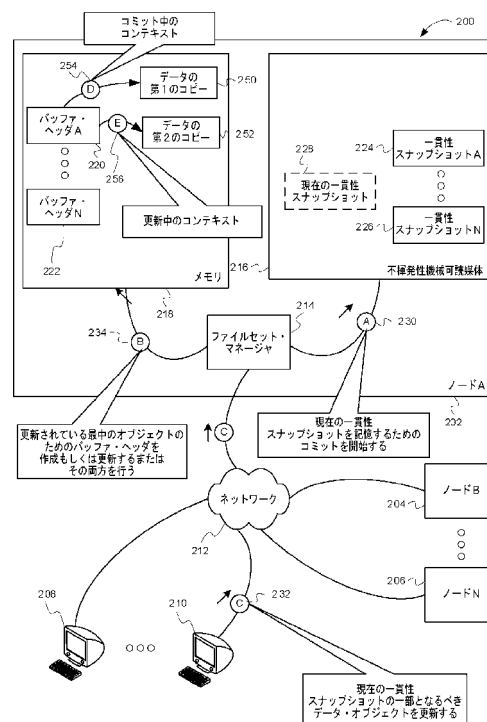
(54) 【発明の名称】 リダイレクト・オン・ライト・ファイル・システムにおける複数のコンテキストを提供する方法、装置およびコンピュータ・プログラム

(57) 【要約】 (修正有)

【課題】リダイレクト・オン・ライト・ファイル・システムにおける複数のコンテキストを提供する。

【解決手段】リダイレクト・オン・ライト・ファイル・システム内の複数のデータ・オブジェクトの現在の一貫性スナップショットを永続的に記憶するためのコミットを開始する。複数のデータ・オブジェクトの各々は、第1のコンテキストを有する、第1のコピーを有する。現在の一貫性スナップショットを記憶するためにコミットしている最中に、複数のデータ・オブジェクトのうちのデータ・オブジェクトに対する更新を受信すると、第1のコピーからデータ・オブジェクトのデータの第2のコピーを作成する。データの第2のコピーは少なくとも2つのコンテキストのうちの第2のコンテキストを有する。データ・オブジェクトに対する更新の受信に応答して、更新に基づきデータ・オブジェクトのデータの第2のコピーを更新する。

【選択図】図2



【特許請求の範囲】**【請求項 1】**

リダイレクト・オン・ライト・ファイル・システム内の複数のデータ・オブジェクトの現在の一貫性スナップショットを不揮発性機械可読媒体に記憶するためのコミットを開始することであって、前記複数のデータ・オブジェクトの各々は、コミット中のコンテキストを有する、前記複数のデータ・オブジェクトのデータの第 1 のコピーを有する、前記開始することと、

他の一貫性スナップショットの世代値に対して一意的である世代値を、前記現在の一貫性スナップショットに付与することと、

前記現在の一貫性スナップショットを記憶するために前記コミットしている最中に、前記複数のデータ・オブジェクトのうちのデータ・オブジェクトに対する更新を受信することと、

前記データ・オブジェクトに対する前記更新を受信することに応答して、

前記データ・オブジェクトのための世代値をインクリメントすることと、

前記データ・オブジェクトの前記世代値から導き出される世代値を前記更新に関連付けることと、

前記データ・オブジェクトの前記データの第 1 のコピーからコピーされる、前記データ・オブジェクトのデータの第 2 のコピーを作成することであって、前記データ・オブジェクトの前記データの第 2 のコピーは更新中のコンテキストを有する、作成することと、

前記データ・オブジェクトの前記データの第 1 のコピーを更新することとは独立に、前記更新に基づき前記データ・オブジェクトの前記データの第 2 のコピーを更新することと、

を含む方法。

【請求項 2】

前記方法は、前記データ・オブジェクトに対する前記更新を受信することに応答して、揮発性機械可読媒体内に、前記データ・オブジェクトに関連付けられるバッファ・ヘッダを作成することをさらに含み、前記バッファ・ヘッダは、前記データ・オブジェクトの前記データの第 1 のコピーを指し示す第 1 のデータ・ポインタおよび前記データ・オブジェクトの前記データの第 2 のコピーを指し示す第 2 のデータ・ポインタを含む、請求項 1 に記載の方法。

【請求項 3】

前記現在の一貫性スナップショットは、前の一貫性スナップショットより後の、前記複数のデータ・オブジェクトに対する更新を含む、請求項 1 に記載の方法。

【請求項 4】

前記現在の一貫性スナップショットを記憶するためのコミットを前記開始することは、一貫性スナップショットを作成する定期的作業に応答するものである、請求項 1 に記載の方法。

【請求項 5】

前記データ・オブジェクトの前記世代値から導き出される前記世代値を前記更新に関連付けることが、前記データ・オブジェクトの前記世代値に等しい前記世代値を前記更新に関連付けることを含む、請求項 1 に記載の方法。

【請求項 6】

リダイレクト・オン・ライト・ファイル・システム内の複数のデータ・オブジェクトの現在の一貫性スナップショットを永続的に記憶するためのコミットを開始することであって、前記複数のデータ・オブジェクトの各々は、異なるコンテキストを有する、前記複数のデータ・オブジェクトのデータの複数のコピーを有するように構成可能であり、前記複数のデータ・オブジェクトの各々は、少なくとも 2 つのコンテキストのうちの第 1 のコンテキストを有する、前記データの少なくとも 2 つのコピーのうちの第 1 のコピーを有する、前記開始することと、

前記現在の一貫性スナップショットを記憶するためにコミットしている最中に、前記複数のデータ・オブジェクトのうちのデータ・オブジェクトに対する更新を受信することと、

前記データ・オブジェクトに対する前記更新の受信に応答して、

前記第 1 のコピーから前記データ・オブジェクトのデータの第 2 のコピーを作成することであって、前記データの第 2 のコピーは前記少なくとも 2 つのコンテキストのうちの第 2 のコンテキストを有する、前記作成すること、および

前記更新に基づき前記データ・オブジェクトの前記データの第 2 のコピーを更新することと、
を含む方法。

10

【請求項 7】

他の一貫性スナップショットの世代値に対して一意的である世代値を、前記現在の一貫性スナップショットに付与することをさらに含む、請求項 6 に記載の方法。

【請求項 8】

前記データ・オブジェクトに対する前記更新の受信に応答して、

前記データ・オブジェクトのための世代値をインクリメントすること、および

前記データ・オブジェクトの前記世代値から導き出される世代値を前記更新に関連付けること、
をさらに含む、請求項 7 に記載の方法。

20

【請求項 9】

前記現在の一貫性スナップショットは、前の一貫性スナップショットより後の、前記複数のデータ・オブジェクトに対する更新を含む、請求項 6 に記載の方法。

【請求項 10】

前記現在の一貫性スナップショットを永続的に記憶するためのコミットを前記開始することは、一貫性スナップショットを作成する定期的作業に応答するものである、請求項 6 に記載の方法。

【請求項 11】

不揮発性機械可読媒体と、

揮発性機械可読媒体と、

プロセッサと、

前記プロセッサ上で実行するように作動可能なファイルセット・マネージャと、
を含む装置であって、

30

前記ファイルセット・マネージャは、

リダイレクト・オン・ライト・ファイル・システム内の複数のデータ・オブジェクトの現在の一貫性スナップショットを前記不揮発性機械可読媒体内に記憶するためのコミットを開始することであって、前記複数のデータ・オブジェクトの各々は、異なるコンテキストを有する、前記複数のデータ・オブジェクトのデータの複数のコピーを有するように構成可能であり、前記複数のデータ・オブジェクトの各々は、前記異なるコンテキストのうちの第 1 のコンテキストを有する、前記データの第 1 のコピーのうちの第 1 のコピーを有し、前記複数のコピーのうちの第 1 のコピーは前記揮発性機械可読媒体内に記憶されるように構成される、前記開始すること、

40

前記現在の一貫性スナップショットを記憶するためにコミットしている最中に、前記複数のデータ・オブジェクトのうちのデータ・オブジェクトに対する更新を受信すること、ならびに

前記データ・オブジェクトに対する前記更新の受信に応答して、

前記第 1 のコピーから前記揮発性機械可読媒体内に前記データ・オブジェクトのデータの第 2 のコピーを作成することであって、前記データの第 2 のコピーは前記異なるコンテキストのうちの第 2 のコンテキストを有する、前記作成すること、および

前記更新に基づき前記データ・オブジェクトの前記データの第 2 のコピーを更新すること、

50

を実行するように構成される、
装置。

【請求項 1 2】

前記ファイルセット・マネージャは、他の一貫性スナップショットの世代値に対して一意的である世代値を、前記現在の一貫性スナップショットに付与するように構成される、請求項 1 1 に記載の装置。

【請求項 1 3】

前記ファイルセット・マネージャは、前記データ・オブジェクトに対する前記更新の受信に応答して、

前記データ・オブジェクトのための世代値をインクリメントすること、および

前記データ・オブジェクトの前記世代値から導き出される世代値を前記更新に関連付けること、

をするように構成される、請求項 1 2 に記載の装置。

【請求項 1 4】

前記現在の一貫性スナップショットは、前の一貫性スナップショットより後の、前記複数のデータ・オブジェクトに対する更新を含む、請求項 1 1 に記載の装置。

【請求項 1 5】

前記現在の一貫性スナップショットを記憶するための前記コミットの前記開始は、一貫性スナップショットを作成する定期的作業に応答するものである、請求項 1 1 に記載の装置。

【請求項 1 6】

複数のデータ・オブジェクトのうちのデータ・オブジェクトのための複数のコンテキストを提供するコンピュータ・プログラムであって、コンピュータに、

リダイレクト・オン・ライト・ファイル・システム内の前記複数のデータ・オブジェクトの現在の一貫性スナップショットを不揮発性機械可読媒体内に記憶するためのコミットを開始することであって、前記複数のデータ・オブジェクトの各々は、コミット中のコンテキストを有する、前記複数のデータ・オブジェクトのデータの第 1 のコピーを有する、前記開始すること、

他の一貫性スナップショットの世代値に対して一意的である世代値を、前記現在の一貫性スナップショットに付与すること、

前記現在の一貫性スナップショットを記憶するために前記コミットしている最中に、前記複数のデータ・オブジェクトのうちの前記データ・オブジェクトに対する更新を受信すること、

前記データ・オブジェクトに対する前記更新を受信することに応答して、

前記データ・オブジェクトのための世代値をインクリメントすること、

前記データ・オブジェクトの前記世代値から導き出される世代値を前記更新に関連付けること、

前記データ・オブジェクトの前記データの第 1 のコピーからコピーされる、前記データ・オブジェクトのデータの第 2 のコピーを作成することであって、前記データ・オブジェクトの前記データの第 2 のコピーは更新中のコンテキストを有する、前記作成すること、および

前記データ・オブジェクトの前記データの第 1 のコピーを更新することとは独立に、前記更新に基づき前記データ・オブジェクトの前記データの第 2 のコピーを更新すること、

を実行させる、コンピュータ・プログラム。

【請求項 1 7】

前記データ・オブジェクトに対する前記更新の受信に応答して、揮発性機械可読媒体内に、前記データ・オブジェクトに関連付けられるバッファ・ヘッダを作成することであって、前記バッファ・ヘッダは、前記データ・オブジェクトの前記データの第 1 のコピーを指し示す第 1 のデータ・ポインタおよび前記データ・オブジェクトの前記データの第 2

10

20

30

40

50

記第 2 のコピーを指し示す第 2 のデータ・ポインタを含む、前記作成することを前記コンピュータに実行させる、請求項 16 に記載のコンピュータ・プログラム。

【請求項 18】

揮発性機械可読媒体内に前記データ・オブジェクトの前記データの前記第 1 のコピーおよび前記データ・オブジェクトの前記データの前記第 2 のコピーを作成することを前記コンピュータに実行させる、請求項 16 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、概してコンピュータの分野に関し、より具体的には、リダイレクト・オン・ライト・ファイル・システムにおけるデータ・バックアップに関する。

10

【背景技術】

【0002】

ファイル・システムは、システム・クラッシュが発生した場合にその内部の一貫性を確保するために種々の方法を用いる。1つのアプローチは、ファイル・システムが数秒ごとにボトム・アップ順でディスク上の新しい位置に変更データを書き込むことである。その内部に記憶されたデータのこれらのビューは一貫性スナップショットと呼ばれる。システム・クラッシュの後、ファイル・システムは、一貫性があることが保証されている、ファイル・システムの最後の一貫性スナップショットの上部から起動する。

20

【発明の概要】

【発明が解決しようとする課題】

【0003】

一貫性スナップショットが書き込まれている間に、ユーザがファイル・システムを新しく変更しようとする可能性がある。一貫性スナップショットが不揮発性機械可読媒体（例えば、ハード・ディスク）上での記憶のためにコミットされるまで、このような変更をブロックすることは容易であろう。しかし、このアプローチは受け入れられるものではない。なぜなら、このようなアプローチは、一貫性スナップショットがユーザに対して非透過的になってしまうからである。特に、このアプローチは、一貫性スナップショットが記憶のためにコミットされている間、全ファイル・システムを数秒ごとにフリーズさせてしまう可能性がある。

30

【課題を解決するための手段】

【0004】

本発明による実施形態には、リダイレクト・オン・ライト・ファイル・システム内の複数のデータ・オブジェクトの現在の一貫性スナップショットを不揮発性機械可読媒体に記憶するためのコミットを開始することであって、複数のデータ・オブジェクトの各々は、コミット中のコンテキストを有する、複数のデータ・オブジェクトのデータの第 1 のコピーを有する、開始することを含む方法が含まれる。方法は、他の一貫性スナップショットの世代値に対して一意的である世代値を、現在の一貫性スナップショットに付与することを含む。方法は、現在の一貫性スナップショットを記憶するためにコミットしている最中に、複数のデータ・オブジェクトのうちのデータ・オブジェクトに対する更新を受信することを含む。方法は、データ・オブジェクトに対する更新を受信することに応答して、データ・オブジェクトのための世代値をインクリメントすることを含む。方法は、データ・オブジェクトに対する更新を受信することに応答して、データ・オブジェクトの世代値から導き出される世代値を更新に関連付けることを含む。同様に、データ・オブジェクトに対する更新を受信することに応答して、方法は、データ・オブジェクトのデータの第 1 のコピーからコピーされる、データ・オブジェクトのデータの第 2 のコピーを作成することを含む。データ・オブジェクトのデータの第 2 のコピーは更新中のコンテキストを有する。同様に、データ・オブジェクトに対する更新を受信することに応答して、方法は、データ・オブジェクトのデータの第 1 のコピーを更新することとは独立に、更新に基づきデータ・オブジェクトのデータの第 2 のコピーを更新することを含む。

40

50

【 0 0 0 5 】

本発明による実施形態には、リダイレクト・オン・ライト・ファイル・システム内の複数のデータ・オブジェクトの現在の一貫性スナップショットを永続的に記憶するためのコミットを開始することであって、複数のデータ・オブジェクトの各々は、異なるコンテキストを有する、複数のデータ・オブジェクトのデータの複数のコピーを有するように構成可能である、開始することを含む方法が含まれる。複数のデータ・オブジェクトの各々は、少なくとも2つのコンテキストのうちの第1のコンテキストを有する、データの少なくとも2つのコピーのうちの第1のコピーを有する。方法は、現在の一貫性スナップショットを記憶するためにコミットしている最中に、複数のデータ・オブジェクトのうちのデータ・オブジェクトに対する更新を受信することを含む。データ・オブジェクトに対する更新の受信に応答して、方法は、第1のコピーからデータ・オブジェクトのデータの第2のコピーを作成することを含む。データの第2のコピーは少なくとも2つのコンテキストのうちの第2のコンテキストを有する。データ・オブジェクトに対する更新の受信に応答して、方法は、更新に基づきデータ・オブジェクトのデータの第2のコピーを更新することを含む。

10

【図面の簡単な説明】

【 0 0 0 6 】

【図1】いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供するクラスタ化ファイル・システム構成の概念図である。

20

【図2】いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供するクラスタ化ファイル・システム構成のより詳細な概念図である。

【図3】いくつかの実施形態による、クラスタ化ファイル・システム内に記憶されるデータ・オブジェクトのためのバッファ・ヘッダの例を示す図である。

【図4】いくつかの実施形態による、データ・オブジェクトの複数の世代に対して一貫性スナップショットをコミットするタイムラインの例を示す図である。

【図5】いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供する作業のフローチャートである。

30

【図6】いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供する作業のフローチャートである。

【図7】コンピュータ・システムの例を示す図である。

【発明を実施するための形態】

【 0 0 0 7 】

添付の図面を参照することによって、本実施形態はより良く理解され、数多くの対象、特徴および利点が当業者に明らかになればよい。

【 0 0 0 8 】

以下に続く記載は、本発明の主題の技法を具体化するシステム、方法、技法、命令シーケンスおよびコンピュータ・プログラムの例を含む。しかし、記載されている実施形態は、これらの具体的詳細を備えずに実施されてもよいことは理解されよう。例えば、例が、ファイル・システムの一部であるデータのためのデュアル・コンテキストに言及していても、いくつかの他の実施形態例はデータのためのコンテキストをいくつかでも（例えば、3つ、4つ、5つ等）構成することができる。同様に、ファイル・システムのための一貫性スナップショットが作成されるように記載される一方で、いくつかの他の実施形態例では、他のレベルにおける一貫性スナップショットが作成されることができる。例えば、ユーザは、ファイルのサブセット、特定のファイル等を、ファイル・システムのための定期スナップショットよりも頻繁に一貫性スナップショット内にバックアップされるように構成することができる。それ故、いくつかの実施形態例はこれらの他のレベルの一貫性スナッ

40

50

ブショットに適用できる。他の例では、記載を分かりにくくしないようにするために、周知の命令インスタンス、プロトコル、構造および技法は詳細に示されていない。

【0009】

多重コンピュータ・システムまたはノード、ならびに永続性記憶資源を含む資源からクラスタが形成される。クラスタの記憶資源にわたってクラスタ化ファイル・システムが実装される。クラスタ記憶資源は、クラスタのノードによる直接アクセスを許すように結合される。記憶資源はノードに直接ケーブルでつながれるか、もしくはネットワーク（例えば、ストレージ・エリア・ネットワーク）を介してアクセス可能とされるか、またはその両方の方法が用いられることができる。

【0010】

クラスタが確立されると、管理者が、クラスタのノードのうちの1つをクラスタ・リーダーとして動作するように構成する。実施形態は、リーダーを自動的に選ぶようにクラスタをプログラムすることもできる。クラスタ・リーダーは、ノードがクライアントなのかサーバなのか、それともクライアントとサーバの両方なのかを示すクラスタ役割データを保持する。サーバがクラスタ化ファイル・システム内のファイルセットを管理する。クラスタ・リーダーは、どのノードがクラスタ化ファイル・システム・マネージャとして動作するかの指示も保持する。クラスタ化ファイル・システム・マネージャはクラスタ化ファイル・システムのためのメタデータを管理する。実施形態によっては、クラスタ化ファイル・システム・マネージャはクラスタのための唯一のサーバである - フェイルオーバー・サーバの責任は負っていない。実施形態によっては、クラスタ化ファイル・システム・マネージャはクラスタ化ファイル・システム内のファイルセットの管理を、サーバである他のノードに任せる。本願明細書において、用語「ファイルセット」はファイルもしくはディレクトリのセットまたはその両方のセットに言及するために用いられる。どのノードがクラスタ内のサーバであるかの指示とともに、クラスタ・リーダーは、サーバまたは「ファイルセット・マネージャ」によって管理されるファイルセットの指示を保持することができる。クラスタ内のノードがクラスタ・リーダーおよびクラスタ化ファイル・システム・マネージャとして動作するように構成されることができる。ノードがクラスタ・リーダーとして動作するのか、サーバとして動作するのか、クライアントとして動作するのか、その他のものとして動作するのかは、クラスタのユーザに透過的であることができる。ノードがクライアントおよびサーバの両方として動作しても、あるいはクライアントが遠隔ノード上にあっても、ユーザは同じ動作を知覚することになる。

【0011】

クラスタ化ファイル・システム・マネージャは、クラスタ化ファイル・システムのファイルのためのiノードの階層としてメタデータを保持することができる。クラスタ化ファイル・システム・メタデータは、クラスタ化された記憶資源の記憶の論理ユニットに関する情報を示す。情報はクラスタ記憶ユニットの位置（例えば、オフセットまたはブロック番号）ならびにエクステントの長さを含むことができる。本記載において、用語「ブロック」は、クラスタ記憶の単位に言及するために用いられる（例えば、4KBブロック）。本記載は、隣接ブロックのセットに言及するために用語「エクステント」も用いる。エクステントの「長さ」に言及するとき、該長さは、エクステントを形成する多数の隣接ブロックに言及している。これらの用語を用いると、4KBブロックを仮定すれば、クラスタ化ファイル・システムは、総計10GBの記憶資源のプールを0ないし2,621,439ブロックと見なす。クラスタ・クライアントがクラスタ記憶の論理ユニットに書き込みをするとき、書き込みを果たすために論理ユニット（例えば、ブロック番号）は記憶仮想化層によって物理的位置（例えば、シークおよびオフセット）に変換される。実施形態はブロックおよびエクステントに限定されるものではないが、クラスタ記憶の単位（例えば、可変長ブロック、ビット長等）についてあり得るあらゆる実装を考慮すれば、記載を分かりにくくしてしまうことになる。

【0012】

実施形態によっては、クラスタ化ファイル・システム・マネージャはクラスタ化ファイ

10

20

30

40

50

ル・システム・メタデータ（「メタデータ」）をiノードの階層データ構造内に保持する。クラスタ化ファイル・システム・マネージャはメタデータのためのルートをクラスタ記憶資源（「クラスタ記憶」）内の既知の位置（すなわち、所定位置）に保持する。一貫性スナップショットを支援するクラスタ内では、一貫性スナップショットのルートを、対応する一貫性スナップショットのルート・メタデータとともに記憶するために、クラスタ記憶内の複数の位置が確保または定義される。ルート・メタデータは、一貫性スナップショットを識別すること、および一貫性スナップショットの保全性を確実にすることを助ける。実施形態は、一貫性スナップショットの進行を追跡する、一貫性スナップショットの時間ベースの識別子（例えば、世代値）、およびデータの保全性を検査するためのルート・チェックサムを用いることができる。実施形態は、ノードがルートを書き込み始める際、第1のルート・チェックサム（「ヘッダ・チェックサム」）を書き込み、ルートが永続性クラスタ記憶にうまく書き込まれた後、第2のルート・チェックサム（「トレーラ・チェックサム」）を書き込むことができる。実施形態はヘッダ・チェックサムおよびトレーラ・チェックサムを用いて、一貫性スナップショットのルートの書き込みが妨げられなかったことを確実にすることができる。障害から回復するためには、位置の各々が調べられ、最新の世代値を持つ位置が選択され、選択された位置によって参照されるその一貫性スナップショットから回復が始まることを可能にする。実施形態は、一貫性スナップショットをいくつでも保存するようにクラスタを構成することができる。

10

【0013】

いくつかの実施形態例は所与のファイル・システム内のデータの一貫性スナップショットを提供するが、このようなスナップショットは、一貫性スナップショットが記憶のためにコミットされている間、入ってくるファイル・システム・トランザクションをブロックしたりまたは遅らせたりしない。それ故、ファイル・システム内に記憶されているデータに対する更新が、同じファイル・システムの一貫性スナップショットの記憶と同時に進行することができる。以下においてさらに記載されるように、この同時実行が可能とされるために、同じデータ・オブジェクトのための少なくとも2つのコンテキストが保持される。

20

【0014】

ファイル・システムの一貫性スナップショットは、一意の世代値に関連付けられる。例えば、世代値は整数値とすることができる。それ故、一貫性スナップショットを記憶するためのコミット（すなわち、該一貫性スナップショットの同期）が始められると、ファイル・システムのための世代はインクリメントされることができる。

30

【0015】

実施形態例によっては、ファイル・システム内のオブジェクト（例えば、データ、ファイル等）に対するあらゆる変更がトランザクションに関連付けられる。トランザクションはファイル・システムの世代に関連付けられ、かくして一貫性スナップショットに関連付けられる。デュアル・コンテキスト構成の実施形態例によっては、ファイル・システム内のオブジェクトが、常に、オブジェクトの最大で2つのコピーとともにキャッシュされる。1つのコピーは、更新中のコンテキストのためのものである。特に、オブジェクトのための更新中のコンテキストは、オブジェクトが更新されている最中であつた（例えば、ユーザがオブジェクトを更新した）後、オブジェクトを有するファイル・システムの一貫性スナップショットが記憶のためにコミットされている間に作成される。オブジェクトの第2のコピーは、コミット中のコンテキストのためのものである。オブジェクトのこのコピーは、一貫性スナップショットの一部として記憶のためにコミットされている最中の/コミットされることになるオブジェクトのコピーである。実施形態例によっては、2つのオブジェクトは、サイズ2の配列として一緒にキャッシュされることができる。同様に、オブジェクトは、配列内の各オブジェクトに関連付けられる世代を記憶する、2つの要素の配列を有する。

40

【0016】

実施形態例によっては、オブジェクトのどのコンテキストを用いるべきかを決定するた

50

めに、トランザクションに関連付けられる世代値または世代番号が、キャッシュされているオブジェクトの世代番号と比較される。正しいオブジェクト配列要素が選択され、オブジェクトの正しいコンテキスト（例えば、更新中のコンテキストまたはコミット中のコンテキスト）を変更する。

【0017】

以下においてさらに記載されるように、一貫性スナップショットは定期的に（例えば、5秒ごとに）作成される。これらの一貫性スナップショットは、以前のバージョンのオブジェクトを回復しようと試みるために作成される。例えば、これらの一貫性スナップショットは、システム・クラッシュの後に、ファイル・システム内に記憶されているオブジェクトを回復するために用いられることができる。実施形態例によっては、以前の一貫性スナップショットがその同期（記憶のためのコミット）を終える前に、一貫性スナップショットの間隔が到達されると、新しい一貫性スナップショットはスキップされる。特に、そのときは、キャッシュされているオブジェクトの第3のコピーが必要となろうから、一貫性スナップショットは取られないであろう。

【0018】

図1は、いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供するクラスタ化ファイル・システム構成の概念図を示す。図示のクラスタはノード103、105、107、109を含む。クラスタは、直接アクセス可能記憶デバイスのプール101、ネットワーク・アクセス可能記憶デバイス113、115、およびネットワーク・インフラストラクチャ111も含む。ノード103、105、107、109はネットワーク・インフラストラクチャ111を介して通信する。ノード103、105、107、109はケーブルを介して記憶デバイス・プール101にアクセスし、ネットワーク・インフラストラクチャ111を介してネットワーク・アクセス可能記憶デバイス113、115にアクセスする。図示のクラスタにおいて、ノード103、105、107、109のうちのいずれでもクラスタのためのクラスタ化ファイル・システム・マネージャとして構成されることができる。クラスタ化ファイル・システム・マネージャはその内部のクラスタ化ファイル・システムのファイルの記憶の様々な側面を管理することができる。例えば、クラスタ化ファイル・システム・マネージャは、クラスタ化ファイル・システムのファイルのためのiノードの階層としてメタデータを保持することができる。実施形態例によっては、クラスタ化ファイル・システム・マネージャの作業のうちの一部またはすべては、異なるノード103、105、107、109に振り分けられることができる。これらの作業のうちの一つは、（以下でさらに記載されるように）クラスタ化ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供することに関連する作業を含む。以下においては、データ・オブジェクトのための複数のコンテキストを提供するこれらの作業は異なるノード103、105、107、109にわたって振り分けられるように記載されているが、一方、いくつかの他の実施形態例では、このような作業はクラスタ化ファイル・システム・マネージャによって遂行されることができる。

【0019】

図2は、いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供するクラスタ化ファイル・システム構成のより詳細な概念図を示す。図2は、図1のノード103、105、107、109を代表することができるノードA 202、ノードB 204およびノードN 206を含むシステム200を示す。図2はノードA 202内の多数の構成要素を示す。図示されていないが、ノードB 204およびノードN 206はそれらの内部に同様の構成要素を含むことができる。

【0020】

実施形態例によっては、システム200は、データが変更される際、リダイレクト・オン・ライト（`redirect-on-write`、ROW）を用いるファイル・システムのデータ・オブジェクトを記憶するように構成される。特に、リダイレクト・オン・ラ

10

20

30

40

50

イトを用いると、変更データのために新しいブロックが割り付けられる。ファイル・システムは1つ以上のファイルセットを含むことができる。実施形態例によっては、ファイル・システム内の各ファイルはiノードを含むことができる。iノードは、ファイル内に記憶されるデータに関する情報またはメタデータを記憶する別個のファイルまたはデータ構造であることができる。例えば、ファイルの部分(例えば、ブロック)ごとに、iノードは、このデータが記憶されるファイルセットのアドレス、ファイルセット識別情報および世代を記憶することができる。特に、ファイルのデータが記憶されるブロックは、異なるファイルセット、およびファイルセットの世代にわたって振り分けられることができる。異なるファイルセット、およびファイルセットの世代は複数の記憶デバイスにわたって振り分けられることができる。図2を参照すると、これらのファイルセットはノードA 202、ノードB 204およびノードN 206のうちのいずれかにおける機械可読媒体内に記憶されること

10

【0021】

システム200は多数のクライアント・デバイス(クライアント・デバイス208およびクライアント・デバイス210として図示)を含む。システム200はネットワーク212を含み、ノードA 202、ノードB 204、ノードN 206、クライアント・デバイス208およびクライアント・デバイス210はネットワーク212を通じて通信可能に共に結合される。

【0022】

ノードA 202は、通信可能に共に結合された、ファイルセット・マネージャ214、不揮発性機械可読媒体216、およびメモリ(例えば、揮発性機械可読媒体)218を含む。ファイルセット・マネージャ214はソフトウェア、ファームウェア、ハードウェアまたはそれらの組み合わせであることができる。例えば、ファイルセット・マネージャ214は、ノードA 202内のプロセッサ(不図示)上で実行するオペレーティング・システムの一部であることができる。不揮発性機械可読媒体216は、すでに作成された多数の一貫性スナップショット(一貫性スナップショットA 224および一貫性スナップショットN 226として図示)を記憶する。不揮発性機械可読媒体216は、その内部に記憶するためにコミットされている途中である現在の一貫性スナップショット228も記憶している最中である。実施形態例によっては、一貫性スナップショットは定期的に(例えば5秒ごとに)作成される。一貫性スナップショットは、所与の時点におけるファイル・システム内のデータ・オブジェクトのスナップショットを含む。実施形態例によっては、一貫性スナップショットは、最後の一貫性スナップショットが記憶のためにコミットされてより後、不揮発性機械可読媒体216に記憶するためにまだコミットされていない、メモリ218内にあるデータ・オブジェクトに対するあらゆる変更(例えば、改変、追加、削除等)を記憶する。これらの一貫性スナップショットは、ファイル・システム内に記憶されている、以前のバージョンのオブジェクトを回復しようと試みるために作成される。例えば、これらの一貫性スナップショットは、システム・クラッシュの後に、ファイル・システム内に記憶されているオブジェクトを回復するために用いられることができる。

20

30

【0023】

メモリ218は多数のバッファ・ヘッダ(バッファ・ヘッダA 220、バッファ・ヘッダN 222等)を記憶する。以下においてさらに記載されるように(図3の記載を参照)、バッファ・ヘッダは、ファイル・システム内に記憶されるデータ・オブジェクトに関する様々なメタデータを記憶する。データ・オブジェクトがアクセスされている、変更されている等の最中であれば、ファイルセット・マネージャ214はメモリ218内のデータ・オブジェクトのためのバッファ・ヘッダを作成する(その中にまだ作成されていなければ)。例えば、ファイルセット・マネージャ214は、現在の一貫性スナップショット228を作成するためにデータ・オブジェクトがアクセスされている最中、何らかのクライアント・デバイス要求に基づきデータ・オブジェクトが変更されている最中などに、バッファ・ヘッダを作成することができる。メモリ218のサイズおよびアクセスされて

40

50

いるデータ・オブジェクトの数に基づき、ファイルセット・マネージャ 2 1 4 は、関連データ・オブジェクトがアクセスされている最中でないバッファ・ヘッダのうちのいくつかをフラッシュするよう求められてよい。それに応じて、ファイルセット・マネージャ 2 1 4 は、データ・オブジェクトのアクセスが生じると、メモリ 2 1 8 内のデータ・オブジェクトのためのバッファ・ヘッダを再作成するよう求められてよい。以下においてさらに記載されるように、バッファ・ヘッダ内のメタデータは、所与のデータ・オブジェクトのために作成されるデータの異なるコピーのためのデータ・ポインタを記憶する。本例では、バッファ・ヘッダ A 2 2 0 は、データの第 1 のコピー 2 5 0 を指し示す第 1 のデータ・ポインタ、およびデータの第 2 のコピー 2 5 2 を指し示す第 2 のデータ・ポインタを有する。メモリ 2 1 8 内に記憶される異なるバッファ・ヘッダのために同様のデータ・ポインタが作成されることができる。

10

【 0 0 2 4 】

実施形態例によっては、ファイル・システム内の同じデータ・オブジェクトのためのデータの複数のコピーが作成される。データの複数のコピーの各々は、異なるコンテキストに関連付けられることができる。実施形態例によっては、データ・オブジェクトが、デュアル・コンテキスト構成のためにそのデータのコピーを 2 つ有することができる。例として、メモリ 2 1 8 は同じデータ・オブジェクトのためのデータの 2 つのコピー（データの第 1 のコピー 2 5 0 およびデータの第 2 のコピー 2 5 2）を記憶する。ファイル・システム内に記憶されるいずれかまたはすべてのデータ・オブジェクトはこのマルチ・コピー、マルチ・コンテキスト構成を含むことができる。図示のように、データの第 1 のコピー 2 5 0 はコミット中のコンテキスト 2 5 4 を有し、データの第 2 のコピー 2 5 2 は更新中のコンテキスト 2 5 6 を有する。同じデータ・オブジェクトのための 2 つのコンテキストはファイル・システム内のデータの一貫性スナップショットを提供するが、このようなスナップショットは、一貫性スナップショットが記憶のためにコミットされている間、入ってくるファイル・システム・トランザクションをブロックしたりまたは遅らせたりしない。それ故、ファイル・システム内に記憶されているデータに対する更新が、同じファイル・システムの一貫性スナップショットの記憶と同時に進行することができる。具体的には、コミット中のコンテキスト 2 5 4 は、現在の一貫性スナップショット 2 2 8 内にこの特定のデータ・オブジェクトを作成するために用いられるデータのコピーに関連付けられる。更新中のコンテキスト 2 5 6 は、現在の一貫性スナップショット 2 2 8 が、（不揮発性機械可読媒体 2 1 6 内に作成される）記憶のためにコミットされている間、データ・オブジェクトに対する更新（例えば、ユーザがデータに変更を行っている）を受け付けるために用いられるデータのコピーに関連付けられる。

20

30

【 0 0 2 5 】

図 2 は多数の作業（作業 2 3 0、作業 2 3 2 および作業 2 3 4）も示す。本例では、ファイルセット・マネージャ 2 1 4 は作業 2 3 0 を遂行し、該作業 2 3 0 において、ファイルセット・マネージャ 2 1 4 は現在の一貫性スナップショットを記憶するためのコミットを開始する。特に、ファイルセット・マネージャ 2 1 4 は現在の一貫性スナップショット 2 2 8 の作成を開始する。作業の一部として、ファイルセット・マネージャ 2 1 4 は、以前の一貫性スナップショットが記憶のためにコミットされて以来、どのようなデータ・オブジェクトが変更されたのかを判定することができる。次に、ファイルセット・マネージャ 2 1 4 は、変更されたデータ・オブジェクトを不揮発性機械可読媒体 2 1 6 内の新しい位置にボトム・アップ順で書き込むことができる。実施形態例によっては、現在の一貫性スナップショット 2 2 8 内に記憶されている最中のデータ・オブジェクトごとに、ファイルセット・マネージャ 2 1 4 はメモリ 2 1 8 内の関連バッファ・ヘッダを作成もしくは更新するまたはその両方を行うことができる（作業 2 3 4 として図示）。メモリ 2 1 8 内にデータ・オブジェクトのための関連バッファ・ヘッダがない場合、ファイルセット・マネージャ 2 1 4 は、このようなデータが、現在の一貫性スナップショット 2 2 8 内に記憶するためにアクセスされている最中に、バッファ・ヘッダを作成する。以下において図 4 ~ 6 を参照してさらに記載されるように、各データ・オブジェクトのためのバッファ・ヘッ

40

50

ダは様々なメタデータ（例えば、世代、コンテキスト、位置、データ・ポインタ）を含む。ファイルセット・マネージャ 214 は、メモリ 218 内にバッファ・ヘッダを作成することの一部としてこのメタデータを更新する。代替的に、バッファ・ヘッダが所与のデータ・オブジェクトのためにメモリ 218 内にすでにインスタンス化されている場合には、ファイルセット・マネージャ 214 はその中のメタデータを更新することができる。例えば、ファイルセット・マネージャ 214 は、データ・ポインタによって参照されている複数のデータのための世代およびコンテキストを定義する様々なフィールドを更新することができる（以下においてさらに記載されるように）。

【0026】

同様に、現在の一貫性スナップショット 228 を記憶するためのコミットの完了より前に、現在の一貫性スナップショット 228 内に含まれるべきデータ・オブジェクトが変更される。本例では、クライアント・デバイス 210 が、現在の一貫性スナップショット 228 の一部であるデータ・オブジェクトのための更新要求を、ネットワーク 212 を通じて送信し、該更新要求はファイルセット・マネージャ 214 によって受信される（作業 232 として図示）。この状況において、ファイルセット・マネージャ 214 は、データの第 1 のコピーからコピーされる、データ・オブジェクト内のデータの第 2 のコピーを作成する（例えばデータの第 1 のコピー 250 およびデータの第 2 のコピー 252 を参照）。同様に、データの第 2 のコピーは、データの第 1 のコピーのために定義されるコンテキストとは別個且つ異なるコンテキストを有する。実施形態例によっては、データの第 2 のコピーは、デュアル・コンテキストを提供するために第 2 のコピーが必要になるまでは作成されない。例えば、同じデータ・オブジェクトを記憶する一貫性スナップショットが作成されている途中である間は、ファイルセット・マネージャ 214 は、データ・オブジェクトに対する更新が要求されるまで第 2 のコピーを作成しない。同様に、ファイルセット・マネージャ 214 はメモリ 218 内のこのデータ・オブジェクトのためのバッファ・ヘッダを作成もしくは更新するまたはその両方を行う。例えば、ファイルセット・マネージャ 214 は、データの第 2 のコピーを指し示すためにバッファ・ヘッダ内の第 2 のデータ・ポインタを更新することができる。同様に、ファイルセット・マネージャ 214 は、データの 2 つの異なるコピーが 2 つの異なるコンテキストを有するように、コンテキストを更新する。ファイルセット・マネージャ 214 の、データ・オブジェクトのための複数のコンテキストを提供する作業のより詳細な記載が、以下において図 5～6 のフローチャートを参照して説明される。

【0027】

図 3 は、いくつかの実施形態による、クラスタ化ファイル・システム内に記憶されるデータ・オブジェクトのためのバッファ・ヘッダの例を示す。バッファ・ヘッダ 300 は、クラスタ化ファイル・システム内に記憶されるデータ・オブジェクトに関連する多数のフィールドを含む。上述されたように、データ・オブジェクトのためのバッファ・ヘッダが、まだそれがメモリ内にない場合に且つ、データ・オブジェクトへのアクセスに応答して、メモリ内に作成される。例えば、ファイルセット・マネージャ 214 は、一貫性スナップショット内のデータ・オブジェクトを記憶するためにデータ・オブジェクトにアクセスすることができる。別の例では、ファイルセット・マネージャ 214 は、データ・オブジェクトを更新する何らかのアプリケーション（例えば、クライアント・デバイス 208、210）に応答してデータ・オブジェクトにアクセスすることができる。バッファ・ヘッダを作成することに加えて、ファイルセット・マネージャ 214 はその内部のフィールド（302～316）にデータを格納することもできる。フィールド 302～304 はこのデータ・オブジェクトのための 2 つの異なる世代値を定義する。最後にコミットされた世代（Last Committed Generation、LCG）フィールド 302 は、このデータ・オブジェクトが一貫性スナップショット内の記憶のためにコミットされた最後の時のこのデータ・オブジェクトのための世代値を定義する。最後に更新された世代（Last Updated Generation、LUG）フィールド 304 は、このデータ・オブジェクトが更新されていた最後の時のこのデータ・オブジェクトのため

10

20

30

40

50

の世代値を定義する。データ・オブジェクトの世代値は、データ・オブジェクトが最初に更新される度に、ただし、データ・オブジェクトが一貫性スナップショットの一部として永続的記憶のためにコミットされるより前に、インクリメントされる。例えば、データ・オブジェクトの現在の世代値が15であるとする。データ・オブジェクトが一貫性スナップショットの一部として永続的記憶のためにコミットされた後に何らかのアプリケーションがデータ・オブジェクトを更新しようとする、と、世代値は16にインクリメントされる。このデータ・オブジェクトのこの世代値は、データ・オブジェクトが一貫性スナップショットの一部として永続的記憶のためにコミットされるまで、16にとどまる。

【0028】

フィールド306～308はこのデータ・オブジェクトのための2つの異なるコンテキスト値を定義する。これらのコンテキスト値は0または1のいずれかに設定される。特に、データ・オブジェクトのためのコンテキストは2つの値の間で反転する（デュアル・コンテキストの一部であるため）。最後にコミットされたコンテキスト（Last Committed Context、LCX）フィールド306は、このデータ・オブジェクトが一貫性スナップショット内の記憶のためにコミットされた最後の時のこのデータ・オブジェクトのためのコンテキストを定義する。最後に更新されたコンテキスト（Last Updated Context、LUX）フィールド308は、このデータ・オブジェクトが更新されていた最後の時のこのデータ・オブジェクトのためのコンテキストを定義する。例えば、データ・オブジェクトが一貫性スナップショットの一部として永続的記憶のためにコミットされた後であるが、データ・オブジェクトに対する更新よりも前に、LCXフィールド306およびLUX308は両方とも同じ値（例えば1）に設定される。その後、何らかのアプリケーションがデータ・オブジェクトを更新しようとする、と、LUXフィールド308は0の値に反転される。その後、このデータ・オブジェクトが再び一貫性スナップショットの一部として永続的記憶のためにコミットされる際、LCXフィールド306は0の値に反転される。フィールド302～308の使用は、以下において図5～6のフローチャートを参照してさらに記載される。

【0029】

物理的位置フィールド310はファイル・システム内のデータ・オブジェクトの物理的位置（例えばブロック番号）を定義する。論理的位置フィールド312は、データ・オブジェクトが記憶される論理的位置を、このデータ・オブジェクトのための関連iノードの位置に基づき定義する。例えば、論理的位置は、このデータ・オブジェクトが記憶される、iノードに加えていくらかのオフセットの物理的位置を含むことができる。

【0030】

データ・ポインタ0・フィールド314は、メモリ218内のデータ・オブジェクトのデータの第1のコピーを指し示している第1のデータ・ポインタ（データ・ポインタ0）を記憶する。データ・ポインタ1・フィールド316は、メモリ218内のデータ・オブジェクトのデータの第2のコピーを指し示している第2のデータ・ポインタ（データ・ポインタ1）を記憶する。上述されたように、データ・オブジェクトのデータの第2のコピーは、データ・オブジェクトのための第2のコンテキストが必要となるまで作成されない。例えば、データ・オブジェクトが一貫性スナップショットの一部として永続的記憶のためにコミットされた後であるが、データ・オブジェクトに対するいかなるその後の更新よりも前においては、データ・オブジェクトのデータのコピーは1つだけ提供されることができる。この状況において、データ・ポインタ0・フィールド314（データの第1のコピーを指し示す）はデータの第1のコピーを指し示し、データ・ポインタ1・フィールド316（データの第2のコピーを指し示す）は位置を指し示さない（例えばヌル（NULL））。データの第2のコピーは、データ・オブジェクトのために第2のコンテキストが必要になった後、データの第1のコピーのコピーから作成される。例えば、データ・オブジェクトが一貫性スナップショット内に記憶されている最中であり、同時に、クライアント・デバイスがデータ・オブジェクトに対する更新を要求しているとする。この状況において、データ・オブジェクトの第2のコピーが作成される。同様に、データ・ポインタ0

10

20

30

40

50

・フィールド 3 1 4 (データの第 1 のコピーを指し示す) はなおデータの第 1 のコピーを指し示し、データ・ポインタ 1・フィールド 3 1 6 (データの第 2 のコピーを指し示す) は、今度は、データ・オブジェクトのデータの第 2 のコピーを指し示すように変更される。フィールド 3 1 4 ~ 3 1 6 の使用は以下において図 5 ~ 6 のフローチャートを参照してさらに記載される。

【0031】

図 4 は、いくつかの実施形態による、データ・オブジェクトの複数の世代に対して一貫性スナップショットをコミットするタイムラインの例を示す。タイムライン 4 0 0 は左から右へ時間が増大する。時点 4 0 2 は、データ・オブジェクトのための世代 N が終了した時刻である。時点 4 0 4 は、同じデータ・オブジェクトのためのより後の世代 (世代 N + 1) が終了した、より後の時刻である。時点 4 0 6 は、同じデータ・オブジェクトのためのより後の世代 (世代 N + 2) が終了した、より後の時刻である。期間 4 0 8 は、(データ・オブジェクトを含む) 一貫性スナップショットが永続的記憶のためにコミットされている期間である。期間 4 0 8 は、世代 N が終了した後の時点 4 0 2 において開始される。コミットの一部として上述されたように、ファイルセット・マネージャ 2 1 4 はデータ・オブジェクトの階層をボトム・アップ順で横断し、子データ・オブジェクトのブロック番号およびチェックサムを収集する。図示のように、期間 4 0 8 内には 2 つのサブ期間 - 期間 4 1 0 および期間 4 1 2 - がある。期間 4 1 0 は、データ・オブジェクトの 1 つのコピーまたはバージョンがメモリ内に存在する期間を含む。例えば、この期間は、データ・オブジェクトが永続的記憶のためにコミットされている最中の時間であって、データ・オブジェクトはまだ変更されていない (例えば、クライアント・デバイス上で実行するアプリケーションによって)、時間を含むことができる。期間 4 1 2 は、データ・オブジェクトの 2 つのコピーまたはバージョンがメモリ内に存在する期間を含む。期間 4 1 2 は、世代 N のための一貫性スナップショットのコミットが依然、行われている間にデータ・オブジェクトが変更されるのに応答して開始される。例えば、この期間は、データ・オブジェクトが永続的記憶のためにコミットされている最中且つ、データ・オブジェクトが変更されている (例えば、クライアント・デバイス上で実行するアプリケーションによって) 最中の時間を含むことができる。換言すると、データ・オブジェクトの第 1 のバージョンは、発行されている最中の世代 N の一貫性スナップショットの一部として存在する。データ・オブジェクトの第 2 のバージョンは、世代 N の一貫性スナップショットの発行の完了より前に現在の世代 N + 1 内でデータ・オブジェクトへの書き込みが行われる場合に備えてもしくはそれが行われる結果、またはその両方の故に存在する。

【0032】

図 5 ~ 6 は、いくつかの実施形態例による、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供する作業のフローチャートを示す。図 5 はフローチャート 5 0 0 を示し、図 6 はフローチャート 6 0 0 を示す。フローチャート 6 0 0 はフローチャート 5 0 0 の続きであり、ポイント A において移行する。フローチャート 5 0 0 ~ 6 0 0 は、分散型構成で行われるものとして記載されており、そこでは、ファイルセット・マネージャ 2 1 4 がその内部の作業を遂行する。他の実施形態例によっては、フローチャート 5 0 0 ~ 6 0 0 の作業は集中型構成で行われ、そこでは、ファイル・システム・マネージャがこのような作業を遂行することができる。フローチャート 5 0 0 ~ 6 0 0 は、データ・オブジェクトのためのデュアル・コンテキストが必要とされる状況の例を示す。特に、この状況例では、特定のデータ・オブジェクト (「データ・オブジェクト A」と呼ばれる) を含む一貫性スナップショットが不揮発性機械可読媒体内の記憶のためにコミットされている最中となっている。なぜなら、データ・オブジェクト A は、前の一貫性スナップショットが記憶のためにコミットされてより後に、変更されているからである。この一貫性スナップショットが記憶のためにコミットされるのと同時に、データ・オブジェクト A をさらに変更する作業がある。例えば、クライアント・デバイス上で実行するアプリケーションがデータ・オブジェクト A を変更することができる。図 1 ~ 3 を参照してフローチャート 5 0 0 ~ 6 0 0 の作業が記載される。フロー

チャート 5 0 0 がまず記載され、その後にフローチャート 6 0 0 の記載が続く。

【 0 0 3 3 】

ファイルセット・マネージャ 2 1 4 が、ファイル・システム内の多数のデータ・オブジェクトを含む現在の一貫性スナップショットを不揮発性機械可読媒体に記憶するためのコミットを開始する (5 0 2)。実施形態例によっては、ファイルセット・マネージャ 2 1 4 は現在の一貫性スナップショットを記憶するために定期的にコミットする (例えば、3 秒、5 秒、1 0 秒等)。それ故、この作業は一貫性スナップショットを作成する定期的作業のうちの 1 つであることができる。図 2 を参照すると、ファイルセット・マネージャ 2 1 4 は現在の一貫性スナップショット 2 2 8 を記憶するためのコミットを開始する。実施形態例によっては、現在の一貫性スナップショット 2 2 8 は、前の一貫性スナップショットより後に変更されたデータ・オブジェクトを含むことになる。データ・オブジェクトに対するそれらの改変はメモリ 2 1 8 内に常駐することができ、そのため、改変は不揮発性機械可読媒体 2 1 6 内の記憶のためにまだコミットされていない。フローチャート 5 0 0 の作業は 5 0 4 において続く。

【 0 0 3 4 】

ファイルセット・マネージャ 2 1 4 は、現在の一貫性スナップショット内に記憶されるべきデータ・オブジェクトのためのバッファ・ヘッダがメモリ内にあるかどうかを判定する (5 0 4)。図 2 を参照すると、ファイルセット・マネージャ 2 1 4 は、現在の一貫性スナップショット 2 2 8 内に記憶されるべきデータ・オブジェクトのためのバッファ・ヘッダがメモリ 2 1 8 内にあるかどうかを判定する。特に、実施形態例によっては、データ・オブジェクトがアクセスされる (読み出される、書き込まれる等) 度に、メモリ 2 1 8 内に関連バッファ・ヘッダが作成される。現在の一貫性スナップショット 2 2 8 内に記憶されるべきデータ・オブジェクトごとのバッファ・ヘッダがメモリ内にすでにある場合は、フローチャート 5 0 0 の作業は 5 0 8 において続く。さもなければ、フローチャート 5 0 0 の作業は 5 0 6 において続く。

【 0 0 3 5 】

ファイルセット・マネージャ 2 1 4 は、(メモリ内にバッファ・ヘッダをまだ有していないデータ・オブジェクトのために) メモリ内にバッファ・ヘッダを作成し更新する (5 0 6)。図 2 を参照すると、ファイルセット・マネージャ 2 1 4 は、メモリ内にバッファ・ヘッダを有していないこれらのデータ・オブジェクトのためにメモリ 2 1 8 内にバッファ・ヘッダを作成する。ファイルセット・マネージャ 2 1 4 はバッファ・ヘッダのフィールドを更新することもできる。図 3 を参照すると、ファイルセット・マネージャ 2 1 4 はこれらのデータ・オブジェクトごとのバッファ・ヘッダのためのこれらのフィールドの値を設定する。ファイルセット・マネージャ 2 1 4 は L C G フィールド 3 0 2 および L U G フィールド 3 0 4 の両方をデータ・オブジェクトのための現在の世代値に設定する。例えば、もし最後にコミットされた一貫性スナップショットが 5 の値を有していれば、ファイルセット・マネージャ 2 1 4 は L C G フィールド 3 0 2 および L U G フィールド 3 0 4 を 5 に設定することになる。コンテキスト・フィールド (3 0 6、3 0 8) は、2 つのコンテキスト (コミット中のコンテキストおよび更新中のコンテキスト) を区別するために 0 または 1 のいずれかに設定される。従って、第 2 のコンテキストが必要とされる場合は、これらの 2 つのコンテキスト・フィールド 3 0 6、3 0 8 は逆の値を有することになる。1 つのコンテキストだけが必要とされる場合は、これらの 2 つのコンテキスト・フィールド 3 0 6、3 0 8 は同じ値を有することになる。この状況においては、データ・オブジェクトのための 1 つのコンテキストだけが必要である。従って、ファイルセット・マネージャ 2 1 4 は L C X フィールド 3 0 6 および L U X フィールド 3 0 8 を同じ値 (例えば 1) に設定する。ファイルセット・マネージャ 2 1 4 は、ファイル・システム内のデータ・オブジェクトの位置 (例えばブロック番号) に基づき物理的位置フィールド 3 1 0 を設定する。ファイルセット・マネージャ 2 1 4 は、このデータ・オブジェクトのための関連 i ノードの位置に基づき論理的位置フィールド 3 1 2 を設定する。例えば、論理的位置は、このデータ・オブジェクトが記憶される、i ノードに加えていくらかのオフセットの

10

20

30

40

50

物理的位置を含むことができる。ファイルセット・マネージャ 2 1 4 は、データの第 1 のコピーが配置されるメモリ 2 1 8 内の位置を指し示すためにバッファ・ヘッダ 3 0 0 内のデータ・ポインタ 0・フィールド 3 1 4 を更新する。この状況は複数のコンテキストを要求していないので、第 2 のデータ・オブジェクトは必要ない。それ故、ファイルセット・マネージャ 2 1 4 はデータ・ポインタ 1・フィールド 3 1 6 を、ヌルを指し示すように更新する。フローチャート 5 0 0 の作業は 5 0 8 において続く。

【 0 0 3 6 】

ファイルセット・マネージャ 2 1 4 は、ファイル・システム内のデータ・オブジェクト A (現在の一貫性スナップショットの一部となるべきデータ・オブジェクトの一部である) を更新するトランザクションを受信する (現在の一貫性スナップショットを記憶するためのコミットがまだ行われている間に) (5 0 8)。図 2 を参照すると、ファイルセット・マネージャ 2 1 4 は、クライアント・デバイス 2 0 8、2 1 0 のうちの 1 つから、データ・オブジェクト A を更新するトランザクションを受信する。例えば、クライアント・デバイス 2 0 8、2 1 0 のうちの 1 つの上で実行するアプリケーションがデータ・オブジェクト A を更新することができる。フローチャート 5 0 0 の作業は 5 1 0 において続く。

10

【 0 0 3 7 】

ファイルセット・マネージャ 2 1 4 は、メモリ内にデータ・オブジェクト A のためのバッファ・ヘッダがあるかどうかを判定する (5 1 0)。図 2 を参照すると、ファイルセット・マネージャ 2 1 4 は、メモリ 2 1 8 内にデータ・オブジェクト A のためのバッファ・ヘッダがあるかどうかを判定する。特に、実施形態例によっては、データ・オブジェクトがアクセスされる (読み出される、書き込まれる等) 度に、メモリ 2 1 8 内に関連バッファ・ヘッダが作成される。メモリ 2 1 8 内にデータ・オブジェクト A のためのバッファ・ヘッダがすでにある場合は、フローチャート 5 0 0 の作業は継続点 A において続く (5 1 8)。さもなければ、フローチャート 5 0 0 の作業は 5 1 2 において続く。

20

【 0 0 3 8 】

ファイルセット・マネージャ 2 1 4 はデータ・オブジェクト A のためにメモリ内にバッファ・ヘッダを作成する (5 1 2)。図 2 を参照すると、ファイルセット・マネージャ 2 1 4 はデータ・オブジェクト A のためにメモリ 2 1 8 内にバッファ・ヘッダを作成する。なぜなら、メモリ 2 1 8 内にはデータ・オブジェクト A のための関連バッファ・ヘッダがないからである。ファイルセット・マネージャ 2 1 4 はバッファ・ヘッダのフィールドにデータを格納することもできる (以下の作業によってさらに記載されるように)。フローチャート 5 0 0 の作業は 5 1 4 において続く。

30

【 0 0 3 9 】

ファイルセット・マネージャ 2 1 4 はデータ・オブジェクト A のためのバッファ・ヘッダ内のデータ・ポインタ 0・フィールドを更新する (5 1 4)。図 2 ~ 3 を参照すると、ファイルセット・マネージャ 2 1 4 は、データの第 1 のコピーが配置されるメモリ 2 1 8 内の位置を指し示すために、バッファ・ヘッダ 3 0 0 内のデータ・ポインタ 0・フィールド 3 1 4 を更新する。フローチャート 5 0 0 の作業は 5 1 6 において続く。

【 0 0 4 0 】

ファイルセット・マネージャ 2 1 4 はデータ・オブジェクト A のためのバッファ・ヘッダ内の物理的位置、論理的位置、L C G フィールドおよび L C X フィールドも更新する。図 2 ~ 3 を参照すると、ファイルセット・マネージャ 2 1 4 はバッファ・ヘッダ 3 0 0 のための物理的位置フィールド 3 1 0、論理的位置フィールド 3 1 2、L C G フィールド 3 0 2、L C X フィールド 3 0 6 を更新する。ファイルセット・マネージャ 2 1 4 は、ファイル・システム内のデータ・オブジェクトの位置 (例えばブロック番号) に基づき物理的位置フィールド 3 1 0 を設定する。ファイルセット・マネージャ 2 1 4 は、このデータ・オブジェクトのための関連 i ノードの位置に基づき論理的位置フィールド 3 1 2 を設定する。例えば、論理的位置は、このデータ・オブジェクトが記憶される、i ノードに加えていくらかのオフセットの物理的位置を含むことができる。ファイルセット・マネージャ 2 1 4 は L C G フィールド 3 0 2 をデータ・オブジェクト A のための現在の世代値に設定す

40

50

る。例えば、もし最後にコミットされた一貫性スナップショットが5の値を有していれば、ファイルセット・マネージャ214はLCGフィールド302を5に設定することになる。コンテキスト・フィールド(306、308)は、2つのコンテキスト(コミット中のコンテキストおよび更新中のコンテキスト)を区別するために0または1のいずれかに設定される。従って、第2のコンテキストが必要とされる場合は、これらの2つのコンテキスト・フィールド306、308は逆の値を有することになる。1つのコンテキストだけが必要とされる場合は、これらの2つのコンテキスト・フィールド306、308は同じ値を有することになる。ファイルセット・マネージャ214はLCXフィールド306を1に設定すると仮定する。LUXフィールド308の設定は以下において記載される。フローチャート500の作業は継続点A(518)において続く。

10

【0041】

継続点A(518)はフローチャート600の継続点A(602)において続く。継続点A(602)より、作業は603において続く。

【0042】

ファイルセット・マネージャ214は、データ・オブジェクトAのためのバッファ・ヘッダ内のLCGフィールドまたはLUGフィールドの値がトランザクションの世代値と一致するかどうかを判定する(603)。図2~3を参照すると、ファイルセット・マネージャ214は、バッファ・ヘッダ300内のLCGフィールド302の値またはLUGフィールド304の値がトランザクションの世代値と一致するかどうかを判定する。トランザクションの世代値は、トランザクションがいつ作成されたかに基づき一貫性世代に設定される。従って、ファイルセット・マネージャ214は、トランザクションに関連付けられるこの世代が、最後にコミットされた世代または最後に更新された世代に等しいかどうかを判定する。一致がなければ、作業は604において続く。さもなければ、作業は616において続く(以下においてさらに記載される)。

20

【0043】

ファイルセット・マネージャ214はデータ・オブジェクトAの第1のコピーからデータ・オブジェクトAの第2のコピーを作成する(604)。図2を参照し、データの第1のコピー250はデータ・オブジェクトAの第1のコピーであるとする、ファイルセット・マネージャ214はデータの第1のコピー250をメモリ218内の異なる位置-データの第2のコピー252-にコピーする。フローチャート600の作業は606において続く。

30

【0044】

ファイルセット・マネージャ214は、データ・オブジェクトAの第2のコピーを指し示すようにバッファ・ヘッダ内の第2のデータ・ポインタを更新する(606)。図2~3を参照すると、ファイルセット・マネージャ214は、メモリ218内のデータ・オブジェクトAの第2のコピーを指し示すようにデータ・ポインタ1・フィールド316を更新する。フローチャート600の作業は608において続く。

【0045】

ファイルセット・マネージャ214はバッファ・ヘッダ内のLUXフィールドを、LCXフィールドの値とは逆の値を有するように更新する(608)。図2~3を参照すると、ファイルセット・マネージャ214はLUXフィールド308を、バッファ・ヘッダ300内のLCXフィールド306の値とは逆である値を有するように更新する。上述されたように、LCXフィールド306およびLUXフィールド308の値は2つの値のうちの1つであることができる。デュアル・コンテキスト状況が生じると(この場合のように)、LCXフィールド306およびLUXフィールド308の値は互いの逆になる。フローチャート600の作業は610において続く。

40

【0046】

ファイルセット・マネージャ214は、トランザクションのための世代値に基づきバッファ・ヘッダ内のLUGフィールドのための世代値を設定する(610)。図2~3を参照すると、ファイルセット・マネージャ214は、トランザクションのための世代値に基

50

づき LUG フィールド 304 のための世代値を更新する（上記の 603 の記載における
ランザクションのための世代値の記載を参照）。フローチャート 600 の作業は 614 に
おいて続く。

【0047】

ファイルセット・マネージャ 214 は、このランザクションに基づきデータ・オブジ
ェクト A の第 2 のコピーを更新する（614）。図 2～3 を参照し、データの第 2 のコピ
ー 252 はデータ・オブジェクト A の第 2 のコピーであるとする、ファイルセット・マ
ネージャ 214 は、データ・ポインタ 1・フィールド 316 内のポインタ値に基づきデー
タの第 2 のコピー 252 を更新する。フローチャート 600 の作業は、フローチャート 6
00 のこの経路に沿い完結している。

10

【0048】

603 に戻り、一致があるとする（はいの判定）、ファイルセット・マネージャ 21
4 は、データ・オブジェクト A のためのバッファ・ヘッダ内の LUX フィールドに関連付
けられる第 1 のデータ・ポインタを用いてデータ・オブジェクト A のコピーを更新する（
616）。この状況においては、ランザクションのための世代が LUG フィールド 30
4 と一致するであろうから、603 において一致があった。図 2～3 を参照し、第 1 のデ
ータ・ポインタがデータの第 1 のコピー 250 を指し示すとする、ファイルセット・マ
ネージャ 214 は、データ・ポインタ 0・フィールド 314 内のポインタ値に基づきデー
タの第 1 のコピー 250 を更新する。フローチャート 600 の作業は、フローチャート 6
00 のこの経路に沿い完結している。

20

【0049】

ファイル・システム内の同じまたは異なるデータ・オブジェクトに対する追加の更新が
引き続き行われることができる。同様に、一貫性スナップショットのコミットの完了後、
ファイルセット・マネージャ 214 は、（一貫性スナップショットを永続的に記憶するた
めにコミットするための定期的な間隔に基づき）追加の一貫性スナップショットをコミッ
トすることができる。

【0050】

当業者によって理解されるように、本発明の主題の態様はシステム、方法またはコンピ
ュータ・プログラムとして具体化されればよい。それ故、本発明の主題の態様は、完全
にハードウェアの実施形態、完全にソフトウェアの実施形態（ファームウェア、常駐ソフト
ウェア、マイクロ・コード等を含む）あるいはソフトウェアおよびハードウェアの態様を
組み合わせた実施形態という形をとってよく、本願明細書においてそれらはすべて広く「
回路」、「モジュール」または「システム」と呼ばれればよい。さらに、本発明の主題の
態様は、1 つ以上のコンピュータ可読媒体（単数または複数）であって、その上に具体化
されるコンピュータ可読プログラム・コードを有する、コンピュータ可読媒体内に具体化
されるコンピュータ・プログラムという形をとってもよい。

30

【0051】

1 つ以上のコンピュータ可読媒体（単数または複数）の任意の組み合わせが利用されて
よい。コンピュータ可読媒体はコンピュータ可読信号媒体またはコンピュータ可読記憶媒
体であればよい。コンピュータ可読記憶媒体は、例えば、電子的、磁氣的、光学的、電磁
的、赤外線または半導体システム、装置またはデバイス、あるいは上述のものの任意の適
当な組み合わせであればよい。ただし、それらに限定されるものではない。コンピュータ
可読記憶媒体のより具体的な例（限定的なリスト）としては以下のもの：1 本以上のワイ
ヤを有する電気接続、ポータブル・コンピュータ・ディスク、ハード・ディスク、ラン
ダム・アクセス・メモリ（random access memory、RAM）、リ
ード・オンリー・メモリ（read-only memory、ROM）、消去可能プロ
グラマブル・リード・オンリー・メモリ（erasable programmable
read-only memory、EPROM またはフラッシュ・メモリ）、光ファ
イバ、ポータブル・コンパクト・ディスク・リード・オンリー・メモリ（compact
disc read-only memory、CD-ROM）、光学式記憶デバイス

40

50

、磁気記憶デバイス、あるいは上述のものの任意の適当な組み合わせ、が挙げられよう。本文書の文脈において、コンピュータ可読記憶媒体とは、命令実行システム、装置またはデバイスによって用いられるまたはそれらと連係して用いられるプログラムを包含または記憶することができる任意の有形媒体であればよい。

【0052】

コンピュータ可読信号媒体とは、例えば、基底帯域内にまたは搬送波の一部として、内部にコンピュータ可読プログラム・コードが具体化される伝搬データ信号を含むものであればよい。このような伝搬信号は、電磁氣的、光学的、またはそれらの任意の適当な組み合わせを含む、ただしそれらに限定されるものではない、様々な形態のいずれを取ってもよい。コンピュータ可読信号媒体とは、コンピュータ可読記憶媒体ではない、命令実行システム、装置またはデバイスによって用いられるまたはそれらと連係して用いられるプログラムを伝達、伝搬または輸送することができる任意のコンピュータ可読媒体であればよい。

10

【0053】

コンピュータ可読媒体上に具体化されるプログラム・コードは、無線、有線、光ファイバケーブル、RF等、または上述のものの任意の適当な組み合わせを含む、ただしそれらに限定されるものではない、任意の適切な媒体を用いて送信されればよい。

【0054】

本発明の主題の態様のための作業を行うためのコンピュータ・プログラム・コードは、Java(R)、Smalltalk(R)、C++または同様のもの等のオブジェクト指向プログラミング言語、ならびに「C」プログラミング言語または同様のプログラミング言語等の従来の手続き型プログラミング言語を含む、1つ以上のプログラミング言語の任意の組み合わせで書き込まれればよい。プログラム・コードは、スタンド・アロン・ソフトウェア・パッケージとして完全にまたは一部分はユーザのコンピュータ上で実行するか、一部分はユーザのコンピュータ上で且つ一部分はリモート・コンピュータ上で実行するか、または完全にリモート・コンピュータもしくはサーバ上で実行すればよい。後者のシナリオでは、リモート・コンピュータは、ローカル・エリア・ネットワーク(local area network、LAN)またはワイド・エリア・ネットワーク(wide area network、WAN)を含む、任意の種類のネットワークを通じてユーザのコンピュータに接続されてもよいし、あるいは接続は、(例えば、インターネット・サービス・プロバイダを利用しインターネットを通じて)外部のコンピュータになされてもよい。

20

30

【0055】

本発明の主題の実施形態による方法、装置(システム)およびコンピュータ・プログラムのフローチャート図もしくはブロック図またはその両方を参照しながら本発明の主題の態様が記載されている。フローチャート図もしくはブロック図またはその両方の各ブロック、ならびにフローチャート図もしくはブロック図またはその両方におけるブロックの組み合わせは、コンピュータ・プログラム命令によって実装されることが理解されよう。これらのコンピュータ・プログラム命令は、汎用コンピュータ、専用コンピュータ、または他のプログラム可能なデータ処理装置のプロセッサに提供されて機械を作り出せばよく、それにより、命令は、コンピュータまたは他のプログラム可能なデータ処理装置のプロセッサを介して実行し、フローチャートもしくはブロック図またはその両方のブロックまたはブロック群において特定される機能群/動作群を実装する手段を生み出す。

40

【0056】

これらのコンピュータ・プログラム命令は、コンピュータ、他のプログラム可能なデータ処理装置または他のデバイスを特定の様式で機能するように仕向けることができるコンピュータ可読媒体内に記憶されてもよく、それにより、コンピュータ可読媒体内に記憶された命令は、フローチャートもしくはブロック図またはその両方のブロックまたはブロック群において特定される機能/動作を実装する命令を含む製造品を作り出す。

50

【 0 0 5 7 】

コンピュータ・プログラム命令は、コンピュータ、他のプログラム可能なデータ処理装置または他のデバイス上にロードされ、一連の作業ステップをコンピュータ、他のプログラム可能な装置または他のデバイス上で遂行させ、コンピュータ実装プロセスを作り出し、それにより、コンピュータまたは他のプログラム可能な装置上で実行する命令は、フローチャートもしくはブロック図またはその両方のブロックまたはブロック群において特定される機能群 / 動作群を実装するためのプロセス群を提供する。

【 0 0 5 8 】

図 7 はコンピュータ・システムの例を示す。コンピュータ・システムはプロセッサ・ユニット 7 0 1 (場合により、複数のプロセッサ、複数のコア、複数のノードを含む、もしくはマルチ・スレッド等を実装する、またはその両方の態様を有する) を含む。コンピュータ・システムはメモリ 7 0 7 を含む。メモリ 7 0 7 はシステム・メモリ (例えば、キャッシュ、S R A M、D R A M、ゼロ・キャパシタ R A M、ツイン・トランジスタ R A M、e D R A M、E D O R A M、D D R R A M、E E P R O M、N R A M、R R A M、S O N O S、P R A M 等のうちの 1 つ以上) であってもよいし、または先にすでに記載された、機械可読媒体の実現可能なもののうちのいずれか 1 つ以上であってもよい。コンピュータ・システムは、バス 7 0 3 (例えば、P C I、I S A、P C I - E x p r e s s、H y p e r T r a n s p o r t (R)、I n f i n i B a n d (R)、N u B u s 等)、ネットワーク・インターフェース 7 0 5 (例えば、A T M インターフェース、イーサネット (R) ・インターフェース、フレーム・リレー・インターフェース、S O N E T インターフェース、無線インターフェース等)、および記憶デバイス (単数または複数) 7 0 9 (例えば、光学的記憶、磁氣的記憶等) も含む。コンピュータ・システムは、リダイレクト・オン・ライト・ファイル・システム内のデータ・オブジェクトのための複数のコンテキストを提供するファイルセット・マネージャ 7 2 5 も含む。これらの機能性のうちのいずれか 1 つは一部 (または完全に)、ハードウェアの形で、もしくはプロセッサ・ユニット 7 0 1 上で、またはその両方の態様で実装されてよい。例えば、該機能性は、特定用途向け集積回路を用いて、プロセッサ・ユニット 7 0 1 内に実装される論理で、周辺デバイスまたはカード上のコプロセッサ内に、等の態様で実装されればよい。さらに、現実のものは、含まれる構成要素がもっと少なくてもよいし、または図 7 には示されていない追加の構成要素 (例えば、ビデオ・カード、オーディオ・カード、追加のネットワーク・インターフェース、周辺デバイス、等) を含んでもよい。プロセッサ・ユニット 7 0 1、記憶デバイス (単数または複数) 7 0 9 およびネットワーク・インターフェース 7 0 5 はバス 7 0 3 に結合される。バス 7 0 3 に結合されるように示されているが、メモリ 7 0 7 はプロセッサ・ユニット 7 0 1 に結合されてもよい。

【 0 0 5 9 】

実施形態は、様々な実装および利用を参照しながら記載されているが、これらの実施形態は例示的なものであること、および本発明の主題の範囲はそれらに限定されないことは理解されよう。多くの変形、変更、追加および改良が可能である。

【 0 0 6 0 】

本願明細書では単一の例として記載されている構成要素、作業または構造に複数の例が提供されてもよい。最後に、様々な構成要素同士、作業同士およびデータ格納同士の間の境界は、いくぶんは自由に決めてよいものであり、特定の作業は特定の例示的構成を背景として示されている。機能性の他の割り振りが想定され、本発明の主題の範囲に含まれてもよい。一般に、構成例において別個の構成要素として提示されている構造および機能性は、組み合わせられた構造または構成要素として実装されてもよい。同様に、単一の構成要素として提示されている構造および機能性は、別個の構成要素として実装されてもよい。これらおよび他の変形、変更、追加および改良は本発明の主題の範囲に含まれてよい。

【 符号の説明 】

【 0 0 6 1 】

2 0 0 システム

10

20

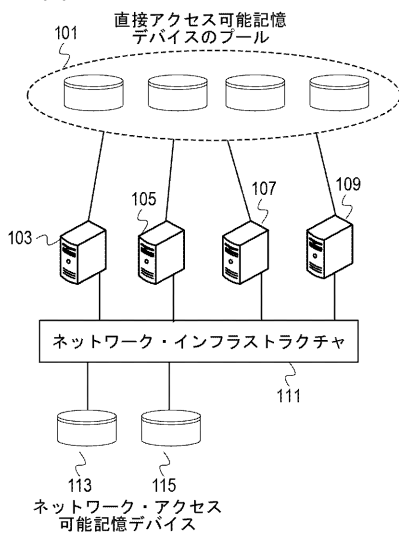
30

40

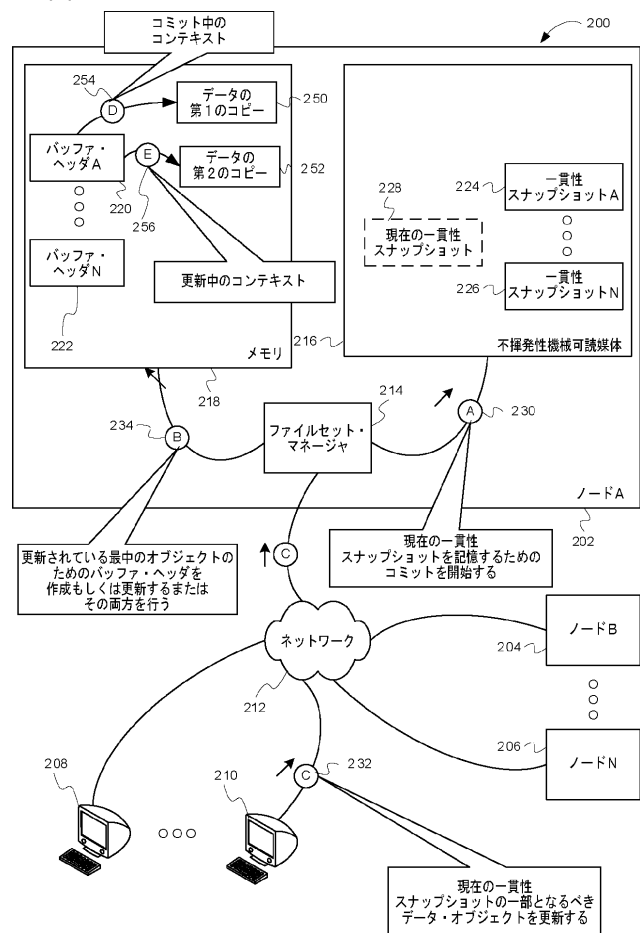
50

- 202 ノード A
- 204 ノード B
- 206 ノード N
- 208、210 クライアント・デバイス
- 212 ネットワーク
- 214 ファイルセット・マネージャ
- 216 不揮発性機械可読媒体
- 218 メモリ
- 220 バッファ・ヘッダ A
- 222 バッファ・ヘッダ N
- 224 一貫性スナップショット A
- 226 一貫性スナップショット N
- 228 現在の一貫性スナップショット
- 230、232、234 作業
- 250 データの第 1 のコピー
- 252 データの第 2 のコピー
- 254 コミット中のコンテキスト
- 256 更新中のコンテキスト

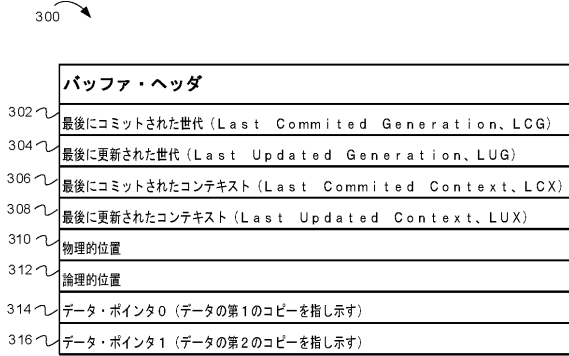
【図 1】



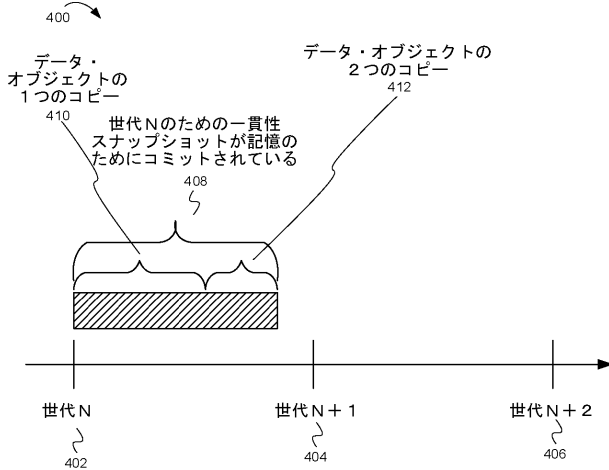
【図 2】



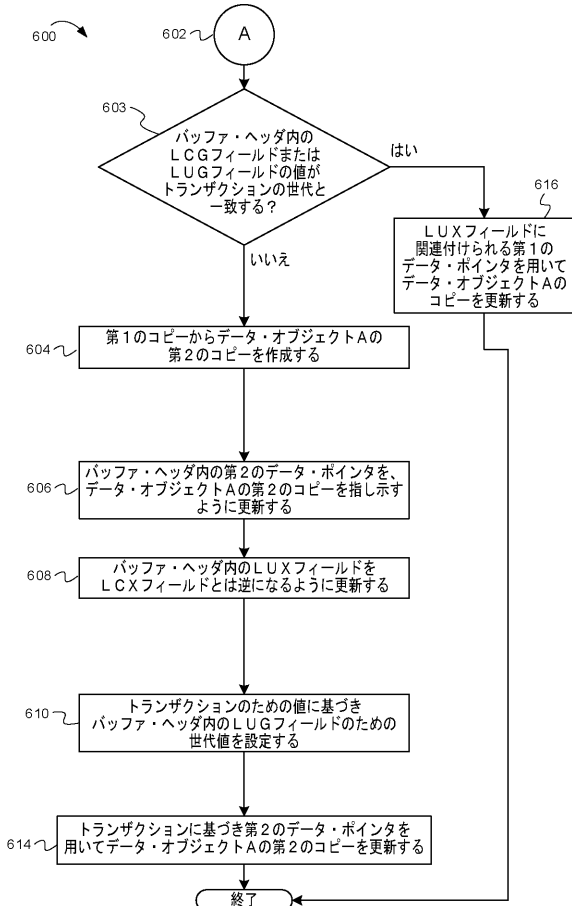
【図 3】



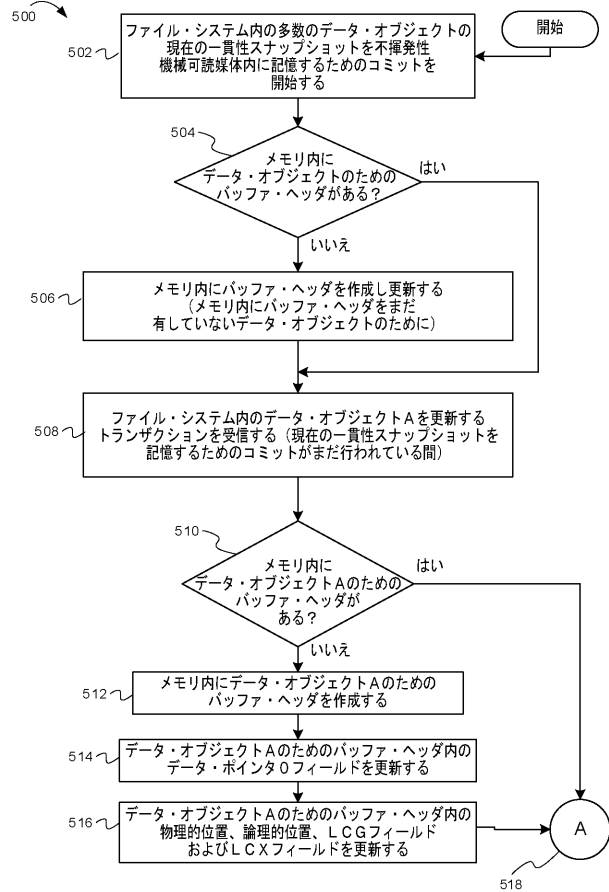
【図 4】



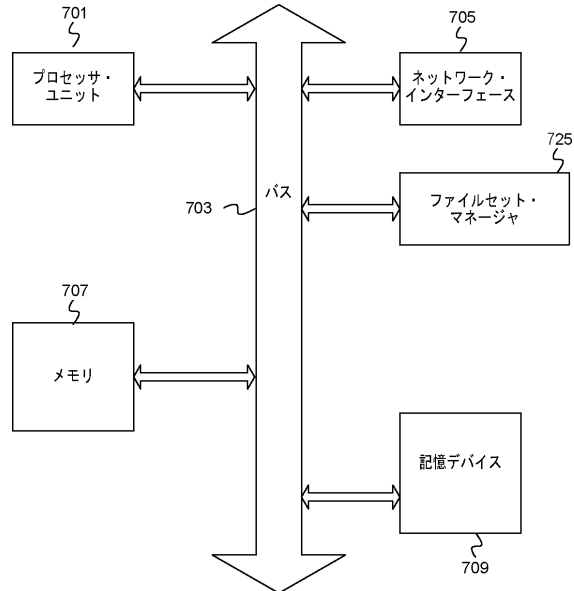
【図 6】



【図 5】



【図 7】



フロントページの続き

- (72)発明者 アンドリュー・ソロモン
アメリカ合衆国 7 8 7 5 8 テキサス州 オースティン バーネット ロード 1 1 5 0 1 エム
ディー 9 5 4 1
- (72)発明者 デイビット・ジョーンズ・クラフト
アメリカ合衆国 7 8 7 5 8 テキサス州 オースティン バーネット ロード 1 1 5 0 1 ジッ
プ 9 0 5
- (72)発明者 マノージ・エヌ・クマール
アメリカ合衆国 7 8 7 5 8 - 3 4 0 0 テキサス州 オースティン バーネット ロード 1 1 5
0 1 4 エイチ - 0 0 4 9 0 5
- (72)発明者 ジャネット・エリザベス・アドキンス
アメリカ合衆国 7 8 7 5 8 - 3 4 0 0 テキサス州 オースティン バーネット ロード 1 1 5
0 1
- (72)発明者 ジュン・チャン
アメリカ合衆国 7 8 7 3 3 テキサス州 オースティン ノース ウェストン エルエヌ 9 0 1