

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6408568号
(P6408568)

(45) 発行日 平成30年10月17日 (2018.10.17)

(24) 登録日 平成30年9月28日 (2018.9.28)

(51) Int. Cl.

F I

G 0 6 F 3/06 (2006.01)
 G 0 6 F 12/00 (2006.01)
 G 0 6 F 13/14 (2006.01)
 G 0 6 F 12/02 (2006.01)

G 0 6 F 3/06 3 0 1 J
 G 0 6 F 12/00 5 1 4 M
 G 0 6 F 3/06 3 0 2 A
 G 0 6 F 13/14 3 2 0 H
 G 0 6 F 12/00 5 9 7 U

請求項の数 20 (全 46 頁) 最終頁に続く

(21) 出願番号 特願2016-521827 (P2016-521827)
 (86) (22) 出願日 平成26年6月20日 (2014.6.20)
 (65) 公表番号 特表2016-524250 (P2016-524250A)
 (43) 公表日 平成28年8月12日 (2016.8.12)
 (86) 国際出願番号 PCT/US2014/043299
 (87) 国際公開番号 W02014/205298
 (87) 国際公開日 平成26年12月24日 (2014.12.24)
 審査請求日 平成29年6月13日 (2017.6.13)
 (31) 優先権主張番号 13/924,567
 (32) 優先日 平成25年6月22日 (2013.6.22)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 314015767
 マイクロソフト テクノロジー ライセン
 シング, エルエルシー
 アメリカ合衆国 ワシントン州 9805
 2 レッドモンド ワン マイクロソフト
 ウェイ
 (74) 代理人 100140109
 弁理士 小野 新次郎
 (74) 代理人 100075270
 弁理士 小林 泰
 (74) 代理人 100101373
 弁理士 竹内 茂雄
 (74) 代理人 100118902
 弁理士 山本 修

最終頁に続く

(54) 【発明の名称】 複数のアクセスメソッドのためのラッチフリーのログ構造化ストレージ

(57) 【特許請求の範囲】

【請求項 1】

少なくとも1つのプロセッサを含むデバイスであって、前記少なくとも1つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含むデバイスを備え、前記データマネージャーは、

任意に選択されたページ指向型アクセスメソッドに、ページデータストレージに対するラッチフリーのアクセスを含む前記ページデータストレージに対するインターフェースアクセスをもたらすように構成されたデータ不透明型インターフェースであって、前記データ不透明型インターフェースは、アトミック操作を介して、キャッシュ層ストレージと、二次ストレージとを含むデータストレージの管理のために共通で利用される間接アドレスマッピングテーブル内のページ状態を表すストレージアドレスエントリーに対して、ラッチフリーのページ更新をもたらすデータ不透明型インターフェースと、

二次ストレージバッファの中に第1ページのページ状態をコピーすることを開始すること、および前記ページ状態の先頭にフラッシュデルタレコードを付加することを開始することに基づいて、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第1ページをフラッシュする操作を開始するよう構成されたページマネージャーであって、前記フラッシュデルタレコードは、二次ストレージ内の第1のページのストレージロケーションを示す二次ストレージアドレスと、呼び出し元に関連付けられた注釈とを含む、ページマネージャーとを含むシステム。

【請求項 2】

前記データ不透明型インターフェースは、前記任意に選択されたページ指向型アクセスメソッドに、ページデータストレージに対するログ構造化アクセスを含む前記ページデータストレージに対するインターフェースアクセスをもたらしように構成される請求項 1 に記載のシステム。

【請求項 3】

キャッシュ層マネージャーをさらに備え、前記キャッシュ層マネージャーは、前記データ不透明型インターフェースに関連付けられた間接アドレスマッピングテーブルに対するテーブル操作を開始するように構成されたマップテーブルマネージャーであって、前記テーブル操作は、前記間接アドレスマッピングテーブル内のエントリーに対してアトミック・コンペア・アンド・スワップ操作を開始して、前記ページデータストレージに関連付けられたページのそれまでの状態を、前記ページの新たな状態で置き換えるステップを含むマップテーブルマネージャーを含む請求項 1 に記載のシステム。

10

【請求項 4】

前記マップテーブルマネージャーは、前記データ不透明型インターフェースに関連付けられた前記間接アドレスマッピングテーブルに対する前記テーブル操作を開始するように構成され、前記間接アドレスマッピングテーブルは、キャッシュ層ストレージと、二次ストレージとを含むデータストレージの管理のために共通で使用される請求項 3 に記載のシステム。

【請求項 5】

20

前記間接アドレスマッピングテーブルは、ページの論理ロケーションを、前記ページの対応する物理ロケーションから分離し、

前記ページデータストレージのユーザーは、前記ページに関する物理ロケーションアドレス値の代わりにページ識別子値を、前記ページデータストレージを参照するデータ構造における別の場所に記憶させる請求項 3 に記載のシステム。

【請求項 6】

前記間接アドレスマッピングテーブル内のエントリーに対するラッチフリーのコンペア・アンド・スワップ操作を使用して、ページフラッシュからもたらされるログ構造化に関連するページロケーション変更を制御するように構成されたログ構造化ストレージ層マネージャーをさらに備える請求項 3 に記載のシステム。

30

【請求項 7】

アップデートデルタレコード操作と、置換アップデート操作に基づいてアップデートを制御するように構成されるレコードマネージャをさらに備える請求項 3 に記載のシステム。

【請求項 8】

第 1 のプロセッサ操作によってページ情報にアクセスすることに先立って、第 1 のエポックに関連づけられた第 1 のエポック登録リスト内で、第 1 のプロセッサ操作の登録を開始するよう構成された第 1 のエポックマネージャーをさらに備える請求項 1 に記載のシステム。

【請求項 9】

40

コンペア・アンド・スワップ (CAS) 操作を介して、マッピングテーブルの中にフラッシュデルタレコードに対するポインターをインストールするステップに基づいて、二次ストレージにページ状態をフラッシュするように構成されたページマネージャーであって、前記フラッシュデルタレコードは、前記 CAS 操作を介して前記マッピングテーブルの中で置き換えられる既存のページ状態の先頭に付加される、ページマネージャーをさらに備える請求項 1 に記載のシステム。

【請求項 10】

前記ページマネージャーは、

前記 CAS 操作が成功したかどうかを判定し、前記 CAS 操作が成功したと判定された場合、二次ストレージフラッシュバッファに前記既存のページ状態を書き込む書込み操

50

作を開始するように構成される請求項 9 に記載のシステム。

【請求項 1 1】

前記ページマネージャーは、

前記 C A S 操作が失敗したかどうかを判定し、前記既存のページにそれまで割り当てられていたストレージスペースを無効にする操作を開始するよう構成される請求項 1 0 に記載のシステム。

【請求項 1 2】

ラッチフリーのアップデート操作を介して、ログ構造化二次ストレージバッファに対してアップデートを制御するよう構成されたバッファーマネージャーをさらに備える、請求項 1 に記載のシステム。

10

【請求項 1 3】

前記バッファーマネージャーは、前記ラッチフリーの操作を介して、前記ログ構造化二次ストレージバッファを同時にアップデートするために、複数のスレッドを可能にするよう構成される、請求項 1 2 に記載のシステム。

【請求項 1 4】

前記バッファーマネージャーは、第 1 の二次ストレージアドレス引数までの、より下位のアドレスを有する、前記ログ構造化二次ストレージバッファにフラッシュされたページが、ログ構造化二次ストレージの中で安定していることを判定するための安定化操作を開始するよう構成される、請求項 1 2 に記載のシステム。

【請求項 1 5】

前記ページマネージャーは、さらに

20

コンペア・アンド・スワップ (C A S) 操作を介して、マッピングテーブルの中にフラッシュデルタレコードのアドレスをインストールすることに基づいて前記ページ状態をアップデートすることを開始することに基づいて、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第 1 のページをフラッシュする操作を開始するよう構成される、請求項 1 に記載のシステム。

【請求項 1 6】

少なくとも 1 つのプロセッサを含むデバイスであって、前記少なくとも 1 つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含むデバイスを備え、前記データマネージャーは、

30

任意に選択されたページ指向型アクセスメソッドに、ページデータストレージに対するラッチフリーのアクセスを含む前記ページデータストレージに対するインターフェースアクセスをもたらすように構成されたデータ不透明型インターフェースであって、前記データ不透明型インターフェースは、アトミック操作を介して、キャッシュ層ストレージと、二次ストレージとを含むデータストレージの管理のために共通で使用される間接アドレスマッピングテーブル内のページ状態を表すストレージアドレスエントリーに対して、ラッチフリーのページ更新をもたらす、データ不透明型インターフェースと、

第 1 のページに関連付けられるページ状態の先頭に部分的スワップデルタレコードを付加することを開始することに基づいて、二次ストレージ内のロケーションに対してキャッシュ層ストレージ内の第 1 のページの一部分をスワップする操作を開始するよう構成されたページマネージャーであって、前記部分的スワップデルタレコードは、前記第 1 のページの欠落した部分の二次ストレージ内のロケーションを示すフラッシュデルタレコードのストレージロケーションを示すメインメモリーアドレスを含む、ページマネージャーとを含む、システム。

40

【請求項 1 7】

少なくとも 1 つのプロセッサを含むデバイスであって、前記少なくとも 1 つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含むデバイスを備え、前記データマネージャーは、

50

二次ストレージ内のロケーションにキャッシュ層ストレージ内の第１のページをフラッシュする操作を、

二次ストレージバッファの中の前記第１のページのページ状態をコピーすることを開始するステップと、

前記ページ状態の先頭にフラッシュデルタレコードを付加することを開始するステップであって、前記フラッシュデルタレコードは、二次ストレージ内の前記第１のページのストレージロケーションと、呼び出し元に関連付けられた注釈とを示す二次ストレージアドレスを含む、ステップと、

コンペア・アンド・スワップ（ＣＡＳ）操作を介して、マッピングテーブルの中にフラッシュデルタレコードのアドレスをインストールすることに基づいて前記ページ状態に対するラッチフリーのアップデートを開始するステップと、に基づいて、開始するように構成されたページマネージャーを含む、システム。

10

【請求項１８】

任意に選択されたページ指向型アクセスメソッドに、ページデータストレージに対するラッチフリーのアクセスを含むページデータストレージに対するインターフェースアクセスをもたらすように構成されたデータ不透明型インターフェースと、

ラッチフリーのアップデート操作を介してログ構造化二次ストレージバッファに対するアップデートを制御するように構成されたバッファーマネージャーとをさらに含む、請求項１７に記載のシステム。

【請求項１９】

20

少なくとも１つのプロセッサを含むデバイスであって、前記少なくとも１つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含むデバイスを備え、

前記データマネージャーは、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第１のページの一部をラッチフリーでスワップする操作を、アトミック操作を介して、前記第１のページに関連付けられたページ状態の先頭に部分的スワップデルタレコードを付加することを開始するステップにより開始するように構成されたページマネージャーであって、前記部分的スワップデルタレコードは、前記第１のページの欠落した部分の二次ストレージ内のロケーションを示すフラッシュデルタレコードのストレージロケーションを示すメインメモリアドレスを含む、ページマネージャーを含む、システム。

30

【請求項２０】

任意に選択されたページ指向型アクセスメソッドに、ページデータストレージに対するラッチフリーのアクセスを含むページデータストレージに対するインターフェースアクセスをもたらすように構成されたデータ不透明型インターフェースをさらに備え、前記データ不透明型インターフェースは、前記任意に選択されたページ指向型アクセスメソッドに、前記ページデータストレージに対するログ構造化アクセスを含むページデータストレージへのインターフェースアクセスをもたらすように構成される、請求項１９に記載のシステム。

【発明の詳細な説明】

【背景技術】

40

【０００１】

[0001]電子デバイスのユーザーは、しばしば、様々なタイプの情報を獲得するのにデータベースシステムにアクセスする必要がある。データ項目を記憶するため、および取り出すための多くの異なる技法が考案されてきた。例えば、一部の最新のハードウェアプラットフォームは、より高いパフォーマンスをもたらそうとしてマルチコアプロセッサ、多層メモリー階層、ならびにフラッシュなどの二次ストレージデバイスなどの最新のハードウェア開発を活用してきた。このことは、潜在的なシステムパフォーマンスを向上させてきたが、システムが、新たに開発されたプラットフォーム態様、および従来のプラットフォーム態様を有効に利用することは困難であった。

【発明の概要】

50

【課題を解決するための手段】

【0002】

[0002] 1つの一般的な態様によれば、システムが、少なくとも1つのプロセッサを含むデバイスを含むことが可能であり、そのデバイスは、その少なくとも1つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含む。データマネージャーは、任意に選択されたページ指向型アクセスメソッドに、ページデータストレージに対するラッチフリー（latch-free）のアクセスを含むページデータストレージに対するインターフェースアクセスをもたらすように構成されたデータ不透明型インターフェース（data opaque interface）を含み得る。

【0003】

[0003]別の態様によれば、システムが、少なくとも1つのプロセッサを含むデバイスを含むことが可能であり、そのデバイスは、その少なくとも1つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含む。データマネージャーは、コンペア・アンド・スワップ（CAS：Compare and swap）操作を介して、マッピングテーブルの中でフラッシュデルタレコードにポインターをインストールすることに基づいて、二次ストレージにページ状態をフラッシュするように構成されたページマネージャーを含むことが可能であり、そのフラッシュデルタレコードは、CAS操作を介してマッピングテーブルの中で置き換えられる既存のページ状態の先頭に付加される。

【0004】

[0004]別の態様によれば、システムが、少なくとも1つのプロセッサを含むデバイスを含むことが可能であり、そのデバイスは、その少なくとも1つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含む。データマネージャーは、コンペア・アンド・スワップ（CAS）操作を介して、二次ストレージバッファーの中に第1のページのページ状態をコピーすることを開始すること、そのページ状態の先頭にフラッシュデルタレコードを付加することを開始することであって、そのフラッシュデルタレコードは、二次ストレージ内の第1のページのストレージロケーションを示す二次ストレージアドレスと、呼び出し元に関連付けられた注釈とを含むこと、およびマッピングテーブルの中にフラッシュデルタレコードのアドレスをインストールすることに基づいてページ状態をアップデートすることを開始することに基づいて、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第1のページをフラッシュする操作を開始するように構成されたページマネージャーを含み得る。

【0005】

[0005]別の態様によれば、システムが、少なくとも1つのプロセッサを含むデバイスを含むことが可能であり、そのデバイスは、その少なくとも1つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含む。データマネージャーは、ラッチフリーのアップデート操作を介してログ構造化二次ストレージバッファーに対するアップデートを制御するように構成されたバッファーマネージャーを含み得る。

【0006】

[0006]別の態様によれば、システムが、少なくとも1つのプロセッサを含むデバイスを含むことが可能であり、そのデバイスは、その少なくとも1つのプロセッサによって実行されるようにコンピューター可読記憶媒体上に実体化された命令を備えるデータマネージャーを含む。データマネージャーは、第1のページに関連付けられたページ状態の先頭に部分的スワップデルタレコードを付加することを開始することであって、その部分的スワップデルタレコードは、第1のページの欠落した部分の二次ストレージ内のロケーションを示すフラッシュデルタレコードのストレージロケーションを示すメインメモリアドレスを含むことに基づいて、二次ストレージ内のロケーションに対してキャッシュ層ストレージ内の第1のページの一部をスワップする操作を開始するように構成されたページマネージャーを含み得る。

【 0 0 0 7 】

[0007]この概要は、詳細な説明において後段でさらに説明される選定された概念を、簡略化された形態で概説するように与えられる。この概要は、主張される主題の重要な特徴または不可欠な特徴を特定することは意図しておらず、主張される主題の範囲を限定するのに使用されることも意図していない。1つまたは複数の実施形態の詳細は、添付の図面、および後段の説明において示される。その他の特徴が、説明および図面、ならびに特許請求の範囲から明白となる。

【図面の簡単な説明】

【 0 0 0 8 】

【図 1】[0008]キャッシュ/ストレージ層に関するアクセスメソッドのための例示的なアーキテクチャー階層化を示す図である。

10

【図 2】[0009]複数のアクセスメソッドのためのラッチフリーのログ構造化ストレージのための例示的なアーキテクチャーを示すブロック図である。

【図 3】[0010]例示的なマッピングテーブルを示す図である。

【図 4 a】[0011]例示的なマッピングテーブル上の例示的なデルタアップデートを示す図である。

【図 4 b】例示的なマッピングテーブル上の例示的なデルタアップデートを示す図である。

【図 5】[0012]例示的な部分的ページスワップアウトおよび例示的な部分的スワップデルタを示す図である。

20

【図 6】[0013]例示的なエポック、およびそれらのエポックのそれぞれのガベージコレクションリストを示す図である。

【図 7 a】[0014]フラッシュ上の例示的なログ構造化ストレージ編成を示す図である。

【図 7 b】フラッシュ上の例示的なログ構造化ストレージ編成を示す図である。

【図 7 c】フラッシュ上の例示的なログ構造化ストレージ編成を示す図である。

【図 8】[0015]例示的なフラッシュバッファ状態を示す図である。

【図 9】[0016]例示的なトランザクションテンプレートを示す図である。

【図 1 0】[0017]例示的なチェックポイントデータを示す図である。

【図 1 1】[0018]複数のアクセスメソッドのためのラッチフリーのログ構造化ストレージのための例示的なシステムを示すブロック図である。

30

【図 1 2 a】[0019]図 1 1 のシステムの例示的な動作を示すフローチャートである。

【図 1 2 b】図 1 1 のシステムの例示的な動作を示すフローチャートである。

【図 1 2 c】図 1 1 のシステムの例示的な動作を示すフローチャートである。

【図 1 2 d】図 1 1 のシステムの例示的な動作を示すフローチャートである。

【発明を実施するための形態】

【 0 0 0 9 】

[0020] I . 概説

ハードウェアプラットフォームの最新の開発は、より高いパフォーマンスをもたらそうとしてマルチコアプロセッサ、多層メモリー階層、ならびにフラッシュなどの二次ストレージデバイスを活用してきた。例えば、中央処理装置 (CPU) 変更は、マルチコアプロセッサ、ならびに複数のレベルのキャッシングが関与するメインメモリーアクセスを含む。例えば、フラッシュストレージ、ならびに容量を低下させるアップデート・インプレース (update-in-place) を行うハードディスクベンダ認識は、ログ構造化をより多く使用することにつながってきた。例えば、クラウドデータセンターは、システムスケールを増大させ、コモディティハードウェアの使用は、高い利用可能性技法により大きな重きを置く。

40

【 0 0 1 0 】

[0021]しかし、潜在的なシステムパフォーマンスが向上する一方で、システムが、これらの最新のプラットフォーム態様を有効に利用することは困難であり得る。例えば、大量のデータにアクセスする複数のユーザーをサポートするデータ指向型システムは、過去に

50

長年、存在してきたとおりのハードウェアのために設計されたソフトウェアアーキテクチャを活用する可能性がある（例えば、それらのシステムは、単一レベルメモリー（プロセッサキャッシングをほとんど行わず、メインメモリーに対して控えめな潜時しか有さない）上で動作し、磁気ディスクにアクセスして、ユニプロセッサをターゲットにすることが可能である）。

【 0 0 1 1 】

[0022]そのアプローチを変えようとする取組みは、環境を改善してきたが、相変わらず、相当な潜在的パフォーマンス利得を逃している。例えば、データ競合にアクセスする際、ブロックングを生じさせるラッチを回避する取組みが行われてきたが、これらの取組みには、スレッドがそのような競合を回避するようにパーティションすることが関与してきた可能性があり、このことは、相当なオーバーヘッドをもたらし得る。例えば、データのアップデート・インプレースを行うことは、メモリーパフォーマンスに悪影響を及ぼす可能性があり、このことは、キャッシュ・ライン・アラインメント（cache line alignment）を配慮すること、およびベクトルにわたる二分探索の代わりにローカルツリーを使用することにつながり得る。しかし、不利に多くのアップデート・インプレースが相変わらず行われ、このことが、例えば、キャッシュ無効化を介して、プロセッサ・キャッシング・パフォーマンスに悪影響を及ぼす。さらに、実施形態が、フラッシュを、フラッシュのより高い毎秒アクセス数、およびより低いアクセス潜時のために活用し始めている。しかし、ランダムなアップデートは、フラッシュ変換層を使用してさえ、比較的高価であり得る。

【 0 0 1 2 】

[0023]J. Levandoski他、「Deuteronomy: Transaction Support for Cloud Data (Deuteronomy: クラウドデータのためのトランザクションサポート)」、Conference on Innovative Data Systems Research (CIDR) (2011年1月)、123~133ページ、およびD. Lomet他、「Unbundling Transaction Services in the Cloud (クラウドにおけるトランザクションサービスのアンバンドリング)」、Conference on Innovative Data Systems Research (CIDR)、2009年が、クラウド環境において一貫性（すなわち、トランザクション）をもたらすための例示的な技法を説明する。本明細書で説明される例示的な技法は、例示的なDEUTERONOMYデータコンポーネント（DC）、および現在のハードウェア上でそのDCのパフォーマンスを最大にすることに焦点を合わせる。例えば、DCが、CRUD（作成、読取り、アップデート、削除）アトミック操作を介してアクセスされるデータの記憶および取出しを管理することが可能である。例えば、DCは、DCの上のソフトウェア層（例えば、DEUTERONOMYトランザクションコンポーネント（TC）および/またはクエリエンジン）を介して融合されて分散システムにされ得るローカル機構を代わりに使用して、非分散型であってもよい。

【 0 0 1 3 】

[0024]本明細書でさらに説明されるとおり、アクセスメソッド（例えば、Bツリー、ハッシング、マルチ属性、時間的、その他）に影響を及ぼし得る、現在のハードウェアによってもたらされる問題が存在するように思われる。さらに、本明細書で説明されるとおり、これらの問題は、ほとんどの（例えば、任意に選択された）アクセスメソッドに適用可能である例示的な一般的な機構を用いて解決され得る。

【 0 0 1 4 】

[0025]例えば、本明細書で説明される例示的な技法によれば、ラッチフリーの技法は、マルチコアプロセッサを用いて有利なプロセッサ利用およびスケーリングを実現するのに利用され得る。例えば、本明細書で説明されるとおり、キャッシュ無効化を少なくするデルタアップデートが、マルチレベルキャッシュベースのメモリーシステムを用いて有利なパフォーマンスを実現するのに利用され得る。例えば、ランダム書込みの限られたパ

パフォーマンス、およびフラッシュ書込み限度を有する書込み制限されたストレージが、ログ構造化を介して克服され得る。

【0015】

[0026]例えば、BW - TREE (例えば、J. Levandoski 他、「The Bw - Tree: A B - tree for New Hardware Platforms (Bw - Tree: 新しいハードウェアプラットフォームのためのBツリー)」、29th IEEE International Conference on Data Engineering (ICDE 2013)、2013年4月8~11日参照)、Bツリーといくぶん似通ったインデックス (例えば、R. Bayer 他、「Organization and Maintenance of Large Ordered Indices (大次数インデックスの編成および管理)」、Acta Informatica, Vol. 1, Issue 3, 1972, 173~189ページ、およびD. Comer、「The Ubiquitous B - tree (ユビキタスBツリー)」、ACM Computing Surveys (CSUR), Vol. 11, Issue 2, 1979年6月、121~137ページ参照) が、これらの例示的な技法を活用することが可能なDCまたはキー・バリューストア (key-value store) の例である。例示的なBW - TREEには、ラッチフリーダムであることおよびログ構造化をより一般的に実現するための技法のためのパラダイムが関与し得る。本明細書で説明される例示的な技法によれば、ラッチフリーの技法およびログ構造技法は、従来のキャッシュ/ストレージサブシステムが、アップデート・インプレースとしてディスクにライトバック (write back) される固定サイズページに対するラッチアクセス (latched access) を扱い得るのといくぶん似通った状態で複数のアクセスメソッドをサポートすることができるキャッシュ/ストレージサブシステムにおいて実施され得る。

【0016】

[0027]本明細書で説明される例示的な技法によれば、LLAMA (Latch - free, Log - structured Access Method Aware) と本明細書で呼ばれ得る例示的なシステムが、最近、開発されたハードウェア環境 (例えば、フラッシュ、マルチコア) のためのキャッシング - ストレージサブシステムを (少なくとも) 含み、ただし、そのような例示的な技法は、最近、開発されたハードウェアだけに限定されないことがデータ処理の分野の業者には理解されよう。

【0017】

[0028]例えば、LLAMAは、プロセッサキャッシュおよび二次ストレージを最適化して、キャッシュ管理とストレージ管理の両方をもたらし任意に選択されたページ指向型アクセスメソッドのためのアプリケーションプログラミングインターフェース (API) をサポートすることが可能である。例えば、キャッシング (CL) 層とストレージ (SL) 層が、ページの論理ロケーションと物理ロケーションを分離する共通のマッピングテーブルを使用することが可能である。例えば、キャッシュ層 (CL) が、キャッシュ層 (CL) のマッピングテーブル上のラッチフリーのコンペア・アンド・スワップ・アトミック状態変化を介してデータアップデートおよび管理アップデート (例えば、インデックス再編成のための) をサポートすることが可能である。

【0018】

[0029]例えば、ストレージ層 (SL) が、ページフラッシュごとにログ構造化することによってもたらされるページロケーション変化を扱うのに同一のマッピングテーブルを使用することが可能である。例えば、ラッチフリーのBW - TREE実施形態 (例えば、順序付けられたBツリー型インデックスの例としてBW - TREEを使用する実施形態) が使用され得る。この脈絡において、「フラッシュする」操作とは、メインメモリー (例えば、キャッシュストレージ) からのページを、そのページを出力バッファにコピーすることによって、二次ストレージに転送することを指すことが可能である。

【0019】

[0030]本明細書で説明される例示的な技法は、ページのロケーションとサイズの両方を

10

20

30

40

50

視覚化し得るマッピングテーブルをもたらすことが可能である。例えば、そのような視覚化は、本明細書でさらに説明されるとおり、メインメモリ設計と安定したストレージ設計（例えば、ログ構造化ストレージ設計）の両方のために利用され得る。

【 0 0 2 0 】

[0031]この脈絡において、「ページ」とは、物理ストレージアドレスを介してアクセスされ得るストレージ内のオブジェクトを指すことが可能である。本明細書で使用される「ページ」は、柔軟性のあるサイズに関連付けられることが可能であり、ストレージの複数の不連続に記憶されたセグメントにわたって分散することが可能なストレージのページ単位を表すことが可能である。ストレージには、揮発性ストレージおよび/または安定したストレージが含まれ得る。

10

【 0 0 2 1 】

[0032]本明細書で説明される例示的な技法は、アクセスメソッド層をキャッシュ/ストレージ管理から分離することが可能である。例として、本明細書で説明される技法は、先行書込みログプロトコル (write-ahead log protocol) を執行するのに使用され得る。例えば、ページをフラッシュするのに先立って、従来のデータベースカーネルは、ページログシーケンスナンバー (LSN) を調べて、トランザクションログにおいてまだ安定していないアップデートが存在するかどうかを判定することが可能である。例えば、LLAMA キャッシュ管理は、例示的なデルタアップデートを活用して、部分的ページを「スワップアウトする」ことが可能である。例えば、LLAMA キャッシュ管理は、二次ストレージ上に既に存在するページの部分（最近のデルタアップデートを含まない）をキャッシュからドロップすることができる。例えば、アクセスメソッド層が、トランザクション・ログ・チェックポインティングのために定期的にフラッシュされる。このため、キャッシュマネージャーは、任意のバッファサイズ制約を満足させるのに十分な候補（場合により、部分的な）ページを見出す。

20

【 0 0 2 2 】

[0033]本明細書で説明される例示的な技法は、相当な数のアクセスメソッド（すなわち、単一のインスタンスだけでない）がこれらの技法を活用することによって可能にするフレームワークを提供することが可能である。さらに、有利な効率をもたらす二次ストレージにデータを書き込むためのログ構造化ストアが実施されることが可能である。したがって、アクセスメソッドは、アクセスメソッドのインデックスのメインメモリ態様に焦点を合わせることが可能であり、本明細書で説明される例示的な技法は、BW-TREE のパフォーマンスメトリックと同様のパフォーマンスメトリックを実現するためのフレームワークをもたらすことが可能である。

30

【 0 0 2 3 】

[0034]例えば、LLAMA などの技法が、その技法のAPIを介して、ラッチフリーのページアップデートをもたらすことが可能であり、このことは、マッピングテーブル上のコンペア・アンド・スワップ (CAS) アトミック操作を介してメインメモリの中で達せられる。

【 0 0 2 4 】

[0035]例えば、キャッシュを管理する際に、LLAMA などの技法は、ページの以前にフラッシュされた部分だけをメモリからドロップすることによって、メインメモリを取り戻すことが可能であり、このため、「汚れた」ページをスワップアウトする場合でさえ、入出力 (I/O) 操作がまったく関与しない。このため、LLAMA などの技法は、その技法のアクセスメソッドユーザーからの入力なしに、その技法のバッファ・キャッシュ・メモリ・サイズを制御できることが可能である。

40

【 0 0 2 5 】

[0036]例えば、二次ストレージの効果的な管理のために、LLAMA などの技法が、ログ構造化を利用することが可能である。例えば、LLAMA などの技法が、部分的ページフラッシュ、および実質的に空きスペースのない、すなわち、実質的に100%ストレージ

50

ジ利用率のページを使用することによって、従来のログ構造化と比べてパフォーマンスを向上させることが可能である。これらは、ページがフラッシュされる際の入出力操作（Ｉ／Ｏ）の回数、および１ページ当たり消費されるストレージの量を減らすことが可能であり、したがって、ログ構造化が使用される場合に経験され得るライトアンプリフィケーション（write amplification）を減らすことが可能である。さらに、すべてのストレージ関連の操作が、完全にラッチフリーであることが可能である。

【 0 0 2 6 】

[0037]例えば、ＬＬＡＭＡなどの技法が、（少なくとも）制限された形態のシステムトランザクションをもたらすことが可能である。この意味で、システムトランザクションは、ユーザーレベルトランザクションではなく、むしろ、ログ構造化ストアを活用して、もっぱら、アクセスメソッドの「プライベート使用」のため（例えば、インデックス構造変更（ＳＭＯ）のため）にアトミック性をもたらす。例えば、このことは、同時に行われるアップデートが続けられる間にインデックスが成長するにつれ、インデックスが適応することを可能にし得る。

10

【 0 0 2 7 】

[0038]例えば、ＢＷ－ＴＲＥＥ構造は、あるタイプのラッチフリーのＢツリー構造を含み得る。例えば、ＢＷ－ＴＲＥＥノードに対するアップデートが、それまでのページ状態の先頭にアップデートデルタを付加することに基づいて実行され得る。このため、ＢＷ－ＴＲＥＥは、複数のスレッドによるページに対する同時に行われるアクセスを許すことが可能であるので、ラッチフリーであり得る。そのようなデルタアップデートは、ページのそれまでの状態を保存するため、向上したプロセッサキャッシュパフォーマンスをもたらすことも可能である。

20

【 0 0 2 8 】

[0039]ＢＷ－ＴＲＥＥを使用する例示的な技法は、やはりラッチフリーであり、Ｂリンクツリー型サイドポインターを使用することが可能なページ分割技法をさらにもたらし得る。分割（および他の構造変更操作）は、メインメモリー内部でも、安定させられる際もともにアトミックであり得る。例えば、アトミックレコードストアが、ＢＷ－ＴＲＥＥアーキテクチャーに基づいて実施され得る。

【 0 0 2 9 】

[0040]本明細書の説明の趣旨を逸脱することなく、本明細書で説明されるラッチフリーのログ構造化ストレージを実現する多くの方法が存在し得ることが、データ処理の分野の業者には理解されよう。

30

【 0 0 3 0 】

[0041] Ｉ Ｉ . 例示的な動作環境

本明細書で説明される特徴は、本明細書の説明の趣旨を逸脱することなく、データ処理の分野の業者によって理解され得る多くの異なる方法において実施され得る例示的な実施形態として与えられる。そのような特徴は、例示的な実施形態特徴としてのみ解釈されるべきであり、それらの詳細な説明だけに限定されるものと解釈されることは意図していない。

【 0 0 3 1 】

40

[0042]図１は、キャッシュ／ストレージ層に関するアクセスメソッドのための例示的なアーキテクチャー階層化を示す。アクセスメソッド層１０２は、図１に示されるとおり、最上位層である。アクセスメソッド層１０２は、中位層であるキャッシュ層１０４と対話する。アプリケーションプログラミングインターフェース（ＡＰＩ）１０６が、アクセスメソッド層１０２とキャッシュ層１０４の間の活動のために使用され得る。例示的なストレージ層１０８が、マッピングテーブル１１０と対話することが可能であり、マッピングテーブル１１０は、キャッシュ層１０４とストレージ層１０８の間で共有され得る。例えば、ＬＬＡＭＡ１１２が、キャッシュ層１０４およびストレージ層１０８を含む。例えば、ストレージ層が、ログ構造化フラッシュストアをサポートすることが可能である。本明細書で説明される例示的な技法によれば、ログ構造化ストアが、フラッシュとディスクス

50

トレージの両方を管理することが可能である。この設計は、アーキテクチャー上、既存のデータベースカーネルと共存することが可能である一方で、スタンドアロン型またはDEUTERONOMY型のアトミックレコードストア(ARS)としても適している。

【0032】

[0043]例えば、LLAMAなどの技法が、キャッシュ/ストレージ層に関するアクセスメソッド実施形態をサポートして、ページ抽象化をサポートすることが可能である。さらに、トランザクションコンポーネント(例えば、DEUTERONOMY型トランザクションコンポーネント)が上に追加され得る。図2は、複数のアクセスメソッドのためのラッチフリーのログ構造化ストレージのための例示的なアーキテクチャーを示すブロック図である。図2に示されるとおり、トランザクションコンポーネント202が、トランザク
10
ションキー・バリューストアをサポートすることが可能であり、アトミックキー・バリューストアを含み得るデータコンポーネント204と一緒に動作することが可能である。図2に示されるとおり、データコンポーネント204は、ラッチフリーの順序付けられたインデックス206および/またはラッチフリーの線形ハッシングインデックス208を含み得る。図2に示されるとおり、データコンポーネント204は、例示的なラッチフリーのログ構造化アクセスメソッドウェア(LLAMA)ストレージエンジン210(例えば、図1のLLAMA112)をさらに含み得る。

【0033】

[0044]例示的なAPI106は、「データ不透明型」であることが可能であり、すなわち、例示的なLLAMA実施形態は、アクセスメソッド(例えば、アクセスメソッド層1
20
02の)が何をページまたはデルタレコードに入れているかを「見る」ことがなく(例えば、検査する、または分析する、または依存することがなく)、アクセスメソッドによってページまたはデルタレコードの中で何が提供されるかとは無関係に動作する。このため、本明細書で説明される例示的なLLAMA実施形態は、前述したとおり、アクセスメソッドによって何が提供されるかに依存しない特定の操作に応答して動作することが可能である。

【0034】

[0045]図3に示されるとおり、ページ302が、メインメモリーキャッシュ312内で、または二次ストレージ314上で、ページ識別子(PID)306を状態308にマップする(例えば、マッピングテーブル304に記憶された「物理アドレス」310を介して)マッピングテーブル304を介してアクセスされ得る。例えば、メインメモリーキャ
30
ッシュ312は、ランダムアクセスメモリー(RAM)を含み得る。例えば、二次ストレージ314は、フラッシュメモリーを含み得る。例えば、ページ302は、オンデマンドで二次ストレージ314からメインメモリーキャッシュ312に読み取られることが可能であり、ページ302は、二次ストレージ314にフラッシュされることが可能であり、ページ302は、キャッシュ312内にある間、ページ状態を変更するようにアップデートされることが可能である。例えば、実質的にすべてのページ状態変更(データ状態と管理状態の両方)は、本明細書で説明される例示的な技法により、アトミック操作として与えられ得る。図3に示されるとおり、例示的な物理アドレス310が、その物理アドレス
40
がフラッシュまたはメモリー(例えば、キャッシュ)に関連付けられているかどうかを示すフラッシュ/メモリーフラグ316(例えば、この例に示されるとおり、1ビットに関する)を、(少なくとも)そのアドレス自体に関するアドレスフィールド318(例えば、この例に示されるとおり、63ビットに関する)と一緒に含み得る。本明細書の説明の趣旨を逸脱することなく、「物理アドレス」を表す多くの方法(例えば、64ビット表現以外の)が存在することが、データ処理の分野の業者には理解されよう。

【0035】

[0046]本明細書で説明される例示的な技法によれば、LLAMAが、LLAMAのAPIを通じて、(例えば、スレッドをブロックすることによって同時に行われるアクセスからページを保護する従来のラッチの代わりに)マッピングテーブル304上でコンペア・
50
アンド・スワップ(CAS)アトミック操作を介してラッチフリーのページアップデート

をもたらすことが可能である。例えば、C A S戦略は、プロセッサ利用率を有利に高めること、およびマルチコアスケールリングを向上させることが可能である。

【 0 0 3 6 】

[0047]本明細書で説明される例示的な技法によれば、キャッシュを管理する際、L L A M Aは、ページの以前にフラッシュされた部分だけをメモリーからドロップすることによって、メインメモリーを取り戻すことが可能であり、このため、「汚れた」ページをスワップアウトする場合でさえ、入出力（I / O）操作をまったく使用しない。このため、L L A M Aなどの例示的なアーキテクチャーは、そのアーキテクチャーのアクセスメソッドユーザーによってページの中に記憶されたデータを検査する必要なしに、そのアーキテクチャーのバッファー・キャッシュ・メモリー・サイズを制御することが可能である（例えば、L L A M Aなどの例示的なアーキテクチャーは、トランザクションおよび先行書込みロギングを意識しないので）。

10

【 0 0 3 7 】

[0048]L L A M Aなどの例示的なアーキテクチャーは、ログ構造化を使用して二次ストレージを管理することが可能である（例えば、ランダムな書込みを回避すること、大きいマルチページバッファーを介する書込みの回数を減らすこと、およびフラッシュメモリーに関するウェアレベリングの利点を提供して）。さらに、L L A M Aなどの例示的なアーキテクチャーが、部分的ページフラッシュ、および実質的に空きスペースのない、すなわち、実質的に100%ストレージ利用率のページを用いてパフォーマンスを有利に向上させる（例えば、従来のログ構造化と比較して）ことが可能である。例えば、これらは、ページがフラッシュされる際のI / Oの回数、および1ページ当たり消費されるストレージの量を減らすことが可能であり、したがって、さもなければ、ログ構造化が使用される場合に経験され得るライトアンプリフィケーションを減らすことが可能である。さらに、実質的にすべてのストレージ関連の操作が、完全にラッチフリーであることが可能である。

20

【 0 0 3 8 】

[0049]さらに、L L A M Aなどの例示的なアーキテクチャーが、（少なくとも）制限された形態のシステムトランザクションをサポートすることが可能である（システムトランザクションに関しては、例えば、D . L o m e t他、「Unbundling Transaction Services in the Cloud（クラウドにおけるトランザクションサービスのアンバンドリング）」、Conference on Innovative Data Systems Research（CIDR）、2009参照）。例えば、システムトランザクションは、ユーザートランザクションではないことが可能であり、むしろ、もっぱら、アクセスメソッドの「プライベート使用」のため（例えば、インデックス構造変更（SMO）（C . M o h a n他、「ARIES / IM : An Efficient and High Concurrency Index Management Method Using Write - Ahead Logging（ARIES / IM : 先行書込みロギングを用いた効率的な高同時実行性インデックス管理方法）」、In Proceedings of the 1992 ACM SIGMOD International Conference on Management of Data（SIGMOD '92）、1992、371~380ページ）のため）にアトミック性をもたらすことが可能である。例えば、トランザクションログとは別個に記録されたシステムトランザクションが効果的であり得るという特性が、データベースカーネルを分解することに対するD E U T E R O N O M Yアプローチの有利な洞察の例である。

30

40

【 0 0 3 9 】

[0050]後段の説明は、L L A M Aなどの例示的なアーキテクチャーを使用する際にアクセスメソッド実施者が経験し得る例示的な操作インターフェースのさらなる説明を、そのアーキテクチャーがどのように使用され得るかのさらなる説明とともに含む。後段の説明は、本明細書で説明される例示的な技法による、例示的なキャッシュ層、およびログ構造

50

化ストレージ層の例示的な設計のさらなる説明を含む。さらに、本明細書で説明される例示的な技法による、アトミック性をもたらすものと理解され得る例示的なシステムトランザクション機構および例示的な測定に関する説明が与えられる。さらに、本明細書で説明される例示的な技法による、システムクラッシュからの例示的なログ構造化ストレージ回復に関する説明が与えられる。

【 0 0 4 0 】

[0051] L L A M A などの例示的なシステムを設計する際、設計目標は、目標が可能な限り「汎用」であることを含むことが可能であり、このことは、ときとして、目標が可能な限り「低レベル」であることにつながる可能性がある。しかし、L L A M A などの例示的なシステムが「汎用」であるには、そのシステムのファシリティを使用する際にアクセスメソッドが何を行うかについて可能な限りわずかしき知らずに、効果的に動作することが望ましい可能性がある。このため、L L A M A などの例示的なシステムの動作は、キャッシュ管理、およびページのアップデートをターゲットとして、「原始的」であり得る。例えば、L L A M A などの例示的なシステムは、S M O（例えば、ページ分割およびページマージ）のために有利に含められ得る原始的トランザクション機構をサポートするいくつかのさらなるファシリティを含み得る。

10

【 0 0 4 1 】

[0052] 本明細書で説明される例示的な技法によれば、L L A M A などの例示的なシステムが、ログシーケンスナンバー（L S N）、先行書込みロギング、またはトランザクションログに関するチェックポイントに関するインターフェースに何も含めないことが可能である。本明細書で説明される例示的な技法によれば、L L A M A などの例示的なシステムは、ユーザー操作に関する冪等性試験（idempotence test）をまったく含まないことが可能である。さらに、本明細書で説明される例示的な技法によれば、L L A M A などの例示的なシステムが、トランザクション回復（例えば、本明細書で説明される例示的な技法により、L L A M A などの例示的なシステムを使用してアクセスメソッドによって扱われ得る）をまったく含まないことが可能である。

20

【 0 0 4 2 】

[0053] 本明細書で説明される例示的な技法によれば、例示的なアクセスメソッドが、ユーザー操作に応答して状態を変更することが可能である。例えば、ユーザーが、レコードを作成すること（C）、読み取ること（R）、アップデートすること（U）、または削除すること（D）（例えば、C R U D 操作）を所望する可能性がある。本明細書で説明される例示的な技法によれば、L L A M A などの例示的なシステムは、これらの操作を直接にはサポートしない可能性がある。むしろ、例示的なアクセスメソッドは、これらの操作を、L L A M A ページの状態に対するアップデートとして実施することが可能である。

30

【 0 0 4 3 】

[0054] 例えば、例示的なアクセスメソッド操作の一部である構造変化が存在する可能性もある。例えば、B W - T R E E ページ分割には、元のページ O に分割デルタを書き込んで、探索者が、新たなページが現時点で、O の中のキーの部分的範囲に関するデータを含むことを知るようにすることが関与し得る。例えば、これらも、L L A M A ページ O に対するアップデートとして扱われ得る。

40

【 0 0 4 4 】

[0055] 本明細書で説明される例示的な技法によれば、L L A M A などの例示的なシステムが、2つの形態のアップデート、例えば、デルタアップデートおよび置換アップデートをサポートすることが可能である。例えば、アクセスメソッドが、ユーザーの所望に応じて、これらの形態のアップデートを活用することを選択することが可能である。例えば、B W - T R E E が、一連のデルタアップデートを行い、何らかの時点で、それらのデルタアップデートを基本ページに適用することによって、そのページを「統合」し、最適化することを決定することが可能である。例えば、B W - T R E E が、その後、置換アップデートを使用して、その新たな基本ページを生成することが可能である。

【 0 0 4 5 】

50

[0056]本明細書で説明される例示的な技法によれば、LLAMAなどの例示的なシステムが、本明細書で説明されるアップデート操作および置換操作の全体にわたって、ページの物理ロケーションに関する情報を二次ストレージの中に保持することが可能であり、したがって、システム100が、本明細書でさらに説明されるとおり、ページがメインメモリーキャッシュからスワップアウトされるような場合にページを再度読み取るために、およびガベージコレクションのために、二次ストレージ・ページ・ロケーション情報を有する。このため、システム100は、以前のページロケーション、および安定ページ状態情報を記憶していることが可能である。

【0046】

[0057]例えば、デルタアップデートが、Update-D(PID, in_ptr, out_ptr, data)として示されることが可能である。例えば、デルタアップデートは、ページのそれまでの状態に対する変更を記述するデルタを先頭に付加することが可能である。例えば、BW-TREEの場合、Update-Dに対する「data」パラメーターは、少なくとも<lsn, key, data>を含むことが可能であり、ここで、lsnは、冪等性を可能にする。例えば、「in_ptr」は、ページのそれまでの状態をポイントし、「out_ptr」は、ページの新たな状態をポイントする。

【0047】

[0058]例えば、置換アップデートが、Update-R(PID, in_ptr, out_ptr, data)として示され得る。例えば、置換アップデートが、ページに関する完全に新たな状態をもたらすことが可能である。Update-Dを使用している場合に保存されたそれまでの状態が、「data」パラメーターによって取って代わられることが可能である。このため、「data」パラメーターは、デルタが「折り込まれた」ページの状態全体を包含する。

【0048】

[0059]例えば、「読取り」が、Read(PID, out_ptr)として示され得る。例えば、読取りは、「out_ptr」を介して、ページのメインメモリー内のアドレスを戻すことが可能である。そのページがメインメモリーの中にある場合、マッピングテーブルエントリが、二次ストレージアドレスを包含することが可能である。例えば、その場合、そのページが、メインメモリーの中に読み込まれることが可能であり、マッピングテーブルが、その新たなメインメモリーアドレスでアップデートされることが可能である。

【0049】

[0060]データ操作をサポートすることに加えて、本明細書で説明される例示的なシステム（例えば、LLAMA）は、ページの存在、ロケーション、および永続性を管理する操作を提供することが可能である。記憶されたデータの量を調整するのに、アクセスメソッドは、そのアクセスメソッドによって管理されるコレクションにページを足す、またはそれらのコレクションからページを引くことが可能である。状態永続性をもたらすのに、アクセスメソッドは、時々、ページを二次ストレージにフラッシュすることが可能である。この永続性を管理するのに、ページは適切に（例えば、ログシーケンス番号（lsn）を用いて）注釈を付けられることが可能である。例えば、ページマネージャーが、ページに対するフラッシュ操作、割当て操作、および解放操作を制御するように構成され得る。

【0050】

[0061]例えば、フラッシュ操作が、Flush(PID, in_ptr, out_ptr, annotation)として示され得る。例えば、Flushは、ページ状態をログ構造化ストア(LSS) I/Oバッファーに入れるようにコピーすることが可能である。Flushは、デルタ（注釈を伴う）をそれまでの状態の先頭に付加するので、メインメモリーに及ぼす影響においてUpdate-Dといくぶん似通っていることが可能である。このデルタには、「フラッシュ」というタグが付けられることが可能である。本明細書で説明される例示的な技法によれば、LLAMAなどの例示的なシステムが、ページが位置するLSS二次ストレージアドレス（フラッシュオフセットと呼ばれる）、およびフ

10

20

30

40

50

ラッシュデルタの中の呼び出し元「`annotation`」を記憶することが可能である。例えば、`Flush`は、戻す際、`I/O`バッファが安定していることをユーザーに保証しない可能性がある。

【0051】

[0062]例えば、バッファーマネージャーが、ラッチフリーのアップデート操作を介してログ構造化二次ストレージバッファに対するアップデートを制御するように構成され得る。このため、例えば、複数のスレッドが、ラッチフリーの操作を介してログ構造化二次ストレージバッファを同時にアップデートすることが可能である。

【0052】

[0063]例えば、「安定させる」操作が、`Mk - Stable (LSS address)`として示され得る。例えば、`Mk__Stable`操作は、`LSS`アドレス引数までの、`LSS`バッファにフラッシュされたページが二次ストレージ上で安定していることを確実にすることが可能である。`Mk__Stable`が戻す際、与えられる`LSS address`、およびすべてのより下位の`LSS`アドレスは、二次ストレージ上で安定していることが確実にされる。

【0053】

[0064]例えば、「高安定」操作が、`Hi - Stable (out - LSS address)`として示され得る。例えば、`Hi__Stable`操作が、二次ストレージ上で現在、安定している最上位の`LSS address`を戻すことが可能である。

【0054】

[0065]例えば、ページマネージャーが、第1のページのページ状態を第2のストレージバッファに入れるようにコピーすることを開始すること、そのページ状態の先頭にフラッシュデルタレコードを付加することを開始することであって、そのフラッシュデルタレコードは、二次ストレージ内の第1のページのストレージロケーションを示す二次ストレージアドレスと、呼び出し元に関連付けられた注釈を含むことに基づいて、キャッシュ層ストレージ内の第1のページを二次ストレージ内のロケーションにフラッシュする操作を開始するように構成され得る。

【0055】

[0066]例えば、バッファーマネージャーが、第1の二次ストレージアドレス引数までの、より下位のアドレスを有する、二次ストレージバッファにフラッシュされたページが、二次ストレージの中で安定していることを判定するための安定化操作を開始するように構成され得る。

【0056】

[0067]例えば、「割り当てる」操作が、`Allocate (out - PID)`として示され得る。例えば、`Allocate`操作は、マッピングテーブルの中で割り当てられた新たなページの`PID`を戻すことが可能である。すべてのそのようなページは、永続的に記憶されていることが可能であり、したがって、`Allocate`は、含められた操作を自動的にフラッシュすることが可能なシステムトランザクション（後段でさらに説明される）の一部として含められることが可能である。

【0057】

[0068]例えば、「解放する」操作が、`Free (PID)`として示され得る。例えば、`Free`操作は、`PID`によって識別されたマッピングテーブルエントリを再使用のために利用可能にすることが可能である。メインメモリーの中で、`PID`は、現在のエボック（後段でさらに説明される）に関する`PID`に関する待ち状態の解放リスト上に置かれ得る。この場合も、アクティブなページは記憶されていることが可能であるため、`Free`は、システムトランザクションの一部として含められることが可能である。

【0058】

[0069]本明細書で説明される例示的な技法によれば、例示的な`LLAMA`システムトランザクションが、構造変更（例えば、`SMO`）に関する相対的な持続性およびアトミック性（すべてまたは無）をもたらすのに使用され得る。例えば、`LSS`、および`LSS`のペ

10

20

30

40

50

ージ指向型レコードが、「ログレコード」として使用され得る。例えば、トランザクション内のすべての操作が、キャッシュ内のページ状態を変更することに加えて、メモリー内 L S S I / O バッファに自動的にフラッシュされ得る。例えば、各 L S S エントリーが、厳密に「ページ」ストアである例示的な L S S に関して、ページの状態を含み得る。

【 0 0 5 9 】

[0070]メインメモリーの中で、トランザクション内のすべてのそのような操作が、後段でさらに説明されるとおり、トランザクションコミットまで、隔離されたままに保たれ得る。例えば、コミット時に、トランザクションにおけるすべてのページ変更が、L S S バッファに自動的にフラッシュされ得る。例えば、中途終了時に、すべての変更が破棄され得る。例えば、システムトランザクションマネージャーが、トランザクションをコミットすること、およびトランザクションを中途終了することを行うように構成され得る。

10

【 0 0 6 0 】

[0071]例えば、システムトランザクションが、L L A M A によってサポートされる操作を介して開始されること、および終了されることが可能である。

[0072]例えば、「トランザクション開始」操作が、T B e g i n (o u t - T I D) として示され得る。例えば、トランザクション I D (T I D) によって識別されたトランザクションが開始され得る。このことには、T I D を、例示的な L L A M A キャッシュ層 (C L) マネージャーによって保持されるアクティブトランザクションテーブル (A T T) に入力することが関与し得る。

【 0 0 6 1 】

20

[0073]例えば、「トランザクションコミット」操作が、T C o m m i t (T I D) として示され得る。例えば、そのトランザクションが、アクティブトランザクションテーブルから取り除かれることが可能であり、そのトランザクションが、コミットされることが可能である。例えば、そのトランザクションにおけるページ状態変更が、マッピングテーブルにインストールされて、L S S バッファにフラッシュされることが可能である。

【 0 0 6 2 】

[0074]例えば、「トランザクション中途終了」操作は、T A b o r t (T I D) として示され得る。例えば、そのトランザクションが、アクティブトランザクションテーブルから取り除かれることが可能であり、変更されたページが、キャッシュの中で「トランザクション開始」にリセットされることが可能であり、変更は、まったくフラッシュされない。

30

【 0 0 6 3 】

[0075]本明細書で説明される例示的な技法によれば、A l l o c a t e および F r e e に加えて、U p d a t e - D 操作が、ページ状態を変更することをトランザクション内で許され得る。例えば、U p d a t e - R は、後段でさらに説明されるとおり、トランザクションを元に戻すことを複雑にする可能性があるので、使用されないことが可能である。

【 0 0 6 4 】

[0076]本明細書で説明される例示的な技法によれば、トランザクション操作はすべて、入力パラメーター、すなわち、T I D および a n n o t a t i o n を有し得る。例えば、T I D は、キャッシュ内のデルタに追加されることが可能であり、a n n o t a t i o n は、トランザクションにおいてアップデートされた各ページに追加されることが可能である (例えば、a n n o t a t i o n がフラッシュされているかのように)。フラッシュバッファの中にインストールされ、コミットされる場合、キャッシュ内のすべてのアップデートされたページは、それらのページのロケーションを記述するフラッシュデルタが先頭に付加されることが可能である (例えば、それらのフラッシュデルタが、トランザクションとは無関係にフラッシュされるかのように)。

40

【 0 0 6 5 】

[0077]B W - T R E E (例えば、J . L e v a n d o s k i 他、「T h e B w - T r e e : A B - t r e e f o r N e w H a r d w a r e P l a t f o r m s (B w - T r e e : 新しいハードウェアプラットフォームのためのBツリー)」、29th

50

IEEE International Conference on Data Engineering (ICDE 2013)、2013年4月8～11日参照)が、ユーザートランザクションがサポートされる(例えば、トランザクションコンポーネント202のために)ことを可能にし得る例示的なキー・バリューストアをもたらすことが可能である。例えば、BW-TREEは、LSNを管理し、先行書込みログ(WAL)プロトコルを執行し、DEUTERONOMYデータコンポーネント(DC)(J. Levandoski他、「Deuteronomy: Transaction Support for Cloud Data (Deuteronomy: クラウドデータのためのトランザクションサポート)」、Conference on Innovative Data Systems Research (CIDR)(2011年1月)、123～133ページ、およびD. Lomet他、「Unbundling Transaction Services in the Cloud (クラウドにおけるトランザクションサービスのアンバンドリング)」、Conference on Innovative Data Systems Research (CIDR)、2009年参照)によって予期されるとおり、チェックポインティング要求に応答することが可能である。本明細書の説明は、BW-TREEが、LLAMAなどの例示的なシステムを使用している場合に、どのように前述のことを実現するかを扱うことを含む。

【0066】

[0078] Update-DおよびUpdate-R LLAMA操作に対する「データ」コンテンツは、キー、LSN、およびキー・バリューストアの「データ部分」を含み得る。例えば、BW-TREEは、このため、これらの操作を介して、キー・バリューストアを実施し、LSNを介して冪等性をもたらし、Update-Dを介してインクリメンタルアップデートを実行し、Update-Rを介してBW-TREEのページ統合を実行し、LLAMA Read操作またはFlush操作を使用して読取りまたは書込みのためにページにアクセスすることが可能である。例えば、システムは、アップデートデルタレコード操作および置換アップデート操作に基づいてアップデートを制御するように構成され得るレコードマネージャーを含み得る。

【0067】

[0079]例えば、アクセスメソッドが、そのアクセスメソッドがアップデート操作を介してLLAMAにもたらすデータの中にLSNを記憶することが可能である。さらに、フラッシュデルタの中に記憶されたFlush操作annotationパラメーターが、ページコンテンツを記述するさらなる情報をもたらすことが可能である。例えば、これらは、BW-TREEが先行書込みロギング(WAL)を執行することを可能にし得る。例えば、安定させる操作(例えば、Mk-Stable)が、ページをフラッシュした後、トランザクションログチェックポインティングのためにアップデートを安定させることが可能である。

【0068】

[0080]例えば、Allocate操作およびFree操作が、例示的なBW-TREE実施形態がその実施形態のツリーを成長させること、または縮小することを可能にし得る。例えば、BeginTrans(例えば、TBegin)およびコミット/中途終了(例えば、TCommit/TAbort)が、構造変更操作(SMO)を実行している際に予期されるアトミック性を可能にし得る。

【0069】

[0081]例えば、アップデート操作(例えば、Update-D/Update-R)は、「ユーザーレベル」データに限定されないことが可能である。例えば、BW-TREEが、Update-Dを使用して、システムトランザクションに関して、後段でさらに説明されるとおり、SMOを実施している際にBW-TREEの「マージ」デルタおよび「分割」デルタを書き込むことが可能である。

【0070】

[0082]本明細書で説明される例示的な技法によれば、キャッシュ層データ操作に関して

10

20

30

40

50

、ページアップデートは、図4に示されるとおり、デルタアップデートであれ、置換アップデート（例えば、図7に関連して後段でさらに説明される）であれ、コンペア・アンド・スワップ（CAS）操作を使用してマッピングテーブル304の中に新たなページ状態ポインター402をインストールすることによって達せられることが可能である。例えば、置換アップデート（例えば、Update-R（PID, in-ptr, out-ptr, data））が、所望される新たな状態と、LSS内のそのページのそれまでの状態のロケーションの両方を含み得る。例えば、新たなアップデートデルタ404（例えば、Update-D（PID, in-ptr, out-ptr, data））が、このLSSロケーションを既に含む、ページ302のそれまでの状態406をポイントする。

【0071】

10

[0083]例えば、そのようなラッチフリーアプローチは、ラッチすることによって生じる遅延を回避することが可能であるが、「楽観的」同時実行制御メソッドと同様に、そのアプローチ独自のペナルティを受ける可能性があり、すなわち、CASが失敗する可能性があり、その場合、アップデートは、再試行される。例えば、CASの操作を適宜、再試行することは、障害が生じた場合に例示的なLLAMA実施形態が示し得るとおり、例示的なLLAMAユーザーに任せられる可能性がある。

【0072】

[0084]本明細書で説明される例示的な技法によれば、データがキャッシュ（例えば、312）内にある場合、ブロックする操作は存在しない可能性があるが、二次ストレージからページを読み取ることは、キャッシュ内にページが出現するのを待つことが関与する可能性がある。マッピングテーブル（例えば、マッピングテーブル304）が、前述したとおり、キャッシュされたページに関してさえ、LSSページをポイントして、ページが、効果的なキャッシュ管理のためにキャッシュとLSSの間で移動されることを可能にする。

20

【0073】

[0085]本明細書で説明される例示的な技法によれば、ページがフラッシュされる際、例示的なLLAMA実施形態が、キャッシュ（例えば、312）内で代表されているものが、LSS（例えば、314）内にあるものと合致することを確実にすることが可能である。このため、フラッシュデルタは、PIDとLSSオフセットの両方をフラッシュデルタの中に含めることが可能であり、そのデルタを、ページ302の先頭に付加することによってLSSバッファおよびキャッシュ（例えば、312）の中に含めることが可能である。

30

【0074】

[0086]本明細書で説明される例示的な技法によれば、例示的なLLAMA実施形態は、デルタアップデートをサポートし得るため、ページ状態は、不連続の断片を含み得る。この特徴をフラッシュする活動と組み合わせることにより、最近のアップデートがキャッシュ内だけに存在し得る一方で、キャッシュ内ページがページの状態の一部をLSSの中に有する（以前にフラッシュされていて）ことがもたらされ得る。このことが生じると、次のフラッシュのストレージ費用を低減することが可能であり得る。

【0075】

40

[0087]このため、例示的なLLAMA実施形態が、以前のフラッシュ以来の変更だけを含むデルタを書き込むことによって、そのようなページをフラッシュすることが可能である。例えば、キャッシュ内の複数のアップデートデルタがすべて、隣接する形態のデルタ（本明細書で「Cデルタ」と呼ばれ得る）を、LSS内のページのリマインダーに対するポインターと一緒に書き込むことによって、フラッシュするために隣接するようにされることが可能である。このため、そのページ全体が、LSS内で、ただし、場合により、いくつかの断片で、アクセス可能であり得る。

【0076】

[0088]本明細書で説明される例示的な技法によれば、Flush操作は、このようにしてときとともにフラッシュされてきたいくつかの部分の有し得るキャッシュされたページ

50

状態を見ることが可能であり、別々の断片およびそれらの断片の L S S アドレスが代表された、キャッシュされたページをもたらす。本明細書で説明される例示的な技法によれば、任意の時点で、F l u s h が、不連続なページ断片のコンテンツを隣接して（かつ冗長に）書き込むことによって、L S S ストレージの中にこれらの断片を一緒にもたすことが可能である。例えば、ユーザーが、異なる読取りアクセス費用のため、L S S がディスクストレージを使用する場合、隣接性を所望する一方で、L S S がフラッシュストレージを使用する場合、断片を別々のままにしておくことに積極的である可能性がある。

【 0 0 7 7 】

[0089]本明細書で説明される技法によれば、ページがフラッシュされる際、システムが、そのフラッシュに先立って、ページのどのような状態がフラッシュされているかを知ることが望ましい可能性がある。例えば、このことは、システムが、ページを単にラッチし、フラッシュを実行することが可能であるので、ラッチを使用して容易に確かめられ得る。しかし、ラッチフリーのアプローチにおいて、システムは、ページに対するアップデートがフラッシュされることを、そのページがフラッシュされることに先立って、防止することに相当な困難を有し得る。例えば、このことは、先行書込みログプロトコルの執行の際、またはそのフラッシュが構造変更の一環として行われる場合に、問題をもたらし得る。例えば、不適切なフラッシュが、それらのフラッシュの C A S を実行する際、失敗することが望ましい可能性がある。このため、本明細書で説明される例示的な技法によれば、C A S の中にフラッシュされるべきページ状態に対するポインターが使用されることが可能であり、そのポインターは、その場合、その特定の状態をキャプチャするに過ぎないことが可能であり、フラッシュが完了する前に状態がアップデートされている場合、失敗することが可能である。しかし、このことは、他の問題を生じる可能性がある。

【 0 0 7 8 】

[0090]本明細書で説明される例示的な技法を研究する際、キャッシュ管理を実行する場合、および L S S にページをフラッシュする場合に有利であり得る強い不変条件(strong invariant)の種類を決定することに困難を経験した。例えば、不変条件は、以下のような特性を含み得る。すなわち、

L S S に正常にフラッシュされたページは、フラッシュされたものとしてキャッシュ内で即時に見られ、そのページのフラッシュされた状態は、すべての後の状態のフラッシュに先行して L S S I / O バッファの中に入っている。フラッシュが失敗したページは、キャッシュ内でフラッシュされたものとして出現せず、L S S を見る際、そのフラッシュが成功しなかったことが明らかとなる。

【 0 0 7 9 】

[0091]例えば、2つの代替のアプローチが、以下を含み得る。すなわち、

1) フラッシュの成功が、C A S をまず実行することによって確実にされ得る。C A S が成功すると、ページが、L S S に書き込まれ得る。例えば、前述のことが行われた場合、競合条件が、信頼できる L S S 回復を損なう可能性がある。例えば、より早期のフラッシュに依存するページが、その後、フラッシュされる可能性があり、ただし、この「後の」フラッシュは、システムクラッシュより前に L S S に書き込まれることに成功する一方で、「より早期の」フラッシュは、完了するのが遅すぎ、安定した L S S 内に出現しない。この状況は、ある種の因果関係を損なう可能性がある。

2) フラッシュされることが所望されるページのページ状態が、キャプチャされ、L S S バッファに書き込まれることが可能である。次に、C A S が試みられることが可能であり、C A S が失敗することが可能である。このため、ページは、システムがクラッシュした場合にそのフラッシュが成功したか、または失敗したかを区別するための指示をまったく伴わずに、L S S に書き込まれる。例えば、様々な時点で L S S に書き込まれた複数のそのようなページが存在し得る。例えば、失敗した C A S よりも早期に L S S 内に出現するページの後の状態が書き込まれ得る。前述したとおり、その状態は、より後に始まったが、より早期のフラッシュより前にバッファースロットを獲得している。

【 0 0 8 0 】

[0092]本明細書で説明される例示的な技法によれば、前述したジレンマが、後段で説明されるとおり解決され得る。例えば、C A Sが十分に早期に実行された場合、ページの状態をログバッファにコピーすることに先立って、そのフラッシュが成功するか否かが判定され得る。このため、例示的なフラッシュ手順が、以下のとおり実行され得る。すなわち、

ステップ1：フラッシュされることが意図されるページの状態を識別する

ステップ2：その状態を書き込むべきL S Sバッファ内のスペースを差し押さえる

ステップ3：C A Sを実行して、そのフラッシュが成功するかどうかを判定する。このことを行うためにフラッシュデルタにおけるL S Sオフセットが獲得される（前述のステップ2において規定されるとおり）

ステップ4：ステップ3が成功した場合、保存されるべき状態をL S Sに書き込む。この状態がL S Sに書き込まれている間、本明細書で説明される例示的なL L A M A技法が、バッファがL S S二次ストレージに書き込まれるのを防止することが可能である。

ステップ5：ステップ3が失敗した場合、「失敗したフラッシュ」を示す指示をバッファ内の確保されたスペースに書き込む。このことは、ストレージを消費し得るが、いずれのフラッシュが成功したか、または失敗したかについての曖昧さを解決する。

【0081】

[0093]この例示的な手順の結果は、L S Sが、回復中、失敗したC A Sの結果であるページを見ない可能性があることである。例えば、このことは、L S S内に後に出現する（「ログ」におけるページの位置の点で）ページが、L S Sログにおいてそのページのより早期のすべてのインスタンスよりも、そのページの後の状態であるという特性を保つこともする。

【0082】

[0094]本明細書で説明される例示的な技法によれば、例示的なL L A M A実施形態が、その実施形態のメモリー制約を満たすようにキャッシュ・アンド・スワップアウト・データを管理することが望ましい可能性がある。例えば、例示的なL L A M A実施形態は、デルタアップデート、置換アップデート、およびフラッシュを意識していることが可能であり、以上のそれぞれを認識することが可能である。しかし、例示的なL L A M A実施形態は、汎用であるものとされる場合、ページのコンテンツについて何も知らない。このため、例示的なL L A M A実施形態は、ページの中にL S Nを保持することによって、アクセスメソッド層がトランザクションをサポートしているかどうかを意識していない。このため、もたらされ得る問題は、例示的なL L A M A実施形態が、L S Nを見ることなく、先行書き込みログプロトコルを執行する可能性がある場合、どのようにキャッシュスペース管理（ページを強制退去させることを含む）を提供し得るかに関する潜在的な疑問を含む。

【0083】

[0095]本明細書で説明される例示的な技法によれば、既にフラッシュされているデータは、キャッシュからドロップされ得る。例えば、ページのアップデート・インプレースが行われるシステムが、最近、アップデートされた、「汚れた」ページをスワップアウトすること（キャッシュからドロップすること）を防止され得る。しかし、デルタアップデートのため、例示的なL L A M A実施形態は、ページのいずれの部分が既にフラッシュされているかを判定することが可能である。例えば、それぞれのそのような部分が、フラッシュデルタを用いて記述されることが可能であり、それらのフラッシュされた部分が、キャッシュから「スワップアウト」されることが可能である。

【0084】

[0096]ページの一部を「スワップアウト」する際、ストレージを単に割当て解除し、再使用することは、スワップアウトされた部分に対するダングリングリファレンスを残し得るので、望ましくない可能性がある。このため、本明細書で説明される例示的な技法によれば、ページのいずれの部分のスワップアウトされているかを記述するデルタが使用され得る。

【0085】

[0097]例えば、完全にスワップアウトされたページに関して、マッピングテーブル 304 内のそのページのメインメモリアドレスが、そのページの最新のフラッシュデルタからの LSS ポインターで置き換えられることが可能である。

【0086】

[0098]図 5 は、例示的な部分的ページスワップアウト、および例示的な部分的スワップデルタを示す。例えば、部分的にスワップアウトされたページに関して、CAS が、「部分的スワップ」デルタレコード 502 を挿入するのに使用され得る。例えば、このデルタレコード 502 は、そのページが部分的にスワップアウトされている（例えば、したがって、そのページのいずれの部分も正常にアクセスされ得ない）ことを示すことが可能であり、そのページの欠落した部分 506 を探し出すための LSS 内のロケーション情報を示すフラッシュデルタレコード 504 をポイントすることが可能である。例えば、「部分的スワップ」デルタ 502 が CAS を用いてインストールされると、ドロップされているページの部分に関するメモリーは、後段でさらに説明されるとおり、例示的なエポック機構を使用して解放され得る。

【0087】

[0099]例えば、ページマネージャーが、第 1 のページに関連付けられたページ状態の先頭に部分的スワップデルタレコードを付加することであって、その部分的スワップデルタレコードは、第 1 のページの欠落した部分の二次ストレージ内のロケーションを示すフラッシュデルタレコードのストレージロケーションを示す二次ストレージアドレスを含むことに基づいて、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第 1 のページの一部分をスワップする操作を開始するように構成され得る。

【0088】

[0100]例えば、ページマネージャーは、エポック機構を使用して、第 1 のページのその部分に関連付けられたキャッシュ層ストレージに関する解放操作を開始するようにさらに構成され得る。

【0089】

[0101]本明細書で説明される例示的な技法によれば、このアプローチは、いくつかの有用な特徴をユーザーに有利に提供することが可能である。例えば、そのような例示的な LLAMA 実施形態のキャッシュ層（例えば、312）が、ページの実際のコンテンツに関する知識なしにメモリーを取り戻すことが可能である。例えば、フラッシュされたページ、およびページのフラッシュされた部分をドロップすることに、I/O 操作がまったく関与しないことが可能である。例えば、部分的にフラッシュされたページをメインメモリーの中に戻すことに、LSS の中に複数の部分を有する完全にフラッシュされたページに関する場合よりも少ない回数の LSS 読取りが関与することが可能である。

【0090】

[0102]例えば、いくつかの例示的なキャッシュ管理戦略が、キャッシュストレージを管理するのに使用され得る（例えば、最も長く使われていない（LRU）、LRU(k)、Clock など（例えば、W. Effelsberg 他、「Principles of database buffer management（データベースバッファ管理の原理）」、ACM Transactions on Database Systems (TODS)、Vol. 9、Issue 4（1984 年 12 月）、560～595 ページ、および E. O'Neil 他、「The LRU-K page replacement algorithm for database disk buffering（データベースディスクバッファリングのための LRU-K ページ交換）」、Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data (SIGMOD '93)、297～306 ページ参照）。これらの例は、さらなる帳簿処理（bookkeeping）を含み得るが、実質的な困難はもたらさないことが可能である。

【0091】

[0103]本明細書で説明される例示的な技法によれば、そのような例示的なラッチフリー

のアプローチを使用して、操作は、ページとページ状態の両方を、ページとページ状態の両方が「ガページ」と指定された後でさえ、検査していることが可能である。例えば、従来の「ラッチ」を使用していない場合、システムは、1) ページ状態全体を置き換え、それまでの状態を、別の操作が読取っている間に割当て解除するUpdate-R操作、または2) マッピングテーブル内のページを、別の操作が検査している間に「解放」するDe-allocate操作のいずれかを防止することに失敗する可能性がある。

【0092】

[0104] 本明細書で説明される例示的な技法によれば、ストレージもPIDも、別の操作がストレージにもPIDにもアクセスしている可能性がまったくなくなるまで、再使用されることを許されないことが可能である。このため、「解放されたリソース」と「再使用可能なリソース」の間の区別が確立され得る。例えば、「解放されたリソース」が、操作によってガページと指定されている。例えば、「再使用可能なリソース」が、解放されており、他のいずれの操作によってもアクセス可能でないことが確実にされることが可能である。例えば、エポックが、割当て解除されたオブジェクトがあまりにも早く再使用されるのを防ぐのに使用され得る（例えば、H. Kung 他、「Concurrent manipulation of binary search trees (バイナリサーチツリーの同時操作)」、ACM Transactions on Database Systems (TODS)、Vol. 5、Issue 3 (1980年9月)、354～382ページ参照)。

【0093】

[0105] 本明細書で説明される例示的な技法によれば、すべての操作は、PIDまたはページ状態にアクセスするのに先立って、現在のエポックEに登録することが可能であり、そのようなアクセスが完了すると、Eを抜けることが可能である。例えば、操作が、解放されたリソースを、現在のエポックのリスト上に常に見込むことが可能であり、現在のエポックは、E（その操作が参加したエポック）であること、または現在のエポックが先に進んでいる場合、後のエポックであることが可能である。例えば、Eのリスト上のいずれのリソースも、Eに登録されたすべての操作が抜けるまで、再使用されないことが可能である。

【0094】

[0106] 例えば、エポックに番号が付けられることが可能であり、時々、新たなエポックE+1が「現在の」エポックになることが可能である。このため、新たな操作が、今やE+1である現在のエポックに引き続き登録することが可能である。例えば、エポック機構不変条件は、以下のとおりである。すなわち、エポックE+1または後のエポックにおけるいずれの操作も、エポックEにおいて解放されたリソースを見ていること、および使用していることはあり得ない。

【0095】

[0107] このため、この不変条件に基づいて、すべての操作がEを抜けると、Eにおいて解放されたリソースにアクセスすることができるアクティブな操作は存在しない。図6は、例示的な2つのエポック602、604、およびそれらのエポック602、604のそれぞれのガページコレクションリスト606、608を示す。図6に示されるとおり、ガページコレクション項目610が、エポック1(602)における「スレッド1」に関連付けられ、ガページコレクション項目612が、エポック1(602)における「スレッド2」に関連付けられ、ガページコレクション項目614が、エポック2(604)における「スレッド3」に関連付けられる。図6に示されるとおり、エポック2(604)のガページコレクションリスト608内のガページコレクション項目616が、エポック1(602)の「スレッド1」に関連付けられる。

【0096】

[0108] 例えば、「スレッド1」および「スレッド2」がエポック1(602)を抜けると、エポック1(602)において解放されたリソース（例えば、ガページコレクション項目610およびガページコレクション項目612）にアクセスすることができるアクテ

ィブな操作は存在しない。

【 0 0 9 7 】

[0109]例えば、第1のエポックマネージャーが、第1のプロセッサ操作によってページ情報にアクセスすることに先立って、第1のエポック登録リスト内に第1のプロセッサ操作を登録することを開始するように構成され得る。

【 0 0 9 8 】

[0110]第1のエポックマネージャーは、第1のプロセッサ操作によって解放された1つまたは複数のリソースを、第1のエポックガベージコレクションリスト内に書き込むように構成され得る。第1のエポックマネージャーは、第1のエポック登録リストが、現在、登録されているプロセッサ操作を含むようになるまで、第1のエポックガベージコレクションリスト内に書き込まれた、書き込まれたリソースの再使用をブロックすることが可能である。

10

【 0 0 9 9 】

[0111]本明細書で説明される例示的な技法によれば、例示的なLLAMA実施形態が、ログ構造化ファイルシステム(LFS)と似通ったログ構造化の様態で二次ストレージ(例えば、フラッシュストレージ)上のデータを編成することが可能である(例えば、M. Rosenblum他、「The Design and Implementation of a Log-Structured File System(ログ構造化ファイルシステムの設計および実施)」、ACM Transactions on Computer Systems(TOCS)、Vol. 10、Issue 1、1992年2月、26~52ページ参照)。このため、各ページフラッシュは、フラッシュ上のページの位置を再配置する。例えば、このことが、本明細書で説明される例示的なマッピングテーブル304を使用するさらなる理由を与えることが可能である。例えば、ログ構造化ストレージが、1ページ当たりの書込みの回数を有利に減らし、それらの書込みを「順次」にすることが可能である。このため、多くのランダムな書込みが、1回の大きなマルチページ書込みに変換され得る。

20

【 0 1 0 0 】

[0112]前述したとおり、「論理ページ」が、基本ページと、そのページに対するアップデートを示す0または1つ以上のデルタレコードとを含むことが可能であり、これにより、ページが、フラッシュされる際に断片でフラッシュに書き込まれることを可能にする。このため、フラッシュ上の論理ページは、潜在的に、ポインターとしてファイルオフセットを使用して一緒にリンクされた異なる物理デバイスブロック上にあるレコードに対応することが可能である。さらに、物理ブロックが、複数の論理ページからのレコードを含み得る。図7aは、フラッシュ314上の例示的なログ構造化ストレージ編成700aを示す。

30

【 0 1 0 1 】

[0113]例えば、論理ページが、フラッシュ上の連鎖の先頭(順次ログ702内のオフセットがマッピングテーブル304から獲得され得る)から開始して、リンクされたレコードをたどることによって、フラッシュ314からメモリー(例えば、RAM312)の中に読み込まれ得る。例えば、フラッシュ314からメモリー312の中に対応する「論理ページ」を読み込むために、オフセット704が、デルタレコード706にアクセスするためにマッピングテーブル304から獲得されて、現在の状態、および基本ページ708が獲得されることが可能である。

40

【 0 1 0 2 】

[0114]例えば、フラッシュ314から対応する「論理ページ」をメモリー312の中に読み込むために、オフセット710が、デルタレコード712にアクセスするためにマッピングテーブル304から獲得されて、デルタおよびリンクが獲得され、第2のデルタレコード714にアクセスが行われ、その後、基本ページ716にアクセスが行われることが可能である。

【 0 1 0 3 】

50

[0115]例えば、フラッシュプロセスが、同一の論理ページの複数のデルタレコードを、それらのデルタレコードと一緒にフラッシュされる場合に有利に統合して、フラッシュ上の隣接するCデルタにすることが可能である。さらに、論理ページが、メモリー内で統合された後にフラッシュされる場合、フラッシュ上で統合されることが可能であり、このことは、ページ読取りパフォーマンスを有利に向上させ得る。

【0104】

[0116]図7bは、第1のストレージオブジェクト746（例えば、図7bにおいて、複数の以前に先頭に付加されたデルタレコードを有する基本ページ742を含む）の物理アドレスを、ページ742の新たな状態744（例えば、ページ742を、以前に先頭に付加されたデルタレコードと統合することからもたらされる）の物理アドレスで置き換えることに基づいて、ページ742のそれまでの状態740をページ742の新たな状態で置き換えることを示す例示的なマッピングテーブル304を示す。

10

【0105】

[0117]例えば、図7cに示されるとおり、ページ742のそれまでの状態740をページ742の新たな状態744で置き換えることは、複数のデルタレコードを統合して、隣接するCデルタ750にすることを含むことが可能であり、そのCデルタ750が、次に、基本ページ742と一緒にフラッシュされ得る。

【0106】

[0118]例えば、ページ742のそれまでの状態740をページ742の新たな状態744で置き換えることは、マッピングテーブル304上のアトミック・コンペア・アンド・スワップ操作を介して、現在のページ742の変更されたバージョンを生成すること、または現在のページ742を置き換えるための別のページを決定すること、および現在のページ742の物理アドレスを、ページ742の新たな状態744（例えば、置換のための変更されたバージョンまたその他のページ）の物理アドレスで置き換えることを含み得る。

20

【0107】

[0119]例えば、図7bの特徴と図7cの特徴の間の区別として、二次ストレージにページを書き込む際、LLAMAは、図7cに示される統合を実行することが可能であるが、図7bの統合を実行することは、Update-Rを実行しているアクセスメソッドに依存する。

30

【0108】

[0120]本明細書で説明される例示的な技法によれば、例示的なLLAMA実施形態が、完全にラッチフリーであり得る。さらに、専用のスレッドは、I/Oバッファをフラッシュするのに使用されないことが可能である。というのは、そうすることが、スレッド作業負荷をバランスがとれた状態に保つことを複雑にする可能性があるからである。このため、すべてのスレッドが、このバッファを管理することに参加することが可能である。例えば、従来のアプローチは、ラッチを使用してきた。しかし、そのような従来の技法は、バッファ内のスペースを割り当てている間、ラッチすることに限られて、データ転送に先立ってそのラッチを解放する可能性があり、データ転送は、その後、並行に進められる可能性がある。

40

【0109】

[0121]本明細書で説明される例示的な技法によれば、例示的なLLAMA実施形態が、本明細書で説明される例示的なシステムの別の場所で行われるとおり、アトミック性のためにCASを代わりに使用して、バッファスペース割当てのための従来のラッチを回避することが可能である。例えば、このことには、CASが実行される状態を定義することが関与する。例えば、バッファ状態の不変の部分は、その部分のアドレス（Base）と、サイズ（Size）とを含み得る。例えば、バッファ内で使用されるストレージの現在の最高水位線が、Baseを基準としたOffsetを用いて追跡され得る。例えば、バッファの使用を求める各要求が、ページフラッシュのためのスペースSizeを確保しようとする作業から始まることが可能である。

50

【 0 1 1 0 】

[0122]本明細書で説明される例示的な技法によれば、バッファー内でスペースを確保するのに、スレッドが、現在の `Offset` を獲得し、`Offset + Size` を計算することが可能である。例えば、`Offset + Size < BufferSize` の場合、要求はバッファー内に記憶され得る。例えば、スレッドは、比較値としての現在の `Offset`、および新たな値としての `Offset + Size` と一緒に `CAS` を発行することが可能である。その `CAS` が成功した場合、`Offset` は、その新たな値に設定されることが可能であり、そのスペースは、確保されることが可能であり、バッファー書込みは、データをバッファーに転送することが可能である。

【 0 1 1 1 】

[0123]本明細書で説明される例示的な技法によれば、このロジックは、バッファー内のスペース割当てを扱うことが可能である。例えば、バッファーに書き込むこと、および複数のバッファーを管理することには、後段でさらに説明される、`CAS` 状態におけるより多くのことが関与する可能性がある。

【 0 1 1 2 】

[0124]二次ストレージにバッファーを書き込む際、`Offset + Size > BufferSize` である場合、スレッドのレコードを保持するのにバッファー内に不十分なスペースしか存在しない。この場合、スレッドは、バッファーを封印して、これにより、バッファーを、もはや使用されるべきでないものにし、二次ストレージに書き込まれるべく準備されたものにすることが可能である。この条件は、フラッシュバッファースタートにおける「`Sealed`」ビットを用いて追跡され得る。例えば、`CAS` が、「`Sealed`」ビットを `F` (例えば、偽) から `T` (例えば、真) に変更することが可能である。例えば、封印されたバッファーが、もはやアップデートされないことが可能であり、封印されたバッファーに出会ったスレッドは、異なる (封印されていない) バッファーを探す。

【 0 1 1 3 】

[0125]本明細書で説明される例示的な技法によれば、封印されたバッファーは、新たなアップデート要求をもはや受け付けないことが可能である。しかし、例示的なシステムは、バッファースペースを獲得することにすべて成功しているそれまでの書込み側が、書込み側のデータをバッファーに転送することを終えていることをまだ保証されていない可能性がある。本明細書で説明される例示的な技法によれば、「`Active`」カウンタが、バッファーにデータを転送している書込み側の数を示すことが可能である。例えば、バッファー内にスペースを確保する際、書込み側の `CAS` は、`Offset`、`Sealed`、および `Active` を表す値を含み得る。例えば、書込み側の `CAS` が、この構造を獲得し、その `CAS` のペイロードサイズを `Offset` に追加し、「`Active`」を 1 だけインクリメントすることが可能であり、`~ Sealed` である場合、この状態をアップデートして、スペースを確保するように `CAS` を実行することが可能である。例えば、書込み側が終了している場合、書込み側は、この状態を再獲得し、「`Active`」を 1 だけデクリメントすることが可能であり、その変更を生じさせるように `CAS` を実行することが可能である。例えば、操作は、失敗した場合、必要に応じて再び行われ得る。

【 0 1 1 4 】

[0126]例えば、バッファーは、`Sealed` であり、かつ `Active = 0` である場合、フラッシュ可能であり得る。例えば、この条件をもたらし書込み側が、`I/O` を開始することを担うことが可能である。例えば、`I/O` が完了すると、バッファーの `Offset` および `Active` ユーザーがともに 0 に設定されることが可能であり、バッファーは、`unSealed` であり得る。

【 0 1 1 5 】

[0127]本明細書で説明される例示的な技法によれば、複数のバッファーに関して、複数のバッファーのセットの中のそれらのバッファーの各々が、前述した状態を有する。図 8 は、例示的な完全なフラッシュバッファースタート 800 を示す。図 8 の例に示されるとおり、バッファー当たりの状態 802 が、32 ビットを含むことが可能であり、次の書込みの

10

20

30

40

50

ためのオフセットに関する24ビット804と、アクティブな書込み側の数に関する7ビット806と、「封印ビット」標識のための1ビット808（例えば、封印されたバッファを示す）を含む。例えば、現在のアクティブなバッファ番号（CURRENT）810が、現在、アクティブなバッファを示すことが可能である（例えば、図示されるとおり、8ビットに関する）。

【0116】

[0128]例えば、バッファは、ラウンドロビンスタイルでアクセスされて、使用されることが可能であり、したがって、1つのバッファが封印されるにつれ（封印ビット標識808によって示される）、本明細書の例示的な技法は、バッファ「リング」の中の次のバッファに進むことが可能である（例えば、CURRENT810を使用して）。本明細書で説明される例示的な技法によれば、CURRENT810が、セットのバッファのうちのいずれが現在、新たな書込み要求を受け付けているかを示すのに使用され得る。

10

【0117】

[0129]本明細書で説明される例示的な技法によれば、現在、アクティブなバッファをSEALする（例えば、「封印ビット」標識808を介して）スレッドは、そのバッファをSEALする際、CURRENT810をアップデートすることとする。例えば、このスレッドは、その後、次のCURRENTバッファを選択することが可能である。例えば、バッファI/Oが完了すると、I/Oスレッドは、バッファを封印解除することが可能であるが、現在のバッファの役割をしている別のバッファが存在し得るので、CURRENT810を設定しない可能性がある。

20

【0118】

[0130]LSSは、ログ構造化ストアであり、したがって、概念的に「付加専用(append only)」である。例えば、LSSの実現には、任意の通常のログ構造化ファイルシステム（LFS）の場合と同様に、ページの新たなバージョンを付加するために絶えずスペースを取り戻すことが関与し得る。例えば、この技法は、本明細書で「クリーニング」（例えば、前掲のM. Rosenblum他参照）と呼ばれ得る。

【0119】

[0131]例示的なページの異なるバージョンは、異なる寿命を有し得るので、再使用することが望ましい可能性がある例示的な「ログ」の古い部分が現在のページ状態を含むことが可能である。例えば、例示的なログのこの「古い」セクションを再使用するのに、依然として最新のページ状態が、ログのアクティブな末尾に移動されることが可能であり、それらの状態をその末尾に付加して、そのより古い部分がその後の使用のためにリサイクルされ得るようにする。例えば、クリーニングすることのこの副作用は、書込みの回数を増加させ得る（このことは、本明細書で「ライトアンプリフィケーション」と呼ばれ得る（例えば、X. - Y. Hu他、「Write amplification analysis in flash-based solid state drives（フラッシュベースの固体状ドライブにおけるライトアンプリフィケーション分析）」、In Proceedings of SYSTOR 2009: The Israeli Experimental Systems Conference（SYSTOR '09）、Article No. 10参照））。

30

40

【0120】

[0132]例えば、クリーニングする作業は、単に編成され得る。例えば、ログは、最も古い部分（例えば、ログの先頭）が「クリーニング」されて、新たなページ状態が書き込まれるログのアクティブな末尾に新たなスペースとして追加される大きい「循環バッファ」として管理され得る。

【0121】

[0133]本明細書で説明される例示的な技法によれば、再配置される各ページは、そのページが再書込みされる際、隣接させられる（例えば、ページが、LSSストアに再付加されると、「再書込みされた」内容は隣接する）。このため、どれほど多くのインクリメン

50

トフラッシュを経たとしても、そのページのすべての部分は、現時点で隣接させられ、このため、L S S内のページのアクセス可能性を有利に最適化する。

【 0 1 2 2 】

[0134]本明細書で説明される例示的な技法によれば、デルタ（本明細書で「再配置デルタ」と呼ばれ得る）に対するC A Sが、そのページに関するマッピングテーブルエントリにおいて実行されて、新たなロケーションをもたらし、そのページのいずれの部分が再配置されたかを記述する（すなわち、その新たなロケーション情報をインストールするようにキャッシュを管理する）ことが可能である。例えば、同時に行われるアップデートまたはフラッシュが、このC A Sを失敗させる可能性があり、その場合、そのC A Sは、再び試みられる。

10

【 0 1 2 3 】

[0135]ストレージ効率は、ログ構造化ストレージシステムに有利な良い影響を及ぼすことが可能である。本明細書で説明される例示的な技法によれば、L S Sに割り当てられた任意の所与の量のスペースに関して、この技法がそのスペースをより効率的に使用するほど、この技法は、クリーニングを実行することが少なくなることが可能であり、このことには、より少ないページ移動が関与し得る。例えば、ページ移動は、ストレージに対するさらなる書込み（例えば、ライトアンプリフィケーション）をもたらす得る。

【 0 1 2 4 】

[0136]潜在的なL S Sストレージ効率に関して、フラッシュされるページ内に空いたスペースは存在しない。例えば、それらのページは、パックされた可変長ストリングとして書き込まれることが可能である（例えば、平均で、従来のB - T R E Eページは、69%だけしか利用されない可能性がある）。さらに、前回のフラッシュ以来のデルタだけが頻繁にフラッシュされる可能性があるため、1回のページフラッシュ当たり消費されるスペースがより少ないことが可能である。さらに、アップデートされたページをキャッシュからスワップアウトすることには、さらなるフラッシュが関与しない。というのは、キャッシュ内のメインメモリーが、それまでにフラッシュされたページの部分に関してだけ取り戻され得るからである。

20

【 0 1 2 5 】

[0137]アクセスメソッドの1つの例示的な態様は、アクセスメソッドが、そのような構造が成長すること、および縮小することを可能にする構造変更操作（S M O）を行うことである。例えば、S M Oは、通常のアップデートが、進行中のS M Oが存在する状況で正しく実行され、アトミック（すべてまたは無）であることが可能であるようにインデックスのアトミック変更を行う方法が存在するものと予期する。例えば、例示的なB W - T R E Eが、そのB W - T R E EのS M Oのための機構としてシステムトランザクションを活用することが可能である。

30

【 0 1 2 6 】

[0138]本明細書で説明される例示的な技法によれば、システムトランザクションの持続性が、ログを介して実現され得る。しかし、本明細書で説明される一部のログは、トランザクションログではなく、例示的なL S S「ページ」ストアであり、このことは、トランザクションシステムが通常、操作をログ記録するだけであり得ることから、いくぶん不十分であるように見える可能性がある。しかし、デルタアップデートを用いて、ページ状態は、前回のページフラッシュ以来のデルタアップデートだけをログ記録することによってログ記録され得る。コミット時に持続性は関与せず、したがって、コミットは、L S Sバッファを「強制」しない。しかし、本明細書で説明される例示的な技法によれば、トランザクションの結果を使用するすべての後の操作は、L S S内のトランザクションコミット後に生じることが確実にされ得る。

40

【 0 1 2 7 】

[0139]本明細書で説明される例示的な技法によれば、非トランザクション操作と同様に、すべてのトランザクション操作も、マッピングテーブル内のページポインターに対するC A Sを介してインストールされ得る。本明細書で説明される例示的な技法は、キャッシ

50

ュ内のコンテンツがL S S内で忠実に表されることを確実にすること、およびL S S内のコンテンツがキャッシュ内で忠実に表されることを確実にすることが可能である。このため、システムトランザクション内の実質的にすべてのアップデートが、フラッシュ操作を含み得る。例えば、すべてのシステムトランザクションアップデートが、L S Sバッファ内に記録されることが可能であり、このため、「ログ記録」されることが可能である。例えば、情報のその2つの表現は等価であることが可能であり、これにより、システムクラッシュが生じた場合に、キャッシュの状態が、L S Sによって安定してキャプチャされた最終のバッファ時点に忠実に再構築され得ることが確実にされる。

【0128】

[0140]この等価性は、S M Oの場合のように、アクションに複数のページが関与する場合、従来、問題がある可能性がある。例えば、B - L I N Kツリーにおけるノード分割S M Oは、新たなページを割り当てることと、その新たなページを参照するその新たなページの兄弟ページリンクポインターをアップデートすることをともに行う。例えば、ラッチベースのシステムにおけるS M Oは、通常、ラッチを使用して隔離をもたらして、マルチページS M Oの内部状態が、そのS M Oが完了するまで、キャッシュマネージャー内で可視でないようにすることが可能である。例えば、ラッチフリーの設計は、アクティブな（したがって、コミットされていない）トランザクションアップデートを隔離する能力が制限され得ることを意味する可能性がある。

【0129】

[0141]本明細書で説明される例示的な技法によれば、例示的なL L A M A実施形態が、ページに対する実質的に任意のアクセスを許すトランザクションインターフェースをもたらすことが可能である（すなわち、任意のページに対する操作がトランザクション内に入れられることが可能である）。しかし、トランザクション中にアップデートされたページが、そのトランザクションの外部の操作によるアクセスから保護されない可能性がある。しかし、本明細書で説明される例示的な技法によれば、完全に一般的な隔離能力が関与しないS M Oが設計され得る。例えば、図9が、S M Oトランザクションをキャプチャするために使用され得る例示的なトランザクションテンプレート900を示す。

【0130】

[0142]例えば、ステップ1（902）で、ページがマッピングテーブル内で割り当てられる、または解放される。ステップ2（904）で、ページが、必要に応じてアップデートされる。ステップ3（906）で、既存のページが、新たなページをインデックスの残りの部分に接続するように、または別のページをアップデートしている間に既存のページを取り除くようにアップデートされる。

【0131】

[0143]本明細書で説明される例示的な技法によれば、ノード分割（図9の例示的なテンプレートを使用する）のための新たなノードは、そのノードがツリーに接続され、トランザクションがコミットされる際、図9のステップ3まで、他のスレッドに可視ではない。このため、そのようなS M Oトランザクションは、アトミック性と隔離をともにもたすことが可能である。

【0132】

[0144]従来のトランザクションシステムといくぶん同様に、アクティブトランザクションテーブル（A T T）と本明細書で呼ばれ得るアクティブトランザクションテーブルが、システムトランザクションに関して保持され得る。例えば、A T Tは、アクティブなシステムトランザクションごとに、そのトランザクションのトランザクションi d（T I D）と、そのトランザクションの最新の操作のメモリーアドレスをポイントする（またはそれ以外で参照する）、そのトランザクションの直前の操作に対するポインター（「I P」（「直前」を表す）と本明細書で呼ばれ得る）を含むエントリーを含み得る。

【0133】

[0145]例えば、B e g i n T r a n s操作（例えば、T B e g i n）が、任意の先行するトランザクションよりも高いトランザクションi d（T I D）を有し、I PがN U L L

10

20

30

40

50

の値に設定された新たなエントリーをA T Tに追加することが可能である。例えば、トランザクション操作の実行が、そのI Pによって識別されるその操作に関するログレコードをポイントする、その操作に関する「ログレコード」を作成することが可能であり、I Pが、その新たな操作を参照するようにアップデートされることが可能である。例えば、このことは、すべての「ログレコード」がメインメモリー内に存在して、トランザクションの操作に関する「ログレコード」をバックリンクする(backlink)役割をすることが可能である。さらに、本明細書で説明される例示的な技法によれば、システムトランザクション内の操作は、マッピングテーブルアップデートを介してキャッシュ状態しか(すなわち、L S Sバッファースタートではなく)変更しないことが可能である。本明細書で説明される例示的な技法によれば、これらのページは、トランザクションコミット時にフラッシュされ得る。本明細書で説明される例示的な技法によれば、トランザクションの終了(コミットまたは中途終了)が生じると、トランザクションは、A T Tから取り除かれ得る。

10

【0134】

[0146]例えば、システムトランザクションマネージャーが、キャッシュ層マネージャーによって保持されるアクティブトランザクションテーブル(A T T)に第1のトランザクションのトランザクション識別子(T I D)を追加するように構成され得る。例えば、トランザクションコミットマネージャーが、A T TからT I Dを取り除くこと、第1のトランザクションに関連付けられたページ状態変更をマッピングテーブルの中にインストールすること、および二次ストレージバッファに第1のトランザクションに関連付けられたページ状態変更をフラッシュすることを開始することに基づいて、第1のトランザク

20

【0135】

[0147]本明細書で説明される例示的な技法によれば、コミット操作の時点で、トランザクションによって変更されたページが、アトミック状態でL S Sバッファにフラッシュされる。例示的な技法として、これらのページ書込みは、L S S内のそのトランザクションに関する開始レコードおよび終了レコードで囲まれ得るが、このことには、中断されたトランザクションを元に戻すクラッシュ回復が関与し得る。例えば、そのような元に戻す回復には、元に戻す情報をL S Sに書き込むことが関与し得る。本明細書で説明される例示的な技法によれば、このことは、後段でさらに説明されるとおり、トランザクションによって変更されたすべてのページのアトミックフラッシュを、コミット時に実行することによって回避され得る。

30

【0136】

[0148]本明細書で説明される例示的な技法によれば、S M Oに依存する後のアクションは、S M Oトランザクションを記述する情報よりもL S Sバッファ内で後に出現する。このため、S M Oの状態が、システムトランザクションを扱っているスレッド以外のスレッドにキャッシュ内で可視になると、それらの他のスレッドは、L S Sにコミットされており、L S Sバッファ内に既に存在するS M Oに依存し得る。

【0137】

[0149]図9に示されるとおり、ステップ3は、「新たなページをインデックスの残りの部分に接続するように、または別のページをアップデートしている間に既存のページを取り除くように既存のページをアップデートする」を示す。このため、本明細書で説明される例示的な技法は、アトミックフラッシュを介してメインメモリー内でアップデートすること(トランザクション状態を可視にすること)とL S Sバッファ内でトランザクションをコミットすることの両方をカプセル化することを、これを実現するのにU p d a t e - Dに関する例示的な「コミット」能力を使用して(すなわち、アップデートをトランザクションコミットと組み合わせて)行うことが可能である。

40

【0138】

[0150]本明細書で説明される例示的な技法によれば、L S Sは、アップデート、およびアップデートのC A Sインストールを、トランザクションにおいて変更されたすべてのページのアトミックフラッシュと組み合わせることによって、トランザクションU p d a t

50

e - D「コミット」操作を可能にし得る。例えば、複数のページのコミット時のこのフラッシュは、個々のページフラッシュの場合と同様に実行され得る。例えば、L S Sバッファースペースが、トランザクションにおいて変更されたすべてのページに割り当てられることが可能である。次に、フラッシュデルタが先頭に付加されたU p d a t e - DデルタをインストールするC A Sが、実行され得る。そのC A Sが成功した場合、トランザクションにおいてアップデートされたページは、L S Sフラッシュバッファーに書き込まれ得る。そのトランザクションに関するすべてのページのフラッシュが完了した後、フラッシュプロセスは、フラッシュバッファーの書き込み側の数をデクリメントすることが可能である。例えば、トランザクションにおけるすべてのページに単一のユニットとしてスペースを割り当て、L S Sバッファー上の書き込み側デクリメントまで保持することが、L S Sストア内のそのトランザクションに関するアトミック性を確実にし得る。

10

【 0 1 3 9 】

[0151]例えば、トランザクションコミットマネージャーが、コンペア・アンド・スワップ(C A S)操作を介して、トランザクションに関連付けられたアップデートデルタレコードをマッピングテーブルの中にインストールするように構成されることが可能であり、そのアップデートデルタレコードの先頭には、フラッシュデルタレコードが付加される。例えば、トランザクションコミットマネージャーが、そのC A S操作が成功したかどうかを判定するように構成され得る。トランザクションコミットマネージャーが、そのC A S操作が成功したと判定した場合、トランザクションコミットマネージャーは、トランザクションにおいてアップデートされたページを二次ストレージフラッシュバッファーに書き込む書き込み操作を開始することが可能である。

20

【 0 1 4 0 】

[0152]本明細書で説明される例示的な技法によれば、そのC A Sが失敗した場合、応答は、他のフラッシュ失敗の場合と同様に進められ得る。例えば、L S Sが、回復中、そのスペースを別のものと混同しないように割り当てられていたスペースが、V O I Dにされ得る。このため、例示的な回復プロセスは、システムトランザクションをまったく意識しないことが可能である。むしろ、システムトランザクションは、例示的なキャッシング層の能力のみであり得る。このため、システムクラッシュまたはシステム再起動の前後でT I D一意性を確実にすることなしに進めることが許容可能であり得る。

【 0 1 4 1 】

30

[0153]本明細書で説明される例示的な技法によれば、中途終了されたシステムトランザクションの操作は、回復が不完全なトランザクションを見ないので、キャッシュ内で元に戻され得る。このため、メインメモリー内で一緒にリンクされた、そのトランザクションに関するログレコードのバックチェーンがたどられることが可能であり、元に戻すことは、そのトランザクションに関するA T Tリスト上の操作の性質に基づいてもたらされ得る。例えば、デルタアップデートは、そのデルタを取り除くことによって元に戻されることが可能であり、割当ては、「解放する」ことで元に戻されることが可能であり、「解放する」は、「解放する」に先立つページの状態にページを復元することによって元に戻されることが可能である。「解放する」を元に戻すことを除いて、操作成功を記述する情報を越えた、これらの操作に関する特別の情報は、まったく所望されない可能性がある。

40

【 0 1 4 2 】

[0154]本明細書で説明される例示的な技法によれば、トランザクション内で行われるアクションは、ストレージおよびマッピングテーブル・ページ・エントリー(P I D)を割り当てること、および解放することを含め、暫定的である。例えば、トランザクション実行中、P I Dが、割り当てられること、または解放されることが可能であり、U p d a t e - Dデルタが、生成されることが可能である。例えば、これらのリソースの管理は、本明細書で説明されるとおり、エポック機構に基づいて実現され得る。例えば、S M Oは、単一のユーザー操作要求内で実行されるので、スレッドは、そのトランザクションの持続時間にわたってそのスレッドのエポック内に留まり得る。

【 0 1 4 3 】

50

[0155]本明細書で説明される例示的な技法によれば、例示的なLLAMA実施形態が、トランザクションコミットまたはトランザクション中途終了に依存して、リソースを取り戻すことが可能である。例えば、コミット操作に関して、Free Page PIDが、現在のエポックに関するPID解放待ち状態リストに追加され得る。例えば、中途終了操作に関して、Allocate Page PIDが、元に戻す間に解放されることが可能であり、PID解放待ち状態リストに同様に追加され得る。例えば、Update-D操作に関して、アップデートデルタは、トランザクションが中途終了された場合、現在のエポックに関するストレージ解放待ち状態リストに追加され得る。

【0144】

[0156]本明細書で説明される「クラッシュ回復」とは、一般に、「トランザクション回復」を指す。本明細書で説明される「チェックポインティング」とは、一般に、トランザクションログを管理するのに使用されるチェックポインティングを指す。むしろ、本明細書で説明される「クラッシュ回復」とは、LSS（例えば、ログ構造化ストア）が、ページのLSSのマッピングテーブル、およびそれらのページの状態を、システムクラッシュの時点に回復する例示的な技法を指し得る。この特定のタイプの回復ステップは、通常、従来のアップデート・インプレース・ストレージ・システムの関心事ではない。

【0145】

[0157]本明細書で説明される「クラッシュ回復」に関して、マッピングテーブルは、ある種の「データベース」と見なされ得る。例えば、このデータベースに対するアップデートは、LSSにフラッシュされるページ状態を含み得る。このため、すべてのページフラッシュは、その「マッピングテーブルデータベース」をアップデートすることが可能である。システムがクラッシュした場合、LSS「ログ」は、フラッシュされたページを、マッピングテーブルをアップデートするやり直しログレコードとして使用して、「マッピングテーブルデータベース」を回復するように再実行され得る。

【0146】

[0158]前述の戦略をサポートして、マッピングテーブルは、LSSアップデートを無期限に保持することを回避するように、周期的にチェックポイント処理されることが可能である。例えば、前述したLFSクリーニング技法が、この目的で使用され得る（すなわち、回復ログを短くして）が、そのような技法は、高速回復のために望ましいとされ得るよりも、相当に大きい回復ログ（LSSログ構造化ストア）を残す可能性がある。

【0147】

[0159]本明細書で説明される例示的な技法によれば、チェックポインティングのために有利な方策が使用され得る。例えば、例示的なLLAMA実施形態が、チェックポイント処理中に完全なマッピングテーブルを、2つの交互のロケーションのうちの1つに非同期的に、インクリメンタルに書き込むことが可能である。図10は、本明細書で説明される例示的な技法による例示的なチェックポイントデータ1000を示す。例えば、2つの交互のロケーションが、様々なエンティティのロケーションに関するその他の「最新の」情報を失う可能性があるシステムクラッシュの後でさえ、システムがそれらのロケーションを知っているように、異なる2つの「よく知られたロケーション」（WKL）として選択され得る。このため、クラッシュの時点で存在していたおりのシステムの状態に関する情報をポイントするポインターが、保存され得る（例えば、WKLを使用して）。例えば、2つのチェックポイントを使用することによって、ユーザーが、「ライブの」チェックポイントのアップデート・インプレースを行わないことが可能である。

【0148】

[0160]例えば、各ロケーションは、図10に示されるとおり、完全なマッピングテーブルに加えて、フラッシュログ1006の中に回復開始位置（RSP）1002およびガベージコレクションオフセットGC1004を記憶することが可能である。例えば、RSP1002は、マッピングテーブル304をコピーすることを開始した時点におけるLSSストア内の終了オフセットを含み得る。例えば、GCオフセット1004は、ガベージコレクション「フロンティア」に印を付けることが可能である。

【 0 1 4 9 】

[0161]本明細書で説明される例示的な技法によれば、より後のチェックポイントは、L S S オフセットが、仮想化されることによって単調に増加するので、より高いR S P 1 0 0 2 を有する。例えば、システムクラッシュの後、最高のR S P 1 0 0 2 で完了したチェックポイント処理が、回復されたマッピングテーブル3 0 4 の状態を初期設定するのに使用され得る。例えば、R S P 1 0 0 2 は、やり直し回復を開始するためのL S S 「ログ」(1 0 0 6) 内の位置を示す。最新の完了したチェックポイント処理を識別するのに、R S P 1 0 0 2 は、マッピングテーブル3 0 4 が完全にキャプチャされてしまうまで、チェックポイントに書き込まれない。このため、現在のチェックポイント処理が完了するまで、前の高いR S P (代替のロケーションからの) が最高のR S P 1 0 0 2 である。

10

【 0 1 5 0 】

[0162]本明細書で説明される例示的な技法によれば、マッピングテーブル3 0 4 をチェックポイントの一部として書き出すことは、キャッシュ内に存在するとおりのマッピングテーブル3 0 4 のバイトごとに対応するコピーではない。例えば、キャッシュされた形態のマッピングテーブル3 0 4 は、キャッシュされたページに関してマッピングテーブルエントリの中にメインメモリーポインターを有するのに対して、本明細書で説明される例示的な所望されるチェックポイントには、それらのページのL S S アドレスをキャプチャすることが関与する。別の例として、現在、割り当てられていないマッピングテーブルエントリが、それらのマッピングテーブルエントリをリスト項目として使用する空きリスト上に保持される。このため、空きマッピングテーブルエントリは、0、または直前の空きマッピングテーブルエントリのアドレス(その0 またはアドレスが空きリストに追加された時刻に基づく時間順序で) を有する。例えば、使用可能な空きリストは、本明細書で説明されるとおり、マッピングテーブルを非同期で「コピー」している間、キャプチャされないことが可能である。例えば、本明細書で説明されるマッピングテーブル3 0 4 のコピーは、非同期で、インクリメンタルに書き込まれ、このことは、通常の実行に対する影響を最小限に抑えるのを助けることが可能である。

20

【 0 1 5 1 】

[0163]本明細書で説明される例示的な技法によれば、例示的なL L A M A 実施形態が、L S S の現在の終了オフセットをR S P 1 0 0 2 としてまず保存することが可能であり、現在のL S S クリーニングオフセットをG C 1 0 0 4 として保存することが可能である。例えば、マッピングテーブル3 0 4 が、スキャンされる(例えば、進行中の操作と同時に) ことが可能であり、各P I D エントリ(最新のフラッシュデルタの中に記憶された) に関するページの最新のフラッシュが識別されることが可能であり、そのL S S アドレスが、そのマッピングテーブル3 0 4 エントリに関する例示的なチェックポイントの中に記憶され得る。例えば、エントリが空いている場合、そのエントリは、チェックポイントコピーにおいてゼロ化され得る。例えば、やり直し回復の終わりに空きリストが再構築され得る。さらに、マッピングテーブル3 0 4 をコピーすることが完了すると、それまでに保存されたR S P 1 0 0 2 およびG C 1 0 0 4 は、安定したチェックポイント領域に書き込まれることが可能であり、これにより、チェックポイント処理が完了する。

30

【 0 1 5 2 】

[0164]本明細書で説明される例示的な技法によれば、回復は、最高のR S P 1 0 0 2 を有するチェックポイント(すなわち、最新の完全なチェックポイント) に関するマッピングテーブル3 0 4 をキャッシュ内3 1 2 にコピーすることによって開始され得る。例えば、ログ1 0 0 6 が、次に、R S P 1 0 0 2 からL S S の終わりまで順方向に読み取られ得る。例えば、出会う各ページフラッシュが、そのページフラッシュがページ読取りの結果であるかのように、キャッシュ3 1 2 内にもたらされ得る。

40

【 0 1 5 3 】

[0165]例えば、ページのコンテンツが読み取られることが可能であり、デルタは、L S S 内のロケーションがフラッシュデルタの中で参照されるように設定されることが可能である。例えば、A l l o c a t e P a g e 操作に出会うと、割り当てられたP I D に関す

50

るマッピングテーブル 304 エントリーが、Allocate Page 操作によって予期されるとおり「空」に初期設定され得る。例えば、Free Page 操作に出会うと、マッピングテーブル 304 エントリーは、ZERO に設定され得る。例えば、LSS クリーナが、チェックポイントから読み取られた GC オフセット (1004) からログに対するガベージコレクションを再開することが可能である。

【0154】

[0166] 本明細書で説明される例示的な技法によれば、回復中、すべての空きマッピングテーブル 304 エントリーが、ZERO に設定され得る。例えば、再構築されたマッピングテーブル 304 がスキャンされ得る。例えば、ZERO エントリーに出会うと、その ZERO エントリーは、空きリストに追加されることが可能であり、空きリストは、スタックとして管理され得る（すなわち、再使用されるべき最初のエントリーは、リストに追加された最後のエントリーである）。これらの例示的な技法によれば、下位の PID は、再使用されることが可能であり（再使用の際の選好として）、このことは、テーブルサイズを、クラスタ化されて、小さく保つ傾向があり得る（少なくとも回復の結果）。さらに、それまでに使用された最高の PID を示す最高水位線が、マッピングテーブル内に保持され得る。例えば、空きリストが尽きると、テーブルの使用されていない部分から PID が追加されて、最高水位線がインクリメントされることが可能である。

【0155】

[0167] 本明細書でさらに説明されるとおり、図 11 は、ラッチフリーのログ構造化ストレージを管理するためのシステム 1100 のブロック図である。システム 1100 は、ハードウェア実施形態、ソフトウェア実施形態、またはハードウェア実施形態とソフトウェア実施形態の組合せとして実現され得ることが、データ処理の分野の業者には理解されよう。図 11 に示されるとおり、システム 1100 が、少なくとも 1 つのプロセッサ 1104 を含むデバイス 1102 を含み得る。デバイス 1102 は、任意に選択されたページ指向型アクセスメソッド 1110 に、ページデータストレージ 1112 に対するラッチフリーのアクセスを含むページデータストレージ 1112 に対するインターフェースアクセスをもたらすように構成され得る。データ不透明型インターフェース 1108 を含み得るデータマネージャー 1106 を含み得る。例えば、ページ指向型アクセスメソッド 1110 は、任意の恣意的なアクセスメソッドであり得る。例えば、ページデータストレージ 1112 は、メインメモリーなどの（少なくとも）揮発性ストレージ、ならびにフラッシュストレージ、および他のタイプのディスクドライブなどを含み得る「二次ストレージ」などのより安定したストレージ（例えば、より不揮発性のストレージ）を含め、任意のタイプのページデータストレージを含み得る。本明細書の説明の趣旨を逸脱することなく、本明細書で説明される技法と一緒に使用され得る多くのタイプのページデータストレージが存在することが、データ処理の分野の業者には理解されよう。

【0156】

[0168] 例示的な実施形態によれば、データマネージャー 1106、またはデータマネージャー 1106 の 1 つまたは複数の部分が、後段でさらに説明されるとおり、有形のコンピューター可読記憶媒体上に記憶され得る実行可能命令を含み得る。例示的な実施形態によれば、コンピューター可読記憶媒体は、分散デバイスを含め、任意の数のストレージデバイス、および任意の数の記憶媒体タイプを含み得る。

【0157】

[0169] この脈絡において、「プロセッサ」は、コンピューティングシステムに関連する命令を処理するように構成された単一のプロセッサ、または複数のプロセッサを含み得る。このため、プロセッサは、並行に、かつ / または分散された状態で命令を処理する 1 つまたは複数のプロセッサを含み得る。デバイスプロセッサ 1104 は、図 11 のデータマネージャー 1106 の外部に描かれるものの、デバイスプロセッサ 1104 は、データマネージャー 1106 および / またはデータマネージャー 1106 の要素のいずれかの内部または外部に配置され得る単一のコンポーネントとして、かつ / または分散ユニットとして実装され得ることが、データ処理の分野の業者には理解されよう。

【 0 1 5 8 】

[0170]例えば、システム 1 1 0 0 は、1つまたは複数のプロセッサ 1 1 0 4 を含む得る。例えば、システム 1 1 0 0 は、1つまたは複数の 1 1 0 4 によって実行可能な命令を記憶する少なくとも1つの有形のコンピューター可読記憶媒体を含むことが可能であり、それらの実行可能な命令は、少なくとも1つのデータ処理装置に、本明細書で説明されるとおり、システム 1 1 0 0 に含まれる様々な例示的な構成要素に関連付けられた操作を実行させるように構成される。例えば、その1つまたは複数のプロセッサ 1 1 0 4 は、少なくとも1つのデータ処理装置の中に含められ得る。本明細書の説明の趣旨を逸脱することなく、本明細書の説明により構成され得るプロセッサおよびデータ処理装置の多くの構成が存在することが、データ処理の分野の業者には理解されよう。

10

【 0 1 5 9 】

[0171]この脈絡において、「コンポーネント」とは、ある操作を実行するように構成され得る命令またはハードウェアを指すことが可能である。そのような命令は、命令のコンポーネントグループ内に含まれることが可能であり、または複数のグループにわたって分散されることが可能である。例えば、第1のコンポーネントの操作に関連するいくつかの命令が、第2のコンポーネント（またはより多くのコンポーネント）の操作に関連する命令のグループに含められてもよい。例えば、本明細書の「コンポーネント」は、単一のエンティティ内に配置され得る命令によって実施され得るある種の機能を指すことが可能であり、または複数のエンティティにわたって散在させられる、もしくは分散させられることが可能であり、他のコンポーネントに関連する命令および/またはハードウェアと重なり合うことが可能である。

20

【 0 1 6 0 】

[0172]例示的な実施形態によれば、データマネージャー 1 1 0 6 は、1つまたは複数のユーザーデバイスに関連して実施され得る。例えば、データマネージャー 1 1 0 6 は、後段でさらに説明されるとおり、サーバーと通信することが可能である。

【 0 1 6 1 】

[0173]例えば、1つまたは複数のデータベースが、データベースインターフェースコンポーネント 1 1 2 2 を介してアクセスされ得る。様々なタイプのデータベース構成（例えば、リレーショナルデータベース、階層データベース、分散データベース）および非データベース構成などの、本明細書で説明される情報を記憶するための多くの技法が存在することが、データ処理の分野の業者には理解されよう。

30

【 0 1 6 2 】

[0174]例示的な実施形態によれば、データマネージャー 1 1 0 6 は、即時の結果などのオブジェクトを記憶することが可能なメモリー 1 1 2 4 を含む得る。この脈絡において、「メモリー」は、データおよび/または命令を記憶するように構成された単一のメモリーデバイスまたは複数のメモリーデバイスを含む得る。さらに、メモリー 1 1 2 4 は、複数の分散ストレージデバイスにわたることが可能である。さらに、メモリー 1 1 2 4 は、複数のプロセッサの間に分散され得る。

【 0 1 6 3 】

[0175]例示的な実施形態によれば、ユーザーインターフェースコンポーネント 1 1 2 6 が、ユーザー 1 1 2 8 とデータマネージャー 1 1 0 6 の間の通信を管理することが可能である。ユーザー 1 1 2 8 は、ディスプレイ 1 1 3 2 および他の入出力デバイスに関連付けられ得る受信デバイス 1 1 3 0 に関連付けられ得る。例えば、ディスプレイ 1 1 3 2 は、内部デバイスバス通信を介して、または少なくとも1つのネットワーク接続を介して受信デバイス 1 1 3 0 と通信するように構成され得る。

40

【 0 1 6 4 】

[0176]例示的な実施形態によれば、ディスプレイ 1 1 3 2 は、フラットスクリーンディスプレイ、印刷形態のディスプレイ、2次元ディスプレイ、3次元ディスプレイ、静的ディスプレイ、移動ディスプレイ、触覚出力、オーディオ出力、およびユーザー（例えば、ユーザー 1 1 2 8）とコミュニケーションをとるための他の任意の形態の出力などの知覚

50

ディスプレイとして実施され得る。

【 0 1 6 5 】

[0177]例示的な実施形態によれば、データマネージャー 1 1 0 6 は、データマネージャー 1 1 0 6 と、少なくとも 1 つのネットワーク 1 1 3 6 を介してデータマネージャー 1 1 0 6 と通信することが可能な他のエンティティとの間のネットワーク通信を管理することが可能なネットワーク通信コンポーネント 1 1 3 4 を含む得る。例えば、ネットワーク 1 1 3 6 は、インターネット、少なくとも 1 つのワイヤレスネットワーク、または少なくとも 1 つの有線ネットワークのうちの少なくとも 1 つを含む得る。例えば、ネットワーク 1 1 3 6 は、データマネージャー 1 1 0 6 のためにデータの伝送をサポートすることが可能なセルラーネットワーク、無線ネットワーク、または任意のタイプのネットワークを含む得る。例えば、ネットワーク通信コンポーネント 1 1 3 4 は、データマネージャー 1 1 0 6 と受信デバイス 1 1 3 0 の間のネットワーク通信を管理することが可能である。例えば、ネットワーク通信コンポーネント 1 1 3 4 は、ユーザーインターフェースコンポーネント 1 1 2 6 と受信デバイス 1 1 3 0 の間のネットワーク通信を管理することが可能である。

10

【 0 1 6 6 】

[0178]例えば、データ不透明型インターフェース 1 1 0 8 が、任意に選択されたページ指向型アクセスメソッド 1 1 1 0 に、ページデータストレージ 1 1 1 2 に対するログ構造化アクセスを含むページデータストレージ 1 1 1 2 に対するインターフェースアクセスをもたらすように構成され得る。

20

【 0 1 6 7 】

[0179]例えば、キャッシュ層マネージャー 1 1 3 8 が、データ不透明型インターフェース 1 1 0 8 に関連付けられた間接アドレスマッピングテーブル 1 1 4 2 に対するテーブル操作を開始するように構成され得るマッピングテーブルマネージャー 1 1 4 0 を含むことが可能であり、テーブル操作は、間接アドレスマッピングテーブル 1 1 4 2 の中のエントリーに対してアトミック・コンペア・アンド・スワップ (C A S) 操作を開始して、ページデータストレージ 1 1 1 2 に関連付けられたページのそれまでの状態を、それらのページの新たな状態で置き換えることを含む。

【 0 1 6 8 】

[0180]例えば、マッピングテーブルマネージャー 1 1 4 0 は、データ不透明型インターフェース 1 1 0 8 に関連付けられた間接アドレスマッピングテーブル 1 1 4 2 に対するテーブル操作を開始するように構成されることが可能であり、間接アドレスマッピングテーブル 1 1 4 2 は、キャッシュ層ストレージ 1 1 4 4 と、二次ストレージ 1 1 4 6 とを含むデータストレージの管理のために共通で使用される。

30

【 0 1 6 9 】

[0181]例えば、間接アドレスマッピングテーブル 1 1 4 2 は、ページの論理ロケーションを、それらのページの対応する物理ロケーションから分離し、ページデータストレージのユーザーは、それらのページの物理ロケーションアドレス値の代わりにページ識別子値を、ページデータストレージを参照するデータ構造における別の場所に記憶させる。

【 0 1 7 0 】

[0182]例えば、アップデートマネージャー 1 1 4 8 が、間接アドレスマッピングテーブル 1 1 4 2 の中のエントリーに対してラッチフリーのコンペア・アンド・スワップ操作を使用して、間接アドレスマッピングテーブル 1 1 4 2 上でアトミック状態変更を行うようデータアップデートおよび管理アップデートを制御するように構成され得る。

40

【 0 1 7 1 】

[0183]例えば、ストレージ層 1 1 4 9 が、間接アドレスマッピングテーブル 1 1 4 2 の中のエントリーに対してラッチフリーのコンペア・アンド・スワップ操作を使用して、ページフラッシュからもたらされるログ構造化に関連するページロケーション変更を制御するように構成され得るログ構造化ストレージ層マネージャー 1 1 5 0 を含む得る。

【 0 1 7 2 】

50

[0184]例えば、バッファーマネージャ１１５１が、ラッチフリーのアップデート操作を介してログ構造化二次ストレージバッファに対するアップデートを制御するように構成され得る。このため、例えば、複数のスレッドが、ラッチフリーの操作を介してログ構造化二次ストレージバッファを同時にアップデートすることが可能である。

【０１７３】

[0185]例えば、バッファーマネージャ１１５１は、第１の二次ストレージアドレス引数までの、下位のアドレスを有する、ログ構造化二次ストレージバッファにフラッシュされたページが、ログ構造化二次ストレージ内で安定していることを判定するための安定操作を開始するように構成され得る。

【０１７４】

10

[0186]例えば、ページマネージャ１１５２が、ページに対するフラッシュ操作、割当て操作、および解放操作を制御するように構成され得る。例えば、ページマネージャ１１５２は、コンペア・アンド・スワップ（ＣＡＳ）操作に基づいて、二次ストレージバッファ内に第１のページのページ状態をコピーすることを開始すること、そのページ状態の先頭にフラッシュデルタレコードを付加することであって、そのフラッシュデルタレコードは、二次ストレージ内の第１のページのストレージロケーションと、呼び出し元に関連付けられた注釈とを示す二次ストレージアドレスを含むこと、およびマッピングテーブルの中にフラッシュデルタレコードのアドレスをインストールすることに基づいてページ状態に対するアップデートを開始することに基づいて、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第１のページをフラッシュする操作を開始するように構成され得る。

20

【０１７５】

[0187]例えば、ページマネージャ１１５２は、第１のページに関連付けられたページ状態の先頭に部分的スワップデルタレコードを付加することであって、その部分的スワップデルタレコードは、第１のページの欠落した部分の二次ストレージ内のロケーションを示すフラッシュデルタレコードのストレージロケーションを示すメインメモリアドレスを含むことに基づいて、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第１のページの部分をスワップする操作を開始するように構成され得る。

【０１７６】

[0188]例えば、システムトランザクションマネージャ１１５４が、トランザクションをコミットすること、およびトランザクションを中途終了することを行うように構成され得る。

30

【０１７７】

[0189]例えば、レコードマネージャ１１５６が、アップデートデルタレコード操作および置換アップデート操作に基づいてアップデートを制御するように構成され得る。

[0190]例えば、エポックマネージャ１１６０が、第１のプロセッサ操作によってページ情報にアクセスすることに先立って、第１のエポックに関連付けられた第１のエポック登録リスト内に第１のプロセッサ操作を登録することを開始するように構成され得る。例えば、第１のプロセッサ操作は、スレッドであり得る。

【０１７８】

40

[0191]例えば、ページマネージャ１１５２は、コンペア・アンド・スワップ（ＣＡＳ）操作を介して、マッピングテーブルの中にフラッシュデルタレコードに対するポインターをインストールすることであって、そのフラッシュデルタレコードは、ＣＡＳ操作を介してマッピングテーブル内で置き換えられる既存のページ状態の先頭に付加されることに基づいて、二次ストレージにページ状態をフラッシュするように構成され得る。

【０１７９】

[0192]例えば、ページマネージャ１１５２は、そのＣＡＳ操作が成功したかどうかを判定し、そのＣＡＳ操作が成功したと判定された場合、二次ストレージフラッシュバッファに既存のページ状態を書き込む書込み操作を開始するように構成され得る。

【０１８０】

50

[0193]例えば、ページマネージャー 1 1 5 2 は、その C A S 操作が失敗したと判定された場合、既存のページにそれまでに割り当てられていたストレージスペースを無効にする操作を開始するように構成され得る。

【 0 1 8 1 】

[0194]本明細書の説明の趣旨を逸脱することなく、多くの異なる技法が、ラッチフリーのログ構造化ストレージシステムのために使用され得ることが、データ処理の分野の業者には理解されよう。

【 0 1 8 2 】

[0195] I I I . フローチャート説明

本明細書で説明される特徴は、本明細書の説明の趣旨を逸脱することなく、データ処理の分野の業者によって理解され得る多くの異なる状態で実施され得る例示的な実施形態として提供される。そのような特徴は、例示的な実施形態特徴としてのみ解釈されるべきであり、それらの詳細な説明だけに限定されるものと解釈されることは意図していない。

【 0 1 8 3 】

[0196]図 1 2 a ~ 図 1 2 d は、例示的な実施形態による、図 1 1 のシステムの例示的な動作を示すフローチャートである。図 1 2 a の例において、ページデータストレージに対するラッチフリーのアクセスを含むページデータストレージに対するインターフェースアクセスが、任意に選択されたページ指向型アクセスメソッドにもたらされ得る (1 2 0 2)。例えば、前述したとおり、任意に選択されたページ指向型アクセスメソッド 1 1 1 0 に、ページデータストレージ 1 1 1 2 に対するラッチフリーのアクセスを含むページデータストレージ 1 1 1 2 に対するインターフェースアクセスをもたらすことが可能なデータ不透明型インターフェース 1 1 0 8。

【 0 1 8 4 】

[0197]例えば、ページデータストレージに対するインターフェースアクセスは、安定したページデータストレージに対するログ構造化アクセスを含み得る (1 2 0 4)。例えば、データ不透明型インターフェース 1 1 0 8 が、前述したとおり、任意に選択されたページ指向型アクセスメソッド 1 1 1 0 に、ページデータストレージ 1 1 1 2 に対するログ構造化アクセスを含むページデータストレージ 1 1 1 2 に対するインターフェースアクセスをもたらすことが可能である。

【 0 1 8 5 】

[0198]例えば、データ不透明型インターフェースに関連付けられた間接アドレスマッピングテーブルに対するテーブル操作が開始されることが可能であり、それらのテーブル操作は、間接アドレスマッピングテーブル内のエントリーに対してアトミック・コンペア・アンド・スワップ操作を開始して、ページデータストレージに関連付けられたページのそれまでの状態を、それらのページの新たな状態で置き換えることを含む (1 2 0 6)。例えば、マップテーブルマネージャー 1 1 4 0 が、前述したとおり、データ不透明型インターフェース 1 1 0 8 に関連付けられた間接アドレスマッピングテーブル 1 1 4 2 に対してテーブル操作を開始することが可能であり、テーブル操作は、間接アドレスマッピングテーブル 1 1 4 2 内のエントリーに対してコンペア・アンド・スワップ (C A S) 操作を開始して、ページデータストレージ 1 1 1 2 に関連するページのそれまでの状態を、それらのページの新たな状態で置き換えることを含む。

【 0 1 8 6 】

[0199]例えば、間接アドレスマッピングテーブルが、図 1 2 b に示されるとおり、キャッシュ層ストレージと、二次ストレージとを含むデータストレージの管理のために共通で使用され得る (1 2 0 8)。例えば、マップテーブルマネージャー 1 1 4 0 が、前述したとおり、データ不透明型インターフェース 1 1 0 8 に関連付けられた間接アドレスマッピングテーブル 1 1 4 2 に対するテーブル操作を開始することが可能であり、間接アドレスマッピングテーブル 1 1 4 2 は、キャッシュ層ストレージ 1 1 4 4 と、二次ストレージ 1 1 4 6 とを含むデータストレージの管理のために共通で使用される。

【 0 1 8 7 】

[0200]例えば、ページの論理ロケーションが、前述したとおり、それらのページの対応する物理ロケーションから分離されることが可能であり、ページデータストレージのユーザーは、それらのページに関する物理ロケーションアドレスの代わりにページ識別子値を、ページデータストレージを参照するデータ構造における別の場所に記憶させる（1210）。例えば、間接アドレスマッピングテーブル1142が、前述したとおり、ページの論理ロケーションを、それらのページの対応する物理ロケーションから分離し、ページデータストレージのユーザーは、それらのページに関する物理ロケーションアドレスの代わりにページ識別子値を、ページデータストレージを参照するデータ構造における別の場所に記憶させる。

【0188】

10

[0201]例えば、データアップデートおよび管理アップデートが、間接アドレスマッピングテーブル内のエントリーに対してラッチフリーのコンペア・アンド・スワップ操作を使用して、間接アドレスマッピングテーブル上でアトミック状態変更を行って、制御され得る（1212）。例えば、アップデートマネージャー1148が、前述したとおり、間接アドレスマッピングテーブル1142内のエントリーに対してラッチフリーのコンペア・アンド・スワップ操作を使用して、間接アドレスマッピングテーブル1142上でアトミック状態変更を行って、データアップデートおよび管理アップデートを制御することが可能である。

【0189】

[0202]例えば、ページフラッシュからもたらされるログ構造化に関連するページロケーション変更が、間接アドレスマッピングテーブル内のエントリーに対してラッチフリーのコンペア・アンド・スワップ操作を使用して、制御され得る（1214）。例えば、ログ構造化ストレージ層マネージャー1150が、前述したとおり、間接アドレスマッピングテーブル1142内のエントリーに対してラッチフリーのコンペア・アンド・スワップ操作を使用して、ページフラッシュからもたらされるログ構造化に関連するページロケーション変更を制御することが可能である。

20

【0190】

[0203]例えば、図12cの例において、第1のエポックに関連付けられた第1のエポック登録リストに第1のプロセッサ操作を登録することが、第1のプロセッサ操作によってページ情報にアクセスすることに先立って開始され得る（1216）。

30

【0191】

[0204]例えば、ページ状態が、コンペア・アンド・スワップ（CAS）操作を介して、マッピングテーブルの中にフラッシュデルタレコードに対するポインターをインストールすることに基づいて、二次ストレージにフラッシュされることが可能であり、そのフラッシュデルタレコードは、CAS操作を介してマッピングテーブル内で置き換えられる既存のページ状態の先頭に付加される（1218）。

【0192】

[0205]例えば、ログ構造化二次ストレージバッファーに対するアップデートが、ラッチフリーのアップデート操作を介して制御され得る（1220）。

[0206]例えば、図12dの例において、コンペア・アンド・スワップ（CAS）操作を介して、二次ストレージ内のロケーションにキャッシュ層ストレージ内の第1のページをフラッシュする操作が、第2のストレージバッファーの中に第1のページのページ状態をコピーすることを開始すること、そのページ状態の先頭にフラッシュデルタレコードを付加することを開始することであって、そのフラッシュデルタレコードは、二次ストレージ内の第1のページのストレージロケーションと、呼び出し元に関連付けられた注釈とを示す二次ストレージアドレスを含むこと、およびマッピングテーブルの中にフラッシュデルタレコードのアドレスをインストールすることに基づいて、そのページ状態に対するアップデートを開始することに基づいて、開始され得る（1222）。

40

【0193】

[0207]例えば、二次ストレージ内のロケーションに対してキャッシュ層ストレージ内の

50

第 1 のページの一部分をスワップする操作が、第 1 のページに関連付けられたページ状態の先頭に部分的スワップデルタレコードを付加することを開始することによって、その部分的スワップデルタレコードは、第 1 のページの欠落した部分の二次ストレージ内のロケーションを示すフラッシュデルタレコードのストレージロケーションを示すメインメモリアドレスを含むことに基づいて、開始され得る (1 2 2 4)。

【 0 1 9 4 】

[0208] 本明細書の説明の趣旨を逸脱することなく、ラッチフリーのログ構造化ストレージシステムのために多くの異なる技法が使用され得ることが、データ処理の分野の業者には理解されよう。

【 0 1 9 5 】

[0209] 顧客プライバシーおよび顧客機密保持が、長年にわたるデータ処理環境における継続的な配慮事項である。このため、ラッチフリーのログ構造化ストレージシステムのための例示的な技法は、そのような分析に関連付けられた、関連付けられたアプリケーションまたはサービスを有する 1 つまたは複数の契約合意 (例えば、「サービス利用規約」TOS 合意) を介して許可を与えたユーザーによって提供されるユーザー入力および / またはユーザーデータを使用することが可能である。例えば、ユーザーが、ユーザーの入力データが送信され、デバイス上に記憶されることに同意を与えることが可能であり、ただし、(例えば、ユーザーが受諾した合意を介して) 各関係者が、送信および / または記憶がどのように行われるか、および記憶が保持される場合、どのようなレベルもしくは持続時間の記憶が保持され得るかを管理することが可能である。

【 0 1 9 6 】

[0210] 本明細書で説明される様々な技法の実施形態は、デジタル電子回路として、またはコンピューターハードウェア、ファームウェア、ソフトウェア、またはそれらの組合せ (例えば、様々な機能を実行する命令を実行するように構成された装置) として実施され得る。

【 0 1 9 7 】

[0211] 実施形態は、純粋な伝播される信号などの純粋な信号として実現されるコンピュータープログラムとして実施され得る。そのような実施形態は、「コンピューター可読伝送媒体」を介して実施されるものとして本明細書に記載され得る。

【 0 1 9 8 】

[0212] 代替として、実施形態は、データ処理装置、例えば、プログラマブルプロセッサ、コンピューター、または複数のコンピューターによって実行されるように、またはデータ処理装置の動作を制御するようにマシン使用可能デバイスもしくはマシン可読ストレージデバイス (例えば、ユニバーサルシリアルバス (USB) ストレージデバイス、テープ、ハードディスクドライブ、コンパクトディスク、デジタルビデオディスク (DVD) などの磁気媒体もしくはデジタル媒体) として実現されたコンピュータープログラムとして実施され得る。そのような実施形態は、「コンピューター可読記憶媒体」または「コンピューター可読ストレージデバイス」を介して実施されるものとして本明細書に記載されることが可能であり、このため、純粋な伝播される信号などの純粋に信号である実施形態と異なる。

【 0 1 9 9 】

[0213] 前述したコンピュータープログラムなどのコンピュータープログラムは、コンパイラー型言語、インタープリター型言語、または機械語を含め、任意の形態のプログラミング言語で書かれることが可能であり、スタンドアロンのプログラムとして、またはモジュール、コンポーネント、サブルーチンとして、あるいはコンピューティング環境において使用されるのに適した他のユニットとして展開されることを含め、任意の形態で展開され得る。コンピュータープログラムは、マシン使用可能デバイス上、またはマシン可読ストレージデバイス (例えば、コンピューター可読媒体) 上の実行可能コード (例えば、実行可能命令) として実体化され得る。前述した技法を実施することが可能なコンピュータープログラムは、1 つの場所における、または複数の場所にわたって分散され、通信ネッ

10

20

30

40

50

トワークによって互いに接続された1つのコンピューターまたは複数のコンピューターの上で実行されるように展開され得る。

【0200】

[0214]方法ステップが、入力データを操作すること、および出力を生成することによって機能を実行するようコンピュータープログラムを実行する1つまたは複数のプログラマブルプロセッサによって実行され得る。1つまたは複数のプログラマブルプロセッサは、命令を並行に実行することが可能であり、かつ/または分散処理のために分散構成で配置されることが可能である。また、本明細書で説明される例示的な機能は、ハードウェア論理コンポーネントによって実行されることも可能であり、装置は、少なくとも部分的に、1つまたは複数のハードウェア論理コンポーネントとして実施されることが可能である。例えば、限定なしに、使用され得るハードウェア論理コンポーネントの例示的なタイプには、フィールド・プログラマブル・ゲート・アレイ(FPGA)、Program-specific Integrated Circuit(ASIC)、Program-specific Standard Product(ASSP)、System-on-a-chipシステム(SOC)、Complex Programmable Logic Device(CPLD)などが含まれ得る。

10

【0201】

[0215]コンピュータープログラムの実行に適したプロセッサには、例として、汎用マイクロプロセッサおよび専用マイクロプロセッサ、ならびに任意の種類のデジタルコンピューターの任意の1つまたは複数のプロセッサが含まれる。一般に、プロセッサは、読取り専用メモリから、またはランダムアクセスメモリから、または読取り専用メモリとランダムアクセスメモリの両方から命令およびデータを受け取る。コンピューターの要素には、命令を実行するための少なくとも1つのプロセッサ、ならびに命令およびデータを記憶するための1つまたは複数のメモリーデバイスが含まれ得る。一般に、コンピューターは、データを記憶するための1つまたは複数の大容量ストレージデバイス、例えば、磁気ディスク、光磁気ディスク、または光ディスクを含むことも可能であり、あるいはそのような大容量ストレージデバイスからデータを受け取るように、もしくはそのような大容量ストレージデバイスにデータを送るように、またはその両方を行うように動作上、結合されることが可能である。コンピュータープログラム命令およびデータを実体化するのに適した情報キャリアには、例として、半導体メモリーデバイス、例えば、EPROM、EEPROM、およびフラッシュメモリーデバイス、磁気ディスク、例えば、内部ハードディスクもしくはリムーバブルディスク、光磁気ディスク、ならびにCD-ROMディスクおよびDVD-ROMディスクを含め、すべての形態の不揮発性メモリーが含まれる。プロセッサおよびメモリーは、専用論理回路によって捕捉されること、または専用論理回路に組み込まれることが可能である。

20

30

【0202】

[0216]ユーザーとの対話を可能にするのに、実施形態は、ユーザーに情報を表示するためのディスプレイデバイス、例えば、陰極線管(CRT)、液晶ディスプレイ(LCD)、またはプラズマモニターと、ユーザーがコンピューターに入力を与えることができるキーボードおよびポインティングデバイス、例えば、マウスもしくはトラックボールとを有するコンピューター上で実施され得る。他の種類のデバイスが、ユーザーとの対話を可能にするのに使用されることが可能であり、例えば、ユーザーに与えられるフィードバックは、任意の形態の知覚フィードバック、例えば、視覚フィードバック、聴覚フィードバック、または触覚フィードバックであり得る。例えば、視覚出力(例えば、視覚的ジェスチャー、ビデオ出力)、オーディオ出力(例えば、音声、デバイスサウンド)、触覚出力(例えば、タッチ、デバイスの動き)、温度、匂いなどを含め(ただし、以上には限定されない)、任意の形態の知覚出力を介して与えられ得る。

40

【0203】

[0217]さらに、ユーザーからの入力、音響入力、音声入力、または触覚入力を含め、任意の形態で受け取られることが可能である。例えば、入力は、視覚入力(例えば、ジェ

50

スチャー、ビデオ入力)、オーディオ入力(例えば、音声、デバイスサウンド)、触覚入力(例えば、タッチ、デバイスの動き)、温度、匂いなどを含め(ただし、以上には限定されない)、任意の形態の知覚入力を介してユーザーから受け取られ得る。

【0204】

[0218]さらに、ナチュラルユーザーインターフェース(NUI)が、ユーザーとインターフェースをとるのに使用され得る。この脈絡において、「NUI」とは、ユーザーが、マウス、キーボード、リモコンなどの入力デバイスによって課せられる人為的な制約を免れて、「自然な」状態でデバイスと対話することを可能にする任意のインターフェース技術を指すことが可能である。

【0205】

[0219]NUI技法の例には、音声認識、タッチ認識およびスタイラス認識、画面上のジェスチャー認識と画面に隣接したジェスチャー認識の両方、エアジェスチャー、頭部および目追跡、音声および発話、視覚、タッチ、ジェスチャー、ならびにマシンインテリジェンスに依拠する技法が含まれ得る。例示的なNUI技術には、すべてよりナチュラルなインターフェースを提供することが可能な、タッチセンシティブディスプレイ、音声および発話認識、意図および目標理解、デプスカメラ(例えば、立体カメラシステム、赤外線カメラシステム、RGB(赤、緑、青)カメラシステム、および以上の組合せ)を使用するモーションジェスチャー検出、加速度計/ジャイロ스코ープを使用するモーションジェスチャー検出、顔認識、3Dディスプレイ、頭部、目、および視線追跡、没入型拡張現実システムおよび仮想現実システム、ならびに電界検出電極を用いた脳活動を検知するための技術(例えば、脳波計(EEG)技法および関連する技法)が含まれ得るが、以上には限定されない。

【0206】

[0220]実施形態は、例えば、データサーバーとしてバックエンドコンポーネントを含む、または、ミドルウェアコンポーネント、例えば、アプリケーションサーバーを含む、またはフロントエンドコンポーネント、例えば、ユーザーが実施形態と対話することができるグラフィカルユーザーインターフェースまたはウェブブラウザを有するクライアントコンピュータを含む、あるいはそのようなバックエンドコンポーネント、ミドルウェアコンポーネント、またはフロントエンドコンポーネントの任意の組合せを含むコンピューティングシステムとして実施され得る。コンポーネントは、任意の形態または媒体のデジタルデータ通信、例えば、通信ネットワークによって互いに接続され得る。通信ネットワークの例には、ローカルエリアネットワーク(LAN)およびワイドエリアネットワーク(WAN)、例えば、インターネットが含まれる。

【0207】

[0221]主題は、構造上の特徴および/または方法上の動作に特有の言い回しで説明されてきたが、添付の特許請求の範囲において規定される主題は、以上に説明される特定の特徴または動作に必ずしも限定されないことを理解されたい。むしろ、以上に説明される特定の特徴および動作は、特許請求の範囲を実施する例示的な形態として開示される。説明される実施形態のいくつかの特徴が、本明細書で記載されるとおり例示されてきたが、今や、多くの変形形態、置換形態、変更形態、および均等形態が当業者には想起されよう。したがって、添付の特許請求の範囲は、実施形態の範囲に含まれるすべてのそのような変形形態および変更形態に及ぶものとする。

【図 1】

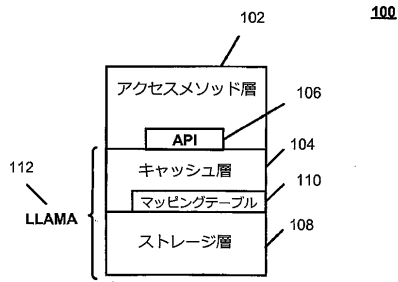


FIG. 1

【図 2】

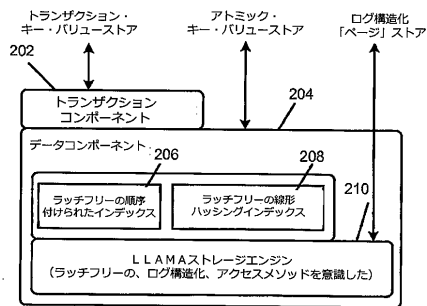


FIG. 2

【図 3】

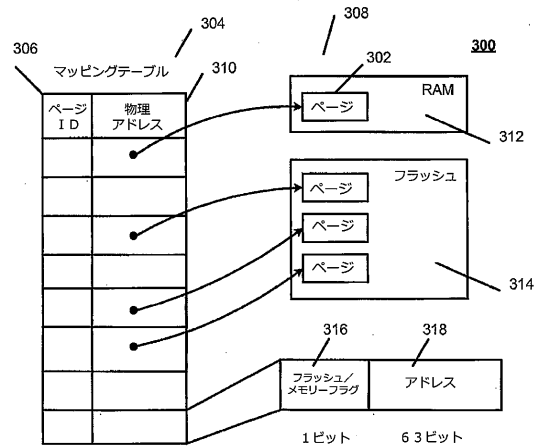


FIG. 3

【図 4 a】

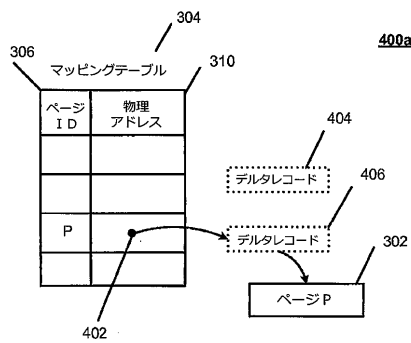


FIG. 4a

【図 4 b】

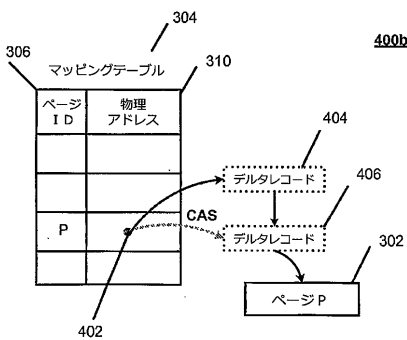


FIG. 4b

【図 5】

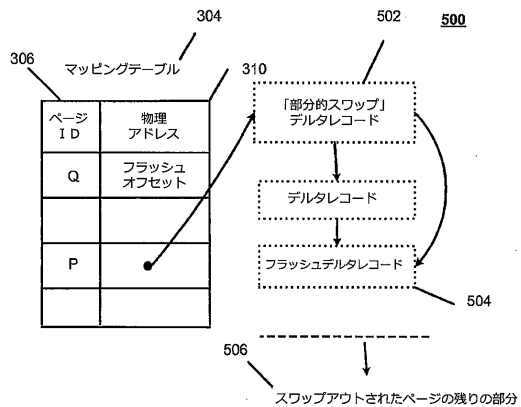


FIG. 5

【 図 7 a 】

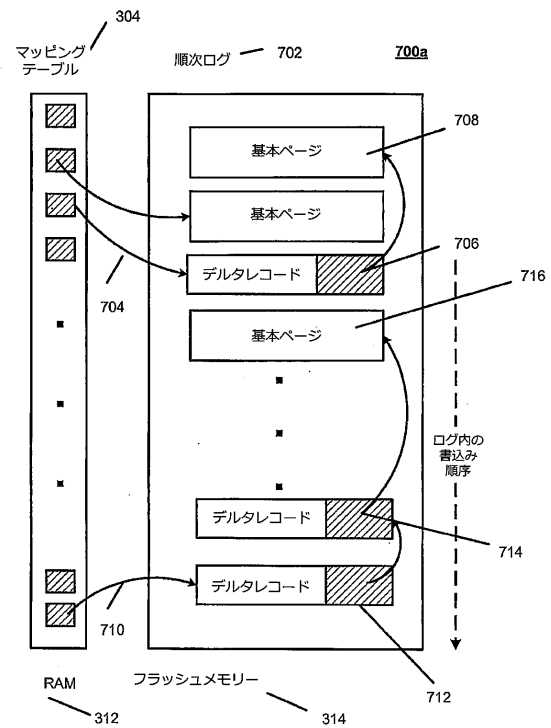


FIG. 7a

【 図 7 c 】

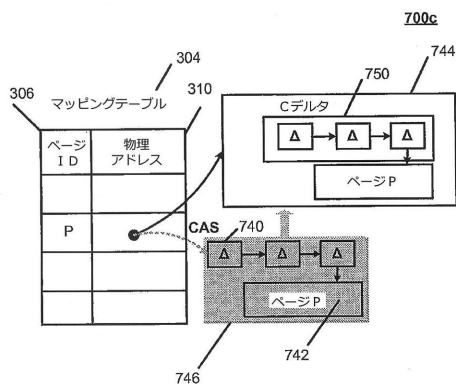


FIG. 7c

【図 8】

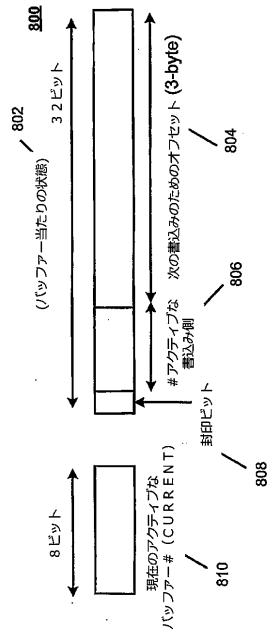


FIG. 8

【図 9】

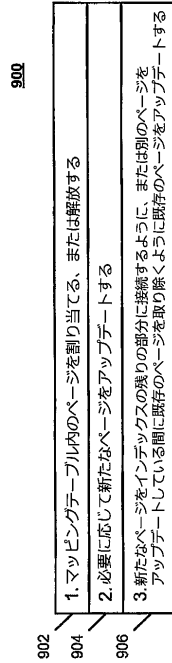


FIG. 9

【図 10】

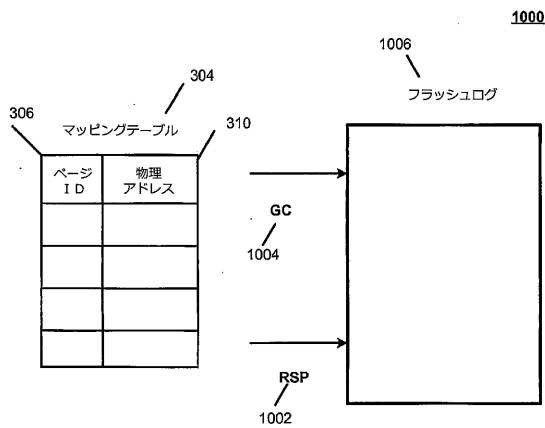


FIG. 10

【図 11】

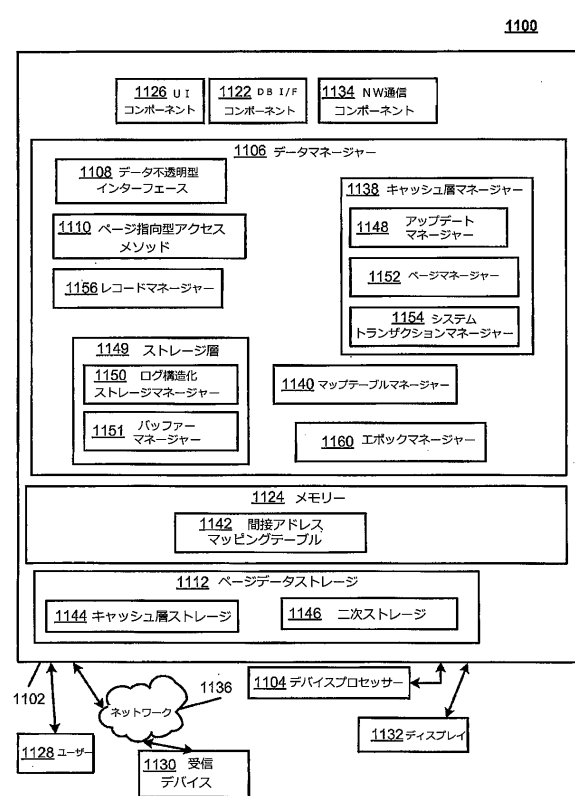


FIG. 11

【図 12 a】

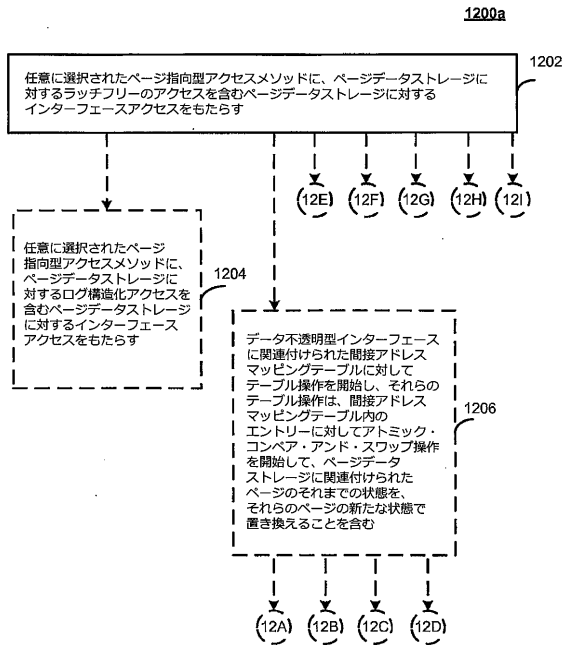


FIG. 12a

【図 12 b】

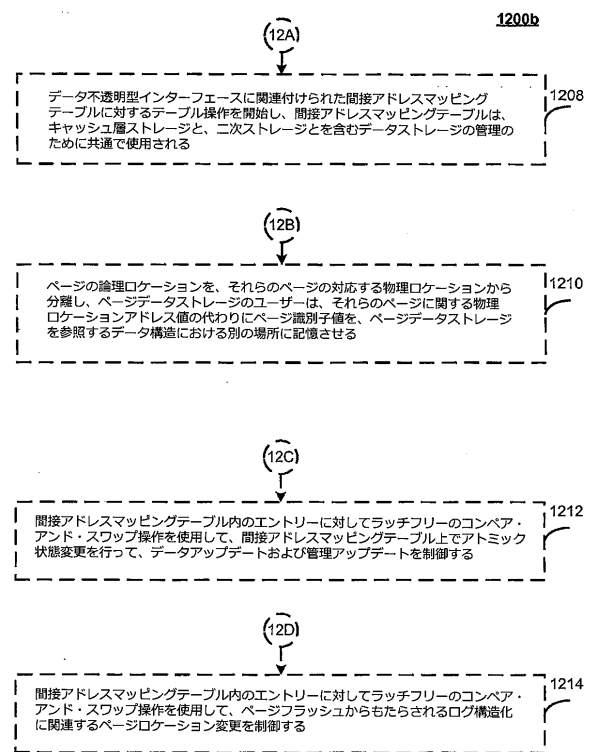


FIG. 12b

【図 12 c】

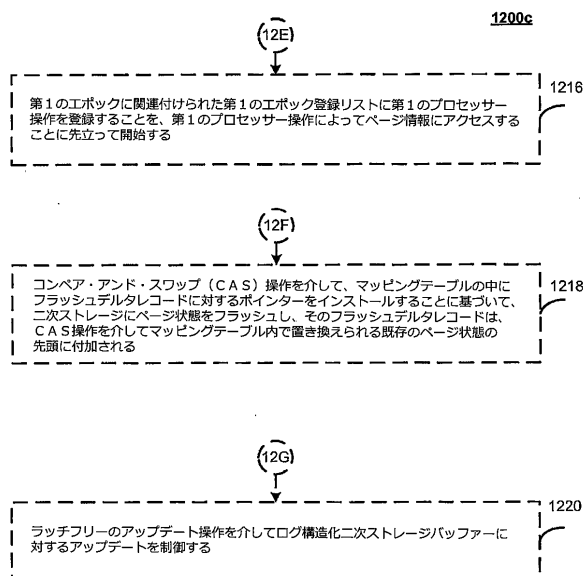


FIG. 12c

【図 12 d】

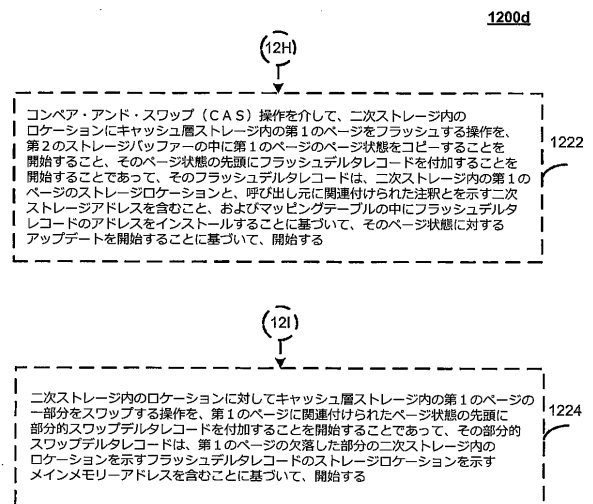


FIG. 12d

フロントページの続き

(51)Int.Cl. F I
G 0 6 F 12/02 5 2 0 A

(74)代理人 100162846

弁理士 大牧 綾子

(72)発明者 ロメット,デーヴィッド・ビー

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9 , レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

(72)発明者 レヴァンドスキ, ジャスティン

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9 , レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

(72)発明者 セングプタ, スディプタ

アメリカ合衆国ワシントン州 9 8 0 5 2 - 6 3 9 9 , レッドモンド, マイクロソフト コーポレーション, エルシーエイ - インターナショナル・パテンツ (8 / 1 1 7 2)

審査官 田中 啓介

(56)参考文献 特表 2 0 0 2 - 5 2 4 8 0 1 (J P , A)

特開 2 0 1 0 - 2 1 8 5 2 9 (J P , A)

米国特許出願公開第 2 0 0 8 / 0 0 6 5 6 7 0 (U S , A 1)

米国特許出願公開第 2 0 0 9 / 0 2 4 0 6 6 4 (U S , A 1)

米国特許出願公開第 2 0 1 0 / 0 1 9 1 7 1 3 (U S , A 1)

特開 2 0 0 7 - 0 1 2 0 6 1 (J P , A)

JUSTIN J LEVANDOSKI; DAVID B LOMET; SUDIPTA SENGUPTA, THE BW-TREE: A B-TREE FOR NEW HARDWARE PLATFORMS, DATA ENGINEERING (ICDE), 2013 IEEE 29TH INTERNATIONAL CONFERENCE ON, [ONLINE], IEEE, 2 0 1 3 年 4 月 8 日, pp.302 - 313, U R L , <http://dx.doi.org/10.1109/ICDE.2013.6544834>

(58)調査した分野(Int.Cl., D B 名)

G 0 6 F 3 / 0 6 - 3 / 0 8

G 0 6 F 1 2 / 0 0 - 1 2 / 1 2 8

G 0 6 F 1 3 / 0 0 - 1 3 / 1 8、1 7 / 3 0