

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5301670号
(P5301670)

(45) 発行日 平成25年9月25日(2013.9.25)

(24) 登録日 平成25年6月28日(2013.6.28)

(51) Int.Cl.

F I

G 0 6 F 3/06 (2006.01)

G 0 6 F 3/06 3 0 4 P

G 0 6 F 3/06 3 0 1 J

請求項の数 19 (全 20 頁)

(21) 出願番号	特願2011-527865 (P2011-527865)	(73) 特許権者	500046438
(86) (22) 出願日	平成21年8月27日 (2009.8.27)		マイクロソフト コーポレーション
(65) 公表番号	特表2012-503250 (P2012-503250A)		アメリカ合衆国 ワシントン州 9805
(43) 公表日	平成24年2月2日 (2012.2.2)		2-6399 レッドモンド ワン マイ
(86) 国際出願番号	PCT/US2009/055198		クロソフト ウェイ
(87) 国際公開番号	W02010/033365	(74) 代理人	100140109
(87) 国際公開日	平成22年3月25日 (2010.3.25)		弁理士 小野 新次郎
審査請求日	平成24年7月13日 (2012.7.13)	(74) 代理人	100075270
(31) 優先権主張番号	12/234,411		弁理士 小林 泰
(32) 優先日	平成20年9月19日 (2008.9.19)	(74) 代理人	100080137
(33) 優先権主張国	米国 (US)		弁理士 千葉 昭男
		(74) 代理人	100096013
			弁理士 富田 博行
		(74) 代理人	100147991
			弁理士 鳥居 健一

最終頁に続く

(54) 【発明の名称】 データ・ストアへの書き込みトラフィックの統合

(57) 【特許請求の範囲】

【請求項 1】

ランダムアクセス媒体上で書き込みログにデータセットを連続的に統合する、機械により実施される方法であって、

前記データセットが前記書き込みログにおいて連続的に統合することに適しているか否かを決定するステップと、

前記データセットが連続的に統合することに適していると決定される場合に、前記ランダムアクセス媒体上の意図した宛先へ書き込む代わりに、前記データセットを前記書き込みログに連続的に書き込むステップと、

ランダムアクセス媒体のボリューム中の論理ブロックアドレスへ前記書き込みログ中のエントリーをマッピングするための再マッピング木を保持するステップと、

(a) 前記書き込みログが前記書き込みログへの新たなエントリーが許されない状態にある場合に、前記ランダムアクセス媒体上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するステップと、

(b) 前記ステップ(a)が実行されない場合であって、前記ランダムアクセス媒体に関する入力/出力動作がなく、前記書き込みログが第1の所定のサイズより大きい場合に、前記ランダムアクセス媒体上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するステップと、

(c) 前記ステップ(b)が実行されない場合であって、前記書き込みログが前記ランダムアクセス媒体に関する前記入力/出力動作に関わらず前記第1の所定のサイズより小

10

20

さい第 2 の所定のサイズよりも大きい場合に、前記ランダムアクセス媒体上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するステップと、

(d) 前記ステップ (c) が実行されない場合であって、前記再マッピング木が第 3 の所定のサイズより大きい場合に、前記ランダムアクセス媒体上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するステップと
を含む方法。

【請求項 2】

前記データセットが前記書き込みログにおいて連続的に統合することに適しているか否かを決定する前記ステップは、

前記データセットが第 4 の所定のサイズより小さいか否かを決定するステップと、

前記データセットが前記第 4 の所定のサイズより小さいと決定される場合に、前記データセットが前記書き込みログにおいて連続的に統合することに適していることを決定するステップと

をさらに含む請求項 1 に記載の機械により実施される方法。

【請求項 3】

一定の時間において、それぞれの全体のチェックポイント又はそれぞれの増分チェックポイントを前記書き込みログへ書き込むステップであって、前記それぞれの全体のチェックポイントの各々又は前記それぞれの増分チェックポイントの各々は前記書き込みログの内容を記述する情報を含む、ステップ

をさらに含む請求項 1 に記載の機械により実施される方法。

【請求項 4】

好ましくないシャットダウンの後に前記書き込みログを回復するステップをさらに含み、前記回復するステップは、

前記書き込みログの内容について記述する情報を各々が含む少なくとも 1 つのチェックポイントが前記書き込みログに含まれるか否かを決定するステップと、

前記少なくとも 1 つのチェックポイントが前記書き込みログに含まれると決定される場合に、前記書き込みログにおいて最後のチェックポイントを見つけるステップと、

前記最後のチェックポイントに少なくとも部分的に基づいて、前記再マッピング木を再構築するステップと、

前記書き込みログについての追加の有効なエントリーが見つからなくなるまで、任意のチェックポイントに含まれる対応する情報を有さない前記書き込みログの第 1 のエントリーに対応する、前記書き込みログのエントリーから読み出しをするステップと、

前記書き込みログの有効なエントリーが見つかる場合に前記再マッピング木を更新するステップと、

前記書き込みログの最後に見つかった有効なエントリーに基づいて前記書き込みログの終端ログ・シーケンス番号を設定するステップと

をさらに含む請求項 1 に記載の機械により実施される方法。

【請求項 5】

少なくとも 1 つのチェックポイントが前記書き込みログに含まれるか否かの決定により、前記書き込みログにチェックポイントが含まれないことが決定される場合に、前記書き込みログの回復は、

前記書き込みログの初めから、前記書き込みログについての追加の有効なエントリーが見つからなくなるまで、前記書き込みログを読み取るステップと、

前記書き込みログの見つかった有効なエントリーに基づいて前記再マッピング木を更新するステップと

をさらに含む請求項 4 に記載の機械により実施される方法。

【請求項 6】

前記書き込みログの前記少なくとも 1 つのエントリーが無効であることを示す無効の記録を、前記書き込みログの少なくとも 1 つのエントリーに書き込むステップをさらに含み、

10

20

30

40

50

前記書き込みログの無効にされたエントリーは、前記ランダムアクセス媒体上のそれぞれの意図した宛先へ排出されない請求項 1 に記載の機械により実施される方法。

【請求項 7】

前記書き込みログは、
前記ランダムアクセス媒体の**ボリューム**内、
前記ランダムアクセス媒体の異なる**ボリューム**内、又は
異なるランダムアクセス媒体の異なる**ボリューム**内
に含まれる請求項 1 に記載の機械により実施される方法。

【請求項 8】

ディスク上に存在する書き込みログにデータセットを連続的に統合するためのシステム 10
であって、

少なくとも 1 つのプロセッサと、
前記少なくとも 1 つのプロセッサに接続されたメモリーと
を備え、前記メモリーは、

データセットのサイズが所定のサイズより小さい場合に、前記データセットが前記書き込みログに連続的に書き込むのに適しているか否かを決定するための命令と、

前記データセットが連続的に書き込むのに適していると決定される場合に、前記ディスク上の意図した宛先へ書き込む代わりに、前記データセットを前記書き込みログに連続的に書き込むための命令と、

前記ディスクの**ボリューム**内の前記書き込みログ中のエントリーと論理ブロックアドレスとの間でマッピングするための再マッピング木を保持するための命令と、 20

(a) 前記書き込みログが前記書き込みログへの新たなエントリーが許されない状態にある場合に、前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と、

(b) 前記命令 (a) が実行されない場合であって、前記ディスクに関する入力 / 出力動作がなく、前記書き込みログが第 1 の所定のサイズより大きい場合に、前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と、

(c) 前記命令 (b) が実行されない場合であって、前記書き込みログが前記ディスクに関する前記入力 / 出力動作に関わらず前記第 1 の所定のサイズより小さい第 2 の所定のサイズよりも大きい場合に、前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と、 30

(d) 前記命令 (c) が実行されない場合であって、前記再マッピング木が第 3 の所定のサイズより大きい場合に、前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と
を含むシステム。

【請求項 9】

前記メモリーは、
前記ディスクの論理ブロックアドレスが前記書き込みログの 1 つ以上のエントリーに対して無効であることを示す無効の記録を、前記書き込みログに書き込むための命令をさらに含む請求項 8 に記載のシステム。 40

【請求項 10】

前記メモリーは、
前記ディスク上の前記それぞれの意図した宛先へ、前記書き込みログの前記エントリーのうち無効にされていないもののみを排出するための命令をさらに含む請求項 9 に記載のシステム。

【請求項 11】

前記再マッピング木は a v l 木である請求項 8 に記載のシステム。

【請求項 12】

前記メモリーは、
時間間隔においてそれぞれの全体のチェックポイント又はそれぞれの増分チェックポイ 50

ントを前記書き込みログへ書き込むための命令をさらに含み、それぞれの全体のチェックポイント又はそれぞれの増分チェックポイントの各々１つは、前記それぞれの全体のチェックポイント又は前記それぞれの増分チェックポイントの前記１つが書き込まれた時における前記書き込みログの内容を記述する情報を含む請求項８に記載のシステム。

【請求項１３】

前記メモリーは、

好ましくないシャットダウンの後に、少なくとも１つのチェックポイントが前記書き込みログに含まれるか否かを決定するための命令と、

少なくとも１つのチェックポイントが前記書き込みログに含まれると決定される場合に、前記書き込みログにおいて最後のチェックポイントを見つけるための命令と、

前記最後のチェックポイントに少なくとも部分的に基づいて、前記再マッピング木を再構築するための命令と、

前記書き込みログについての追加の有効なエントリーが見つからなくなるまで、任意のチェックポイントに含まれる情報を有さない前記書き込みログの第１のエントリーに対応する、前記書き込みログのエントリーから読み出しをするための命令と、

前記書き込みログの有効なエントリーが見つかった後、前記再マッピング木を更新するための命令と

をさらに含む請求項８に記載のシステム。

【請求項１４】

少なくとも１つのプロセッサに対する命令を記録した有形の機械読み取り可能な媒体であって、前記命令は、

フラッシュベースの記憶装置又はディスクに書き込まれるよう意図されたデータが前記フラッシュベースの記憶装置又は前記ディスクの**ボリューム**内に含まれる書き込みログに連続的に書き込むのに適しているか否かを決定するための命令と、

前記データが前記書き込みログに連続的に書き込むのに適していると決定される場合に、前記データを前記書き込みログに連続的に書き込むための命令と、

前記書き込みログ中のエントリーと前記フラッシュベースの記憶装置又は前記ディスクの論理ブロックアドレスとの間でマッピングするための再マッピング木を保持するための命令と、

(a) 前記書き込みログが前記書き込みログへの新たなエントリーが許されない状態にある場合に、前記フラッシュベースの記憶装置又は前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と、

(b) 前記命令 (a) が実行されない場合であって、前記フラッシュベースの記憶装置又は前記ディスクに関する入力 / 出力動作がなく、前記書き込みログが第１の所定のサイズより大きい場合に、前記フラッシュベースの記憶装置又は前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と、

(c) 前記命令 (b) が実行されない場合であって、前記書き込みログが前記フラッシュベースの記憶装置又は前記ディスクに関する前記入力 / 出力動作に関わらず前記第１の所定のサイズより小さい第２の所定のサイズよりも大きい場合に、前記フラッシュベースの記憶装置又は前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と、

(d) 前記命令 (c) が実行されない場合であって、前記再マッピング木が第３の所定のサイズより大きい場合に、前記フラッシュベースの記憶装置又は前記ディスク上のそれぞれの意図した宛先へ前記書き込みログのエントリーを排出するための命令と

を含む有形の機械読み取り可能な媒体。

【請求項１５】

前記命令は、

前記書き込みログにチェックポイントを書き込む命令をさらに含み、前記チェックポイントの内容は前記書き込みログの内容を記述する請求項１４に記載の有形の機械読み取り可能な媒体。

【請求項 16】

前記命令は、

前記チェックポイントに少なくとも部分的に基づいて、シャットダウンの後に前記再マッピング木を回復する命令をさらに含む請求項 15 に記載の有形の機械読み取り可能な媒体。

【請求項 17】

前記再マッピング木を回復する命令は、

前記書き込みログの次の有効なエントリーについて、前記チェックポイントに含まれる対応する情報を有さない前記書き込みログの第 1 のエントリーに対応する、前記書き込みログのエントリーから開始して、前記書き込みログを読み出す命令と、

前記書き込みログの前記次の有効なエントリーに基づいて前記再マッピング木を更新する命令と

をさらに含む請求項 16 に記載の有形の機械読み取り可能な媒体。

【請求項 18】

前記命令は、

前記書き込みログの少なくとも 1 つのエントリーに対応する、前記フラッシュベースの記憶装置又は前記ディスクの論理ブロックアドレスを示す無効の記録を、前記書き込みログに書き込む命令と、

前記書き込みログの無効にされたエントリーを、前記フラッシュベースの記憶装置又は前記ディスク上の意図した宛先へ排出しない命令と

をさらに含む請求項 14 に記載の有形の機械読み取り可能な媒体。

【請求項 19】

前記フラッシュベースの記憶装置又は前記ディスクに書き込まれるよう意図されるデータがストリームに含まれるか否かを決定する命令と、

前記データが前記ストリーム中にあると決定される場合に、前記データが前記書き込みログへ書き込むのに適さないことを決定する命令と

をさらに含む請求項 14 に記載の有形の機械読み取り可能な媒体。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、データ・ストアへの書き込みトラフィック (write traffic) の統合 (集約、aggregation) に関する。

【背景技術】**【0002】**

[0001]例えばディスク・ドライブ及びフラッシュベースの記憶装置などのデータ・ストアは、データを連続的に書き込む場合に最も効率的であり、ランダムに散乱されたデータを書き込む場合に非常に非効率的である。さらに、ディスク・ドライブは、連続するデータを読み取る場合に最も効率的であり、ランダムに散乱されたデータを読み取る場合に非常に非効率的である。典型的な現在のディスク・ドライブは、ディスク・ドライブがディスク上の任意の位置を求めるのに要する時間におよそ 700 キロバイト (KB) のデータを移動させることができる。技術が進歩するにつれて、ディスク・ドライブは同じ期間中により大きな量のデータを移動させることができる。

【発明の概要】**【発明が解決しようとする課題】****【0003】**

[0002]ほとんどのデータ転送は 700 KB よりずっと小さい。その結果、ディスク・ドライブは、空いていない (non-idle) 期間中にかなりの時間を使ってディスク上の位置を求めることになりかねない。

【課題を解決するための手段】**【0004】**

[0003]この概要は、詳細な説明において以下にさらに述べられる概念の選択を単純化された形式で紹介するために提供される。この概要は、特許請求された主題の重要な特徴又は不可欠な特徴を識別するようには意図されず、特許請求された主題の範囲を限定するために使用されるようにも意図されない。

【 0 0 0 5 】

[0004]処理装置及び機械により実施される方法は、もともとはある量のランダムアクセス媒体 (a volume of random access medium) に書き込まれるように意図されたデータを、書き込みログに対して連続的に統合するために提供されてもよい。処理装置は、データが書き込みログへの書き込みに適しているか否かを決定してもよい。データが書き込みログに書き込むのに適していると決定される場合、処理装置は、書き込みログに対してデータを連続的に統合し又は書き込んでよい。そうでない場合、データはもともと意図されていた宛先に書き込まれてもよい。書き込みログのエントリーは、多くの条件のうち少なくとも1つが生じると、もともと意図されていた宛先へ排出されてもよい。当該条件は次のものを含み得るが、これらに限定されない：書き込みログが、新たなログ・エントリーが許されない状態にあること；書き込みログが存在するランダムアクセス媒体に関する入力又は出力動作がないこと；書き込みログが少なくとも第1の所定の量を満たしていること；ランダムアクセス媒体の入力又は出力動作に関わらず、書き込みログが少なくとも第2の所定の量を満たしていること；及び、再マッピング木 (remapping tree) が所定のサイズより大きいこと。再マッピング木は、書き込みログの1つ以上のエントリーへランダムアクセス媒体のアドレス範囲をマッピングしてもよい。

【 0 0 0 6 】

[0005]幾つかの実施例において、それぞれのチェックポイントは時間間隔で書き込みログに書き込まれてもよい。それぞれのチェックポイントの各々は、書き込みログの内容について記述する情報を含んでもよく、そのため、再マッピング木は好ましくないシャットダウンの後に再構築されてもよい。それぞれのチェックポイントに含まれる情報は、再マッピング木のノードから得られてもよい。

【 0 0 0 7 】

[0006]上述の及び他の利点及び特徴が得られる方法について記載するために、より具体的な特定の説明が以下に述べられ、添付された図面中に示されるその具体的な実施例への言及によって与えられる。これらの図面が典型的な実施例を描いているに過ぎず、したがって、その範囲に限定していると考えるべきでないことを理解し、実施例は添付の図面の使用を通じて追加的な特異性及び詳細をもって記載され説明される。

【図面の簡単な説明】

【 0 0 0 8 】

【図1】[0007]本開示の主題と一致する実施例を実施するために使用され得る例示的な処理装置の機能ブロック図を示す。

【図2】[0008]本開示の主題と一致する実施例において書き込みログを格納するための例示的な円形の記憶装置を示す。

【図3】[0009]ある量のランダムアクセス媒体における例示的な書き込みログを示す。

【図4】[0010]図3に示される例示的な書き込みログの詳細な図を示す。

【図5】[0011]本開示の主題と一致する再マッピング木の例示的な構造を示す。

【図6】[0012]書き込みログのエントリーへマッピングされる、ランダムアクセス媒体上に格納されたデータを上書きする効果を示す。

【図7】[0013]ランダムアクセス媒体に書き込むための受信された書き込み要求、書き込みログへの書き込みに関する書き込み完了の指示を処理するための例示的な処理を示すフローチャートである。

【図8】ランダムアクセス媒体に書き込むための受信された書き込み要求、書き込みログへの書き込みに関する書き込み完了の指示を処理するための例示的な処理を示すフローチャートである。

【図9】[0014]書き込みログ中の例示的なエントリー及びエントリーのフラッシング (fl

10

20

30

40

50

ushing)を示す。

【図10】[0015]書き込みログのフラッシュされたエントリーを記録する例示的な処理を示すフローチャートである。

【図11】書き込みログのフラッシュされたエントリーを記録する例示的な処理を示すフローチャートである。

【図12】[0016]書き込みログに書き込み、書き込みログがどれだけ満たされるかに基づいて書き込みログの状態を変更する例示的な処理のフローチャートである。

【図13】[0017]書き込みログを排出するべきか否かを決定する例示的な処理のフローチャートである。

【図14】[0018]ランダムアクセス媒体へ書き込みログを効率的に排出する排出書き込み計画(drain write plan)を生成する例示的な処理のフローチャートである。

【図15】ランダムアクセス媒体へ書き込みログを効率的に排出する排出書き込み計画を生成する例示的な処理のフローチャートである。

【図16】ランダムアクセス媒体へ書き込みログを効率的に排出する排出書き込み計画を生成する例示的な処理のフローチャートである。

【図17】[0019]規則的な期間で書き込みログにチェックポイントを書き込む例示的な処理のフローチャートである。

【図18】[0020]好ましくないシャットダウンの後に書き込みログを回復する例示的な処理のフローチャートである。

【発明を実施するための形態】

【0009】

[0021]実施例が以下に詳細に説明される。具体的な実施例が説明されるが、これが説明の目的のみのためになされることが理解されるべきである。関連する技術の当業者であれば、本開示の主題の趣旨及び範囲から逸脱することなく他のコンポーネント及び構成が使用され得ることを認識するであろう。

概観

[0022]例えば、ディスク、フラッシュベース記憶装置、又は他のランダムアクセス媒体などのランダムアクセス媒体に書き込む要求を受信する方法及びシステムが提供される。その後、書き込み要求は、当該要求が書き込みログに連続的に統合するのに適しているか否かを決定するために分析されてもよい。当該要求は、当該要求がランダムアクセス媒体に所定の量未満のデータを書き込むためのものである場合に、書き込みログに連続的に統合するのに適していると決定されてもよい。幾つかの実施例において、700KB未満であり得る調整可能なしきい値データをランダムアクセス媒体に書き込む要求は、書き込みログに連続的に統合するのに適していると決定されてもよい。

【0010】

[0023]書き込みログは、もともと意図されていた宛先としての同一の量のランダムアクセス媒体内、もともと意図されていた宛先としての異なる量のランダムアクセス媒体内、又はもともと意図されていた宛先としての異なる量の異なるランダムアクセス媒体内に存在してもよい。書き込みログは多くのエントリーを含んでもよく、各エントリーはそれぞれのエントリー及び対応するペイロードデータについて記述するメタデータを有している。ランダムアクセス媒体へのデータ書き込みの要求が書き込みログへの書き込みに適していると決定される場合、データは書き込みログの最後に加えられてもよい。再マッピング木は、書き込みログの外部で、書き込みログの対応するエントリーへ、ランダムアクセス媒体の位置をマッピングしてもよいし、その逆を実行してもよい。

【0011】

[0024]排出は、書き込みログ中のエントリーがランダムアクセス媒体の量の意図した宛先に移動される処理である。書き込みログは、多くの条件の下でランダムアクセス媒体の意図した位置へ排出されてもよい。例えば、幾つかの実施例において、条件は次のものを含んでもよい：書き込みログが、書き込みログへの新たなエントリーが許されない状態にあること；書き込みログが第1の所定の量の全体より多いこと；ランダムアクセス媒体に

関する入力／出力動作がなく、書き込みログが第２の所定の量の全体より多く、当該第２の所定の量が第１の所定の量より少ないこと；及び、再マッピング木（remapping tree）が第３の所定のサイズより大きいこと。排出中に、隣接した書き込みは単一の排出書き込みへとマージされ、排出書き込みはシーク時間（seek time）を最小化するように命じられてもよい。

【 0 0 1 2 】

[0025] 幾つかの実施例において、チェックポイントは調整可能な時間間隔で書き込みログに書き込まれてもよい。チェックポイントは、書き込みログの内容を要約する情報を含んでもよい。好ましくないシャットダウンは予期しないシャットダウン又はシステム障害かもしれない。好ましくないシャットダウン又はシステムの正常なシャットダウンの後、書き込みログに書き込まれた最後のチェックポイントに含まれる情報は、再マッピング木の復元によって好ましくないシャットダウン又は正常なシャットダウンから回復するために使用することができる。

10

例示的な処理装置

[0026] 図１は、本開示の主題と一致する実施例を実施するために使用され得る例示的な処理装置１００の機能ブロック図である。処理装置１００は、デスクトップパソコン（ＰＣ）、ノートブックもしくはラップトップＰＣ、サーバー又は他の処理装置であってもよい。処理装置１００はバス１１０、メモリー１３０、読み取り専用メモリー（ＲＯＭ）１４０、プロセッサ１２０及び記憶装置１５０を含んでもよい。バス１１０は、処理装置１００のコンポーネント間の通信を可能にしてもよい。

20

【 0 0 1 3 】

[0027] プロセッサ１２０は、命令を解釈し実行する少なくとも１つの従来のプロセッサ又はマイクロプロセッサを含んでもよい。メモリー１３０は、プロセッサ１２０による実行のための情報及び命令を格納するランダム・アクセス・メモリー（ＲＡＭ）又は別のタイプの動的記憶装置であってもよい。メモリー１３０は、プロセッサ１２０による命令の実行中に使用される、一時的な変数又は他の媒介情報（intermediate information）を格納してもよい。ＲＯＭ１４０は、プロセッサ１２０のための静的情報及び命令を格納する、従来のＲＯＭ装置又は別のタイプの静的記憶装置を含んでもよい。記憶装置１５０はハードディスク及び対応するドライブ、フラッシュベースの記憶装置、又は、プロセッサ１２０のためのデータ及び／又は命令を格納するための他のタイプのランダム・アクセス・データ記憶装置もしくは媒体を含んでもよい。

30

【 0 0 1 4 】

[0028] 処理装置１００は、例えば、メモリー１３０、ＲＯＭ１４０、記憶装置１５０又は他の媒体などの有形な機械可読媒体に含まれる命令のシーケンスをプロセッサ１２０が実行することに応答して、機能を実行してもよい。そのような命令は、通信インターフェース（図示せず）を介して別の機械可読媒体又は別個の装置からメモリー１３０へと読み込まれてもよい。

書き込みログ

[0029] 本開示の主題と一致する実施例において、書き込みログは、ランダムアクセス媒体を対象とした書き込みを連続的に統合する（集約する）ための、例えば、ディスク又は他のランダムアクセス媒体などの、ランダムアクセス媒体内の円形の記憶領域であってもよい。図２は例示的な書き込みログ２００を示す。書き込みログ２００は、開始エントリー２０２及び終了エントリー２０４を含む複数のエントリーを含んでもよい。終了エントリー２０４が書き込みログ２００の最新のエントリーであり得る一方、開始エントリー２０２は書き込みログ２００の最も古いエントリーであり得る。すなわち、新たなエントリーは終了エントリー２０４の後に書き込みログ２００に追加されてもよく、終了エントリー２０４は追加されたエントリーのうち最新のものに対応するように更新されてもよい。

40

【 0 0 1 5 】

[0030] 図３は、本開示の主題と一致するランダムアクセス媒体３００上に格納され得る例示的な構造を示す。ランダムアクセス媒体３００は、ランダムアクセス媒体３００上に

50

記録されたブロック記憶構造 302 及び書き込みログ 304 を含んでもよい。ブロック記憶構造 302 は、ファイルシステムによって使用するべきでないランダムアクセス媒体 300 の領域に関する情報を含んでもよい。当該領域は、ある量のランダムアクセス媒体 300 を対象としたデータを連続的に統合するために書き込みログ 304 について確保された領域に対応してもよい。幾つかの実施例において、書き込みログ 304 は、もともと意図されていた宛先とは異なる量のランダムアクセス媒体 300 内、又は異なる量の異なるランダムアクセス媒体内に含まれてもよい。

【0016】

[0031] 書き込みログ 304 は、書き込みログ 304 の始めと終わりにそれぞれ配置され得る、第 1 の制御領域 306 及び第 2 の制御領域 308 を含んでもよい。第 1 及び第 2 の制御領域は、書き込みログ 304 の開始エントリーに関する情報及び書き込みログ 304 の物理的終端 (physical end) に関する情報を含むがこれらに限定されない、書き込みログ 304 に関する情報を含み得る。幾つかの実施例において、開始エントリーに関する情報は、書き込みログ 304 の開始ログ・シーケンス番号 (LSN) を含んでもよく、書き込みログ 304 の物理的終端に関する情報はラッピング (wrapping) LSN を含んでもよい。幾つかの実施例は、第 1 の制御領域 306 及び第 2 の制御領域 308 内の最後のチェックポイントの位置に関するポインター又は他の情報などの情報を含んでもよい。

【0017】

[0032] 書き込みログ 304 の各エントリーは LSN に対応する。LSN はラップ・カウント (wrap count) 部分及びオフセット部分を含んでもよい。ラップ・カウント部分は、書き込みログ 304 の物理的な先頭 (physical beginning) に戻るラッピング (wrapping) の発生回数であってもよい。オフセット部分は、書き込みログ 304 の物理的な先頭からの、セクター又は測定の他の尺度 (unit of measurement) のオフセットであってもよい。ラッピング LSN は、書き込みログ 304 の物理的な先頭に対するラッピング及びラップ・カウントの増加が生じ得る LSN であってもよい。例えば、ラッピング LSN が <ラップ・カウント>400 であり、400 が書き込みログ 304 の物理的な先頭からのオフセットである場合、その後現在の LSN が例えば 2 . 400 に達した場合、書き込みログ 304 の物理的な先頭に対するラッピングが生じてよいし、現在の LSN が 3 . 0 になってもよい。

【0018】

[0033] 書き込みログ 304 は多くのエントリーを含んでもよい。図 4 に示されるように、エントリーの各々はエントリー・メタデータ 402 及びエントリー・ペイロードデータ 404 を含んでもよい。エントリー・メタデータ 402 は、それぞれのエントリーに関する情報を含んでもよい。例えば、エントリー・メタデータ 402 は、それぞれのエントリーのサイズに関する情報、及びある量のランダムアクセス媒体 300 における意図した宛先の論理ブロックアドレス (LBA) のほか、異なる又は他の情報を含んでもよい。さらに、エントリー・メタデータ 402 は、それぞれのエントリーが有効か否かを示すチェックサム又は他のデータを含んでもよい。幾つかの実施例において、エントリー・メタデータ 402 及びエントリー・ペイロードデータ 404 は整列されたセクターであってもよい。

【0019】

[0034] 再マッピング木は、例えばメモリー又は RAM などの動的記憶装置において保持されてもよい。再マッピング木は、例えば、ある量のランダムアクセス媒体中の LBA などの意図した宛先を、書き込みログ中のエントリーへマッピングしてもよいし、その逆を実行してもよい。再マッピング木は AVL 木として構築されてもよい。AVL 木は、任意のノードの子サブツリー (child subtrees) の高さがわずか 1 ずつ異なってもよい、自己均衡 (self-balancing) の二分探索木である。

【0020】

[0035] 図 5 は例示的な再マッピング木を示す。例示的な再マッピング木に見られるように、すべての子サブツリーの高さは、同じレベルで始まる他の子サブツリーに対してわず

10

20

30

40

50

かに 1 だけ異なる。もちろん、他の実施例において、再マッピング木は異なるデータ構造かもしれない。再マッピング木の各ノードは、例えば、意図した宛先の L B A、書き込みログ、書き込みログにおける対応する L S N、エントリーの種類について記述する情報のほか、他の又は異なる情報などの情報を含んでもよい。

【 0 0 2 1 】

[0036]ランダムアクセス媒体の L B A が 2 度書き込まれている場合、第 2 の書き込みは第 1 の書き込みを無効にする。同様に、L B A への第 1 の書き込みが書き込みログにリダイレクト (redirect) され、L B A への第 2 の書き込みもまた書き込みログにリダイレクトされる場合、第 1 の書き込みに対応する書き込みログ中のエントリーは、第 2 の書き込みに対応する書き込みログ中のエントリーによって取って代わられてもよい。これは内部上書き (internal overwrite) と呼ばれてもよい。

10

【 0 0 2 2 】

[0037]しかし、第 1 の書き込みが書き込みログにリダイレクトされる一方で第 2 の書き込みがある量のランダムアクセス媒体中の L B A に直接書き込まれる場合、第 1 の書き込みに対応する書き込みログのエントリーはもはや有効でない。第 1 の書き込みに対応するエントリーがもはや有効ではないことを示す無効の記録が書き込みログに追加されてもよい。

【 0 0 2 3 】

[0038]図 6 は、ある量のランダムアクセス媒体 (a volume of a random access medium) 6 0 0 を示す。書き込みログ 6 0 2 は容量 6 0 0 に含まれていてもよい。第 1 の書き込みは書き込みログ 6 0 2 にリダイレクトされ、書き込みログ 6 0 2 のエントリー 6 0 4 として記録されてもよい。第 1 の書き込みは容量 6 0 0 における宛先 6 0 8 に対応する意図された宛先を有してもよい。第 2 の書き込み 6 0 6 は宛先 6 0 8 に直接書き込まれてもよい。したがって、エントリー 6 0 4 は今無効かもしれない。無効の記録 6 1 0 はエントリー 6 0 4 が今無効であることを示すために書き込みログ 6 0 2 に追加されてもよい。その後、書き込みログ 6 0 2 のエントリーが容量 6 0 0 のそれぞれの意図した宛先へ排出されるべきである場合、エントリー 6 0 4 は宛先 6 0 8 へ排出されなくてもよい。というのは、そうすることによって宛先 6 0 8 が上書きされ、宛先 6 0 8 の内容を破壊するからである。代わりに、エントリー 6 0 4 が単に廃棄されてもよい。

20

例示的な手続

30

[0039]図 7 は、受信された書き込み要求が書き込みログに連続的に統合するのに適しているか否かを決定し、そうであるならば、書き込みログへ書き込み要求を転送する、例示的な処理を示すフローチャートである。この処理は、受信された書き込み要求が書き込みログに連続的に統合するのに適しているか否かを決定する処理装置から始まってもよい (ステップ 7 0 4)。

【 0 0 2 4 】

[0040]図 8 は、図 7 のステップ 7 0 4 を行う例示的な処理を示すフローチャートである。その処理は、受信された書き込み要求がストリームに向けたものであるか否かを決定する処理装置から始まってもよい (ステップ 8 0 2 ; 図 8)。処理装置は、読み書きのパターンが連続するか否かを決定するために読み書き要求の容量オフセット (volume offsets) を監視することによってストリームを検出してよい。読み書きのパターンが連続的であると決定される場合、処理装置はストリームを検出し、当該ストリームを書き込みログへ向かわせなくてもよい。ストリームが検出される場合、処理装置は、受信された書き込み要求が書き込みログにリダイレクトするのに適していないことを示してもよく (ステップ 8 0 8)、処理は完了してもよい。

40

【 0 0 2 5 】

[0041]ステップ 8 0 2 中、受信された書き込み要求がストリームに向けたものでないと処理装置が決定する場合、処理装置は、書き込まれるべきデータ、又はデータセットが、7 0 0 K B 又は別のサイズなどの所定のサイズ未満か否かを決定してもよい (ステップ 8 0 4)。データ、又はデータセットが、所定のサイズ以上である場合、処理装置は、書き

50

込み要求が書き込みログにリダイレクトするのに適していないことを示してもよく（ステップ 808）、処理は完了してもよい。

【0026】

[0042]ステップ 804 中、書き込まれるべきデータ、又はデータセットが、所定のサイズ未満であると処理装置が決定する場合、処理装置は、書き込み要求が書き込みログにリダイレクトするのに適していることを示してもよい（ステップ 806）。その後、処理は完了してもよい。

【0027】

[0043]図 7 に戻り、ステップ 704 中に、書き込み要求が書き込みログに連続的に統合するのに適していると処理装置が決定する場合、その後、処理装置は、書き込みログの状態が無効のみか否かを決定して、無効の記録のみが書き込みログに書き込まれてもよいことを示してもよい（ステップ 706）。書き込みログの状態が無効のみでない場合、処理装置は、書き込みログへのエントリーとして、データ、又はデータセットを書き込んでよい（ステップ 708）。

10

【0028】

[0044]書き込みログへの書き込みがうまく完了する場合、処理装置は、ある量のランダムアクセス媒体中の意図した宛先を書き込みログ中のエントリーへマッピングするために再マッピング木を更新してもよい（ステップ 716）。

【0029】

[0045]ステップ 704 中に、書き込まれるべきデータ、又はデータセットが、書き込みログに連続的に統合するのに適しないと処理装置が決定する場合、データ、又はデータセットは、ランダムアクセス媒体の容量中の意図した宛先に書き込まれてもよい（ステップ 710）。その後、処理装置は、再マッピング木を参照することにより、意図した宛先の LBA のうちのいずれかが書き込みログの 1 つ以上のエントリーに対応するか否かを決定してもよい（ステップ 712）。意図した宛先の LBA のうちのいずれかが書き込みログの 1 つ以上のエントリーに対応する場合、書き込みログの 1 つ以上のエントリーに対応する意図した宛先の無効を示すために、1 つ以上の無効の記録が書き込みログに追加されてもよい（ステップ 714）。

20

【0030】

[0046]書き込みログへの書き込みは順序がばらばらに完了してもよい。例えば、図 9 は、書き込みログのエントリー及びエントリーの各々についてのそれぞれの LSN を示す。影付きのエントリーは、保留中の（pending）書き込みを有していない（つまり、エントリーについての書き込みログへのそれぞれの書き込みが完了している）エントリーに対応する。影付きでないエントリーは保留中の書き込みを有するエントリーに対応する。したがって、図 9 によれば、LSN 0.256 によって示されたエントリーは保留中の書き込みを有する。

30

【0031】

[0047]フラッシュされた（排出された、flushed）LSN は書き込みログのエントリーに対応する LSN であり、そのため、対応するエントリーに先行する書き込みログのすべてのエントリーは保留中の書き込みを有さない。開始 LSN は、排出のために考慮されるか又はチェックポイントに含められるべき、書き込みログの第 1 のエントリーに対応する LSN である。書き込みログの排出及びチェックポイントの生成は、ともに以下に述べられるが、フラッシュされた LSN を超えたエントリーを考慮しなくてもよい。フラッシング（flushing）LSN はフラッシングのための目標エントリーに対応する LSN である。終了 LSN（end LSN）は、書き込みログの論理終端（logical end）に対応する LSN である。

40

【0032】

[0048]図 10 及び 11 は、書き込みログにおいてフラッシュされた LSN を記録する例示的な処理を示すフローチャートである。その処理は、ログへの任意の書き込みが生じる前にフラッシング LSN を終了 LSN へと初期化し（ステップ 1002）、フラッシュさ

50

れた L S N を開始 L S N へ初期化する (ステップ 1 0 0 4) 処理装置から始まってよい。現在のフラッシング・カウント (flushing count) は、L S N を排出する前に発行された書き込みログへの未完了の書き込みの数であってもよい。現在のフラッシング・カウントは、0 に初期化されてもよい (ステップ 1 0 0 6)。次のフラッシング・カウントは、L S N を排出した後に発行された書き込みログへの未完了の書き込みの数であってもよい。次のフラッシング・カウントは、0 に初期化されてもよい (ステップ 1 0 0 8)。

【 0 0 3 3 】

[0049]次に、書き込みの完了は、書き込みログへの書き込みに関して受信されてもよい (ステップ 1 0 1 0)。処理装置は、L S N を排出する前に書き込みの完了が書き込みログ・エントリーに向けられたものであるか否かを決定してもよい (ステップ 1 0 1 2)。10
そうでなければ、その後、処理装置は次のフラッシング・カウントをデクリメントし (減少させ) てもよく、それは、L S N を排出した後のエントリーに関する保留中の書き込みの数を示してもよい (ステップ 1 0 1 4)。その後、処理装置はステップ 1 0 1 0 - 1 0 1 2 を繰り返してもよい。

【 0 0 3 4 】

[0050]ステップ 1 0 1 2 中に、書き込みの完了がフラッシング L S N に対応するエントリーに先行する書き込みログ・エントリーに向けたものであることを処理装置が決定する場合、処理装置は現在のフラッシング・カウントをデクリメントしてもよい (ステップ 1 0 1 6)。

【 0 0 3 5 】

[0051]その後、処理装置は、現在のフラッシング・カウントがゼロに等しいか否かを決定して、フラッシング L S N に対応するエントリーに先行するすべての書き込みログ・エントリーが完了したことを示してもよい (ステップ 1 1 0 2 ; 図 1 1)。現在のフラッシング・カウントが 0 に等しい場合、処理装置はフラッシュされた L S N をフラッシング L S N に設定してもよい (ステップ 1 1 0 4)。その後、処理装置は、フラッシング L S N を終了 L S N に等しくなるように設定してもよい (ステップ 1 1 0 6)。代替的に、処理装置は、例えば、フラッシュされた L S N の後の所定の数のエントリーにすぎないエントリーに対応する L S N などの、終了 L S N 以外の L S N に等しくなるようにフラッシング L S N を設定してもよい。処理装置は、その後、次のフラッシング・カウントに等しくなるように現在のフラッシング・カウントを設定してもよい (ステップ 1 1 0 8)。その後、30
次のフラッシング・カウントは、0 に設定されてもよい (ステップ 1 1 1 0)。代替的に、フラッシング L S N が終了 L S N 以外の L S N に設定される場合、その後、次のフラッシング・カウントは、フラッシング L S N に対応するエントリーに先行するエントリーについての保留中の書き込みの数を示す値に設定されてもよい。処理装置は再びステップ 1 0 1 0 (図 1 0) を実行してもよい。

【 0 0 3 6 】

[0052]図 1 2 は、書き込みログに書き込む例示的な処理のフローチャートである。その処理は、終了 L S N に対応するエントリーにおいて書き込みログに書き込むための書き込み要求を発行する処理装置から始まってよい (ステップ 1 2 0 2)。次に、フラッシング・カウントがインクリメントされてもよく (ステップ 1 2 0 4)、終了 L S N が更新されてもよい (ステップ 1 2 0 6)。その後、処理装置は、書き込みログが所定量全体より40
大きいかな否か又は再マッピング木が所定のサイズより大きいかな否かを決定してもよい (ステップ 1 2 0 8)。そうであるならば、その後、書き込みログの状態はディスエーブル (disabled) に設定されてもよい (ステップ 1 2 1 0)。書き込みログが所定量全体より大きくない場合、書き込みログの状態はイネーブル (enabled) に設定されてもよく、その結果、無効の記録のほか他の種類のエントリーを含むエントリーが、書き込みログに加えられてもよい (ステップ 1 2 1 2)。

【 0 0 3 7 】

[0053]図 1 3 は、書き込みログの排出を開始する例示的な処理を示すフローチャートである。その処理は、処理装置が電源で動作しているかな否かを決定する処理装置から始まっ50

てもよい(ステップ1302)。そうであるならば、電力を節約するために、排出は行われなくてもよく、処理は終了してもよい。

【0038】

[0054]ステップ1302中に、処理装置が電源で動作していないと処理装置が決定する場合、処理装置は、書き込みログの状態が無効のみであるか又はディスエーブルであるかを決定してもよい(ステップ1304)。書き込みログの状態が無効のみであるか又はディスエーブルである場合、処理装置はランダムアクセス媒体の意図した宛先へ書き込みログを排出してもよい(ステップ1314)。

【0039】

[0055]ステップ1304中に、書き込みログの状態が単に無効ではなく且つディスエーブルでないと処理装置が決定する場合、処理装置は書き込みログが第1の所定量全体より多いか否かを決定してもよい(ステップ1308)。幾つかの実施例において、第1の所定量全体は全部で67%かもしれない。他の実施例において、第1の所定量全体は別の適切な値かもしれない。書き込みログが第1の所定量全体を超える場合、処理装置は書き込みログを排出してもよい(ステップ1314)。

【0040】

[0056]書き込みログが第1の所定量全体以下である場合、処理装置は、ランダムアクセス媒体に関して入出力がないか否か、及び、書き込みログが第2の所定量全体より大きいと処理装置が決定する場合、処理装置は書き込みログを排出してもよい(ステップ1314)。

【0041】

[0057]ランダムアクセス媒体に関して入出力動作があるか、又は書き込みログが第2の所定量全体より大きくないと処理装置が決定する場合、処理装置は再マッピング木が第3の所定のサイズより大きいと処理装置が決定する場合、処理装置は書き込みログを排出してもよい(ステップ1314)。

【0042】

[0058]図14-16は、本開示の主題と一致する実施例において、書き込みログを排出する例示的な処理を示すフローチャートである。その処理は、書き込みログのエントリーに対応するノードの組について再マッピング木をスキャンして、排出する、処理装置から始まってよい(ステップ1402)。その後、処理装置は、書き込みログの始めにおける記録が無効にされたか否かを決定してもよい(ステップ1404)。書き込みログの始めにおける記録が無効にされている場合、処理装置は、無効にされていない再マッピング木において、開始LSNを最低のLSNへと進めてもよい(ステップ1406)。

【0043】

[0059]その後、処理装置は、再マッピング木の複数のノード(例えば、20個のノード又は別の適切な数のノード)を読み取ってもよく、排出書き込み計画(drain write plan)を生成してもよい(ステップ1408)。フラッシュされたLSNに先行するLSNに対応する書き込みログのエントリーのみが排出されてもよい。その後、処理装置は、排出書き込みがアクティブなボリュームの書き込みとオーバーラップするか否かを決定してもよい(ステップ1410)。排出書き込み及びアクティブなボリュームの書き込みがランダムアクセス媒体の少なくとも1つの同じLBAに対するものである場合、排出書き込みは、アクティブなボリュームの書き込みとオーバーラップする。オーバーラップが検出される場合、処理装置はボリューム書き込みが完了するのを待ってもよく(ステップ1412)、処理装置は、ステップ1402を再び実行することによって処理を再び開始してもよい。

【0044】

[0060]ステップ1410中に、排出書き込みがアクティブなボリュームの書き込みとオ

10

20

30

40

50

ーオーバーラップしないと処理装置が決定する場合、処理装置は、隣接した排出書き込みを単一の排出書き込みへとマージして、ランダムアクセス媒体に関して入力及び出力を低減してもよい（ステップ1502；図15）。その後、処理装置は、シーク時間を最小化するために、ボリューム・オフセットによって、排出書き込みを実行してもよい（ステップ1504）。排出書き込みが完了した後、処理装置は開始LSN及び再マッピング木を更新してもよい（ステップ1506）。

【0045】

[0061]次に、処理装置は、書き込みログがディスエーブルになっていたために排出が生じたか否かを決定してもよい（ステップ1508）。そうであるならば、その後、処理装置は、書き込みログが空であるか否かを決定してもよい（ステップ1510）。書き込みログが空である場合、処理は終了してもよい。そうでなければ、処理装置は再びステップ1402を実行してもよい。

【0046】

[0062]ステップ1508中に、書き込みログがディスエーブルになっていたために排出が生じなかったと処理装置が決定する場合、処理装置は、ランダムアクセス媒体に関して入出力動作がないために排出が生じなかったか否かを決定してもよい（ステップ1512）。そうであるならば、処理装置は、書き込みログの排出に関連する活動以外にランダムアクセス媒体に関して入出力動作があったか否かを決定してもよい（ステップ1514）。そうであるならば、その後、処理は終了してもよい。そうでなければ、処理装置は、書き込みログが比較的空であるか否かを決定してもよい（ステップ1516）。書き込みログが、例えば10%又は別の適切な値など、所定量全体より小さい場合、書き込みログは、比較的空であると決定されてもよい。

【0047】

[0063]ステップ1516中に、書き込みログが比較的空ではないと処理装置が決定する場合、処理装置は再びステップ1402（図14）を実行してもよい。

[0064]ステップ1512中に、ランダムアクセス媒体に関して入出力動作がないために排出が生じていないことを処理装置が決定する場合、処理装置は再マッピング木のサイズが所定量より小さいか否かを決定してもよい（ステップ1602；図16）。そうであるならば、処理は終了してもよい。そうでなければ、処理装置は再びステップ1402（図14）を実行してもよい。

【0048】

[0065]図17は、書き込みログにチェックポイントを書き込む例示的な処理のフローチャートである。チェックポイントは、時間間隔における書き込みログのすべてのエントリーについて記述する情報を含んでもよい。全体のチェックポイントは大きくてもよい。大量のデータを書き込むと、処理装置は遅くなり、チェックポイントの書き込み中にユーザーの経験に対してネガティブな影響を与えるかもしれない。幾つかの実施例において、全体のチェックポイントはより小さな増分（インクリメンタル、incremental）チェックポイントの組に分割されてもよい。増分チェックポイントの各々は以前の増分チェックポイントを指してもよい。以前の増分チェックポイントを検討することにより、実質的に、全体のチェックポイントに対応する情報が提供され得る。開始LSNに対応する書き込みログ・エントリーの前に書き込みログ・エントリーが存在しないので、以前のチェックポイントの検討は、開始LSNに対応する書き込みログ・エントリーにおいて終了する。全体のチェックポイントは、チェックポイントのリンクされたリスト中のたった1つのノードを備えた増分チェックポイントと等価であってもよい。以下の文脈においては、全体のチェックポイントが増分チェックポイントの特別な場合として見ることができるので、チェックポイントは増分チェックポイントを指す（参照する）。

【0049】

[0066]処理は、フラッシュされたLSNよりまだ前であって以前のチェックポイントの後にログ範囲について記述する書き込みログにチェックポイントを書き込むための時機（良いタイミング、right moment）を待つ処理装置から開始してもよい（ステップ1702

）。例えば、当該時機は、シャットダウンが開始されようと、不良セクターがチェックポイントされていない（non-checkpointed）書き込みログ・スペースに存在しようと、又は他の条件が存在しようと、以前のチェックポイント以来チェックポイントされていない書き込みログ・スペースに基づいて、決定されてもよい。次に、処理装置は再マッピング木をスキャンしてもよく、書き込みログ中のエントリーへある量のランダムアクセス媒体のLBAをマッピングする（ステップ1704）。その後、処理装置は、再マッピング木を要約するチェックポイントを作成してもよく、終了LSNに対応する書き込みログのエントリーにチェックポイントを書き込んでよく、終了LSNを更新してもよい（ステップ1706）。その後、処理装置はステップ1702 - 1706を繰り返してもよい。フラッシュされたLSNの前に書き込みログのエントリーに関する情報のみがチェックポイントに記録されてもよいことに留意すべきである。

10

【0050】

[0067]好ましくないシャットダウンは予期しないシャットダウン又はシステムクラッシュであり得る。図18は、好ましくないシャットダウンから回復する例示的な処理を示すフローチャートである。幾つかの実施例において、リブート後にスキャンするログ・スペースの量を最小化するために正常なシステム・シャットダウンの前にチェックポイントを書き込むことができることを除いては、正常なシャットダウンは、好ましくないシャットダウンと同じ方法で処理されてもよい。

【0051】

[0068]処理は、書き込みログの第1の制御領域又は第2の制御領域のいずれかから開始LSNを得て、開始LSNが有効か否かを決定する処理装置から開始してもよい（ステップ1802）。処理装置は、開始LSNに対応する書き込みログ・エントリーを読み取ること及び書き込みログ・エントリーのチェックサムが有効であるか否かを決定することにより、又は期待値もしくは予測可能な値について書き込みログ・エントリーの別のフィールドをチェックすることにより、又は他の方法により、開始LSNが有効か否かを決定してもよい。開始LSNが有効でないと決定される場合、書き込みログは回復不能であると考えられ、その影響に対する通知が提供されてもよい（ステップ1804）。

20

【0052】

[0069]開始LSNが有効であると決定される場合、処理装置はいずれかのチェックポイントが書き込みログに存在するか否かを決定してもよい（ステップ1806）。幾つかの実施例において、最後のチェックポイントに対するポインターは、書き込みログの第1の制御領域及び第2の制御領域に格納されてもよい。そのような実施例において、ポインターが得られてもよく、一連のチェックポイントのうち最後のチェックポイントが読み取られて有効にされ（validated）てもよい。他の実施例において、処理装置は、書き込みログの他の有効なエントリーが見つからなくなるまで、チェックポイントのエントリーを探索するために、開始LSNに対応するエントリーにおいて開始して、書き込みログをスキャンしてもよい。

30

【0053】

[0070]少なくとも1つのチェックポイントが書き込みログに存在する場合、処理装置は書き込みログにおいて一連のチェックポイントのうち最後のチェックポイントを見つけてもよい（ステップ1808）。その後、処理装置は、開始LSNに対応する書き込みログ・エントリーがヒットするまですべてのチェックポイントを検討することにより、一連のチェックポイントに基づいて再マッピング木を復元（再構築）してもよい（ステップ1810）。スキャンLSNは、その後、チェックポイントのうちのいずれによっても記述されない第1の書き込みログ・エントリーに対応するように設定されてもよい（ステップ1811）。

40

【0054】

[0071]ステップ1806中に、チェックポイントが書き込みログに存在しないと処理装置が決定する場合、処理装置はスキャンLSNを開始LSNに設定してもよい（ステップ1820）。ステップ1810又はステップ1820を実行した後、処理装置はスキャン

50

L S Nに対応した後の、次の有効なエントリーを読み取ってもよい（ステップ１８１２）。幾つかの実施例において、書き込みログの１つ以上の無効なエントリーが書き込みログの有効なエントリーの間に存在してもよい。その後、処理装置は、有効なエントリーが見つかったか否かを決定してもよい（ステップ１８１４）。有効なエントリーが見つかった場合、処理装置は、再マッピング木を更新し（ステップ１８１６）、再びステップ１８１２ - １８１４を実行してもよい。ステップ１８１４中に、書き込みログの有効なエントリーが見つからなかったと処理装置が決定する場合、処理装置は、書き込みログの最後に見つかった有効なエントリーに対応するように終了L S Nを設定してもよい（ステップ１８１８）。その後、処理は終了してもよい。

結論

10

【0072】本願の主題は構造的特徴及び／又は方法論のステップに特有の言葉で記載されたが、添付の特許請求の範囲における主題が上述の具体的な特徴又はステップに必ずしも限定されないことが理解されるべきである。むしろ、上述の具体的な特徴及びステップは、請求項の実施のための例示的な形式として開示されるものである。

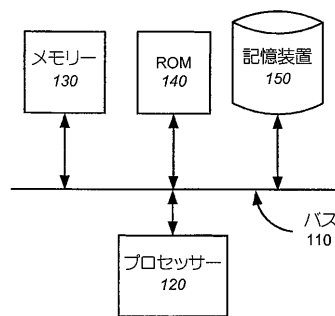
【0055】

【0073】上記の記載は特定の詳細を含んでいるかもしれないが、特許請求の範囲をいかなるようにも限定するものとして解釈されない。記載された実施例の他の構成は本開示の範囲の一部である。さらに、本開示の主題と一致する実施例は図７、８及び１０ - １８に関して記載されたようなものより多くの又は少ないステップを有してもよく、又は示された順序とは異なる順序でステップを実施してもよい。したがって、添付された特許請求の範囲及びその法的な均等物は、いかなる所与の具体的な例でもなく、本発明を規定するものである。

20

【図１】

100



【図２】

200

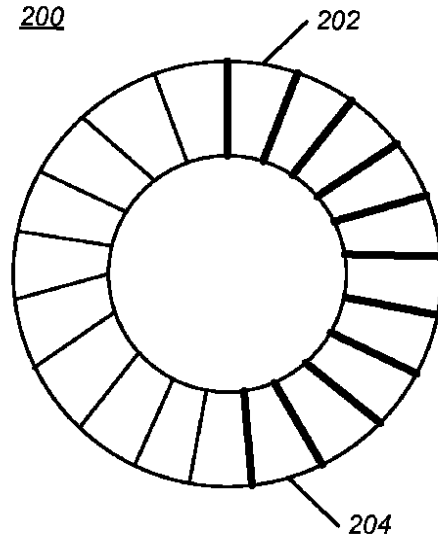
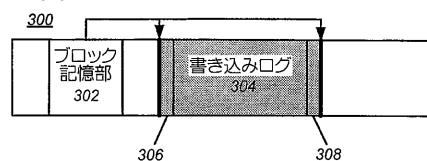
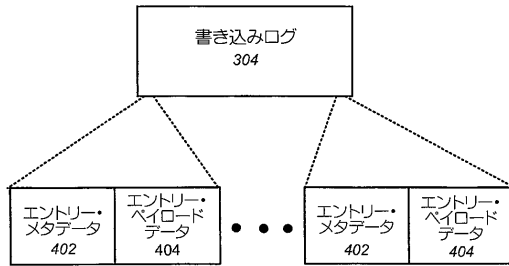


FIG. 2

【図３】



【図 4】



【図 5】

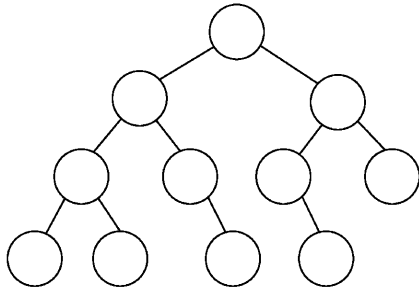


FIG. 5

【図 6】

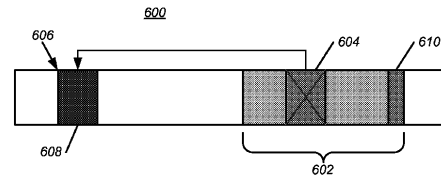
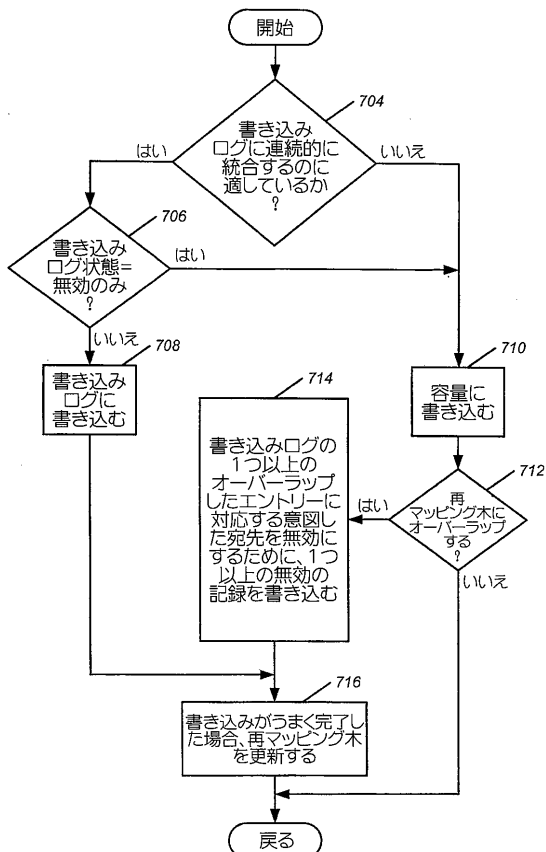
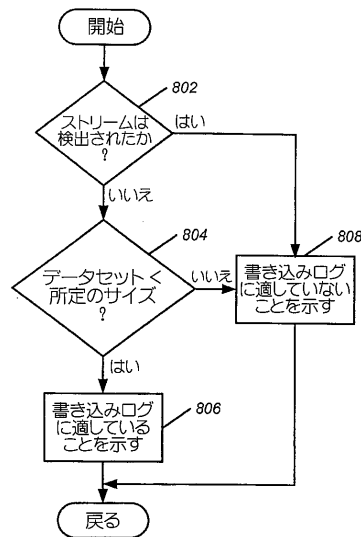


FIG. 6

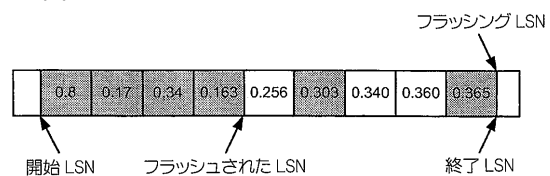
【図 7】



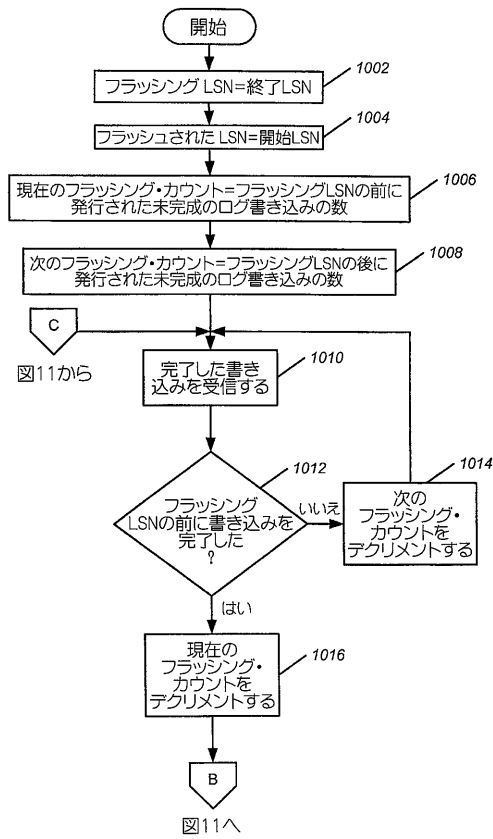
【図 8】



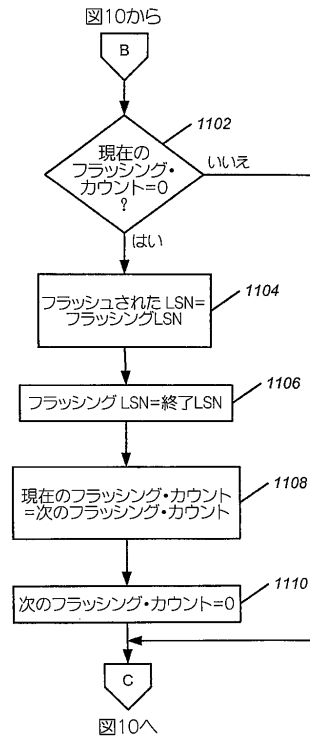
【図 9】



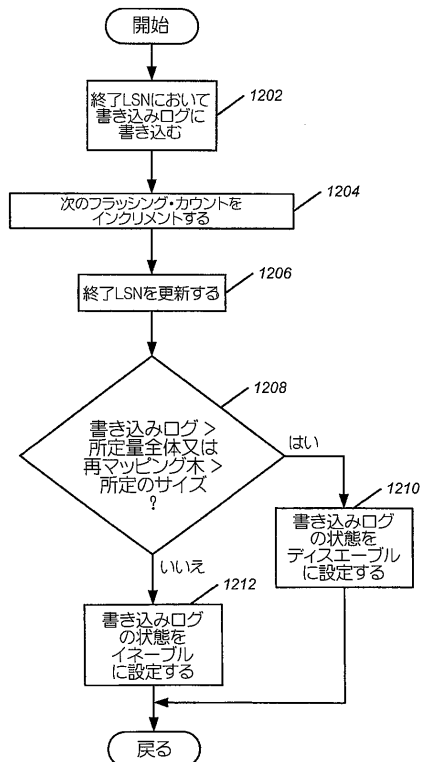
【図10】



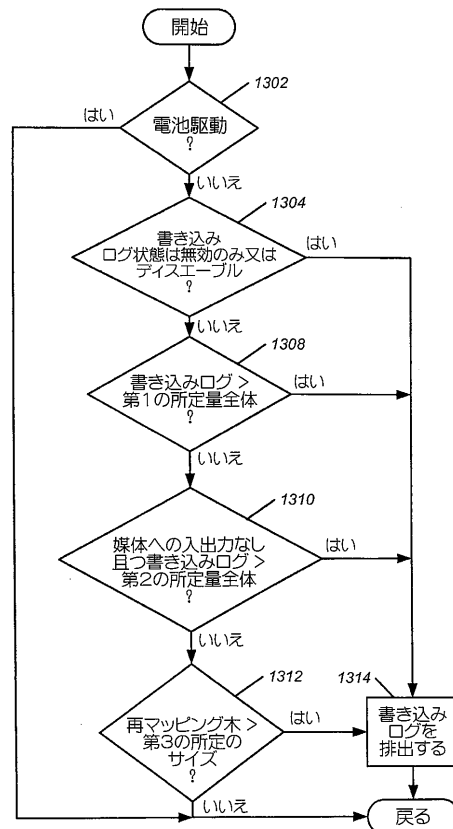
【図11】



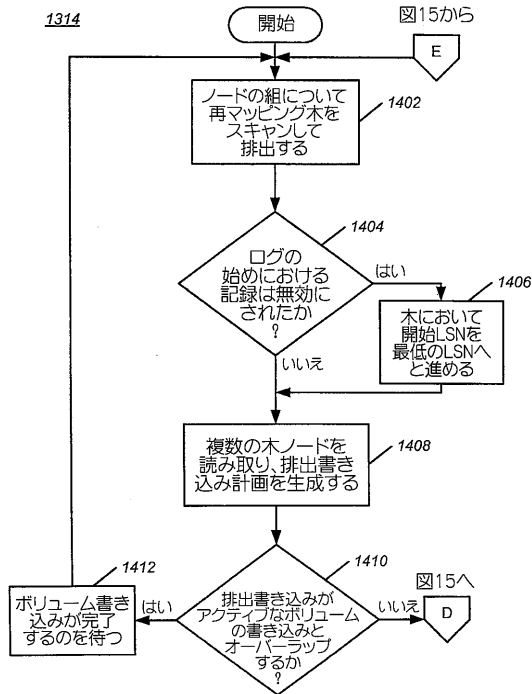
【図12】



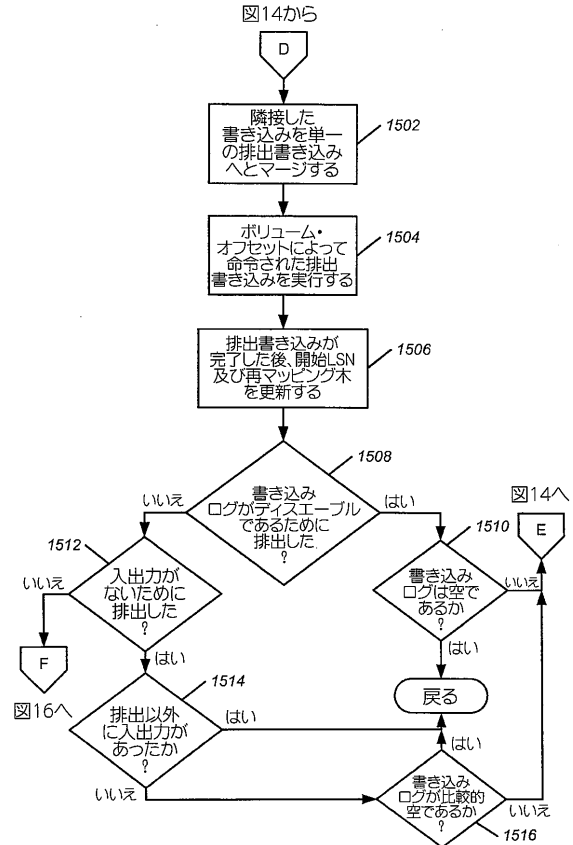
【図13】



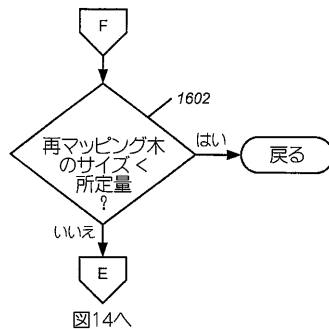
【図14】



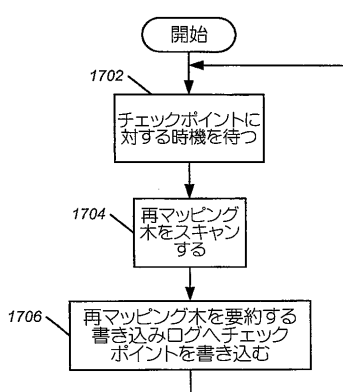
【図15】



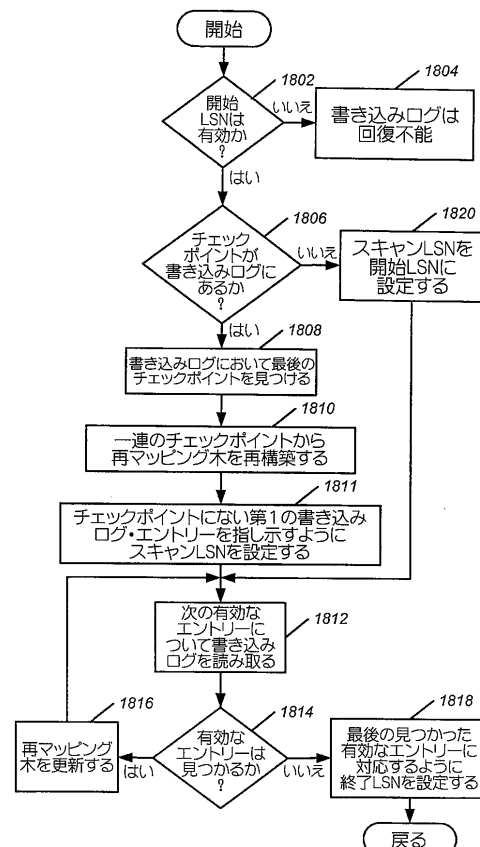
【図16】



【図17】



【図18】



フロントページの続き

- (72)発明者 コン, シ
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, ワン・マイクロソフト・ウェイ, 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・パテント
- (72)発明者 ブレンダー, スコット
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, ワン・マイクロソフト・ウェイ, 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・パテント
- (72)発明者 メーラ, カラン
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, ワン・マイクロソフト・ウェイ, 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・パテント
- (72)発明者 モス, ダレン・ジー
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, ワン・マイクロソフト・ウェイ, 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・パテント
- (72)発明者 ティプトン, ウィリアム・アール
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, ワン・マイクロソフト・ウェイ, 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・パテント
- (72)発明者 ヴェルマ, スレンドラ
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, ワン・マイクロソフト・ウェイ, 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・パテント

審査官 坂東 博司

- (56)参考文献 特開平 0 6 - 3 2 4 8 1 7 (J P , A)
特開平 1 1 - 3 1 6 6 9 9 (J P , A)
特開平 0 5 - 2 3 3 3 9 2 (J P , A)
特開 2 0 0 1 - 2 4 3 0 2 1 (J P , A)
特開 2 0 0 7 - 2 3 3 8 9 6 (J P , A)
特開 2 0 0 2 - 3 2 3 9 5 9 (J P , A)

- (58)調査した分野(Int.Cl. , D B 名)
G 0 6 F 3 / 0 6