

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6445049号
(P6445049)

(45) 発行日 平成30年12月26日 (2018.12.26)

(24) 登録日 平成30年12月7日 (2018.12.7)

| | | | | | |
|-------------------|--------------|------------------|------|-------|------|
| (51) Int. Cl. | | F I | | | |
| G06F 12/00 | 12/00 | (2006.01) | G06F | 12/00 | 531J |
| G06F 11/34 | 11/34 | (2006.01) | G06F | 12/00 | 531R |
| | | | G06F | 11/34 | 176 |

請求項の数 4 (全 18 頁)

| | | | |
|---------------|------------------------------|-----------|-----------------------------------|
| (21) 出願番号 | 特願2016-570238 (P2016-570238) | (73) 特許権者 | 000005108 |
| (86) (22) 出願日 | 平成27年1月20日 (2015.1.20) | | 株式会社日立製作所 |
| (86) 国際出願番号 | PCT/JP2015/051343 | | 東京都千代田区丸の内一丁目6番6号 |
| (87) 国際公開番号 | W02016/117022 | (74) 代理人 | 110001678 |
| (87) 国際公開日 | 平成28年7月28日 (2016.7.28) | | 特許業務法人藤央特許事務所 |
| 審査請求日 | 平成28年11月28日 (2016.11.28) | (72) 発明者 | 友田 敦 |
| | | | 東京都千代田区丸の内一丁目6番6号 株 株式会社日立製作所内 |
| | | (72) 発明者 | 磯田 有哉 |
| | | | 東京都千代田区丸の内一丁目6番6号 株 株式会社日立製作所内 |
| | | (72) 発明者 | 牛嶋 一智 |
| | | | 東京都千代田区丸の内一丁目6番6号 株 株式会社日立製作所内 |

最終頁に続く

(54) 【発明の名称】 ログの管理方法及び計算機システム

(57) 【特許請求の範囲】

【請求項1】

プロセッサとメモリとストレージ装置とを有する計算機システムで、前記プロセッサが所定の処理を実行し、前記ストレージ装置に前記処理の内容を含むログを格納するログの管理方法であって、

前記プロセッサが、前記所定の処理の内容を含むログを生成する第1のステップと、

前記プロセッサが、前記ログを前記ストレージ装置のログファイルに書き込む第2のステップと、

前記プロセッサが、前記ログを格納するログ領域の終端を決定して前記ストレージ装置のログファイルに書き込む第3のステップと、

前記プロセッサが、前記所定の処理の対象となるデータのバックアップデータを生成して、前記ストレージ装置に格納する第4のステップと、

前記プロセッサが、前記バックアップデータと前記ログファイルを読み込んで、前記ログを前記バックアップデータに適用して、前記ログファイルに書き込まれたログ領域の終端までリカバリを実行する第5のステップと、

を含み、

前記第2のステップは、

当該ログが正常であるか否かを示すログヘッダ識別子を当該ログに設定するステップと、

、

当該ログのサイズを示す情報を当該ログに設定するステップと、を含み、

前記第 5 のステップは、

前記ログヘッダ識別子を読み込んで当該ログヘッダ識別子が正常であるか否かを判定するステップと、

前記ログヘッダ識別子が不正である場合には、前記サイズを読み込んで次のログを探索するステップと、

前記ログヘッダ識別子が正常なログの場合には、前記ログを前記バックアップデータに適用してリカバリを行うステップと、

を含むことを特徴とするログの管理方法。

【請求項 2】

プロセッサとメモリとストレージ装置とを有する計算機システムで、前記プロセッサが所定の処理を実行し、前記ストレージ装置に前記処理の内容を含むログを格納するログの管理方法であって、

前記プロセッサが、前記所定の処理の内容を含む固定長のログを生成する第 1 のステップと、

前記プロセッサが、前記ログを前記ストレージ装置のログファイルに書き込む第 2 のステップと、

前記プロセッサが、前記ログを格納するログ領域の終端を決定して前記ストレージ装置のログファイルに書き込む第 3 のステップと、

前記プロセッサが、前記所定の処理の対象となるデータのバックアップデータを生成して、前記ストレージ装置に格納する第 4 のステップと、

前記プロセッサが、前記バックアップデータと前記ログファイルを読み込んで、前記ログを前記バックアップデータに適用して、前記ログファイルに書き込まれたログ領域の終端までリカバリを実行する第 5 のステップと、

を含み、

前記第 2 のステップは、

当該ログが正常であるか否かを示すログヘッダ識別子を当該ログに設定するステップを含み、

前記第 5 のステップは、

前記ログヘッダ識別子を読み込んで当該ログヘッダ識別子が正常であるか否かを判定するステップと、

前記ログヘッダ識別子が不正である場合には、予め設定されたサイズを読み込んで次のログを探索するステップと、

前記ログヘッダ識別子が正常な場合には、前記ログを前記バックアップデータに適用してリカバリを行うステップと、

を含むことを特徴とするログの管理方法。

【請求項 3】

プロセッサとメモリとストレージ装置とを有する計算機システムであって、

前記計算機システムは、

所定の処理を実行して前記処理の内容を含むログを生成し、前記ストレージ装置に前記ログを格納し、前記所定の処理の対象となるデータのバックアップデータを生成して前記ストレージ装置に格納する実行部と、

前記ログを格納するログ領域の終端を決定して前記ストレージ装置のログファイルに書き込み、リカバリの際には前記バックアップデータと前記ログファイルを読み込んで、前記ログを前記バックアップデータに適用するリカバ리를、前記ログファイルに書き込まれたログ領域の終端まで行うログ管理部と、

を有し、

前記実行部は、

当該ログが正常であるか否かを示すログヘッダ識別子と、当該ログのサイズを示す情報を当該ログに設定し、

前記ログ管理部は、

10

20

30

40

50

前記ログヘッダ識別子を読み込んで当該ログヘッダ識別子が正常であるか否かを判定し、前記ログヘッダ識別子が不正である場合には、前記サイズを読み込んで次のログを探索し、前記ログヘッダ識別子が正常なログの場合には、前記ログを前記バックアップデータに適用してリカバリを行うことを特徴とする計算機システム。

【請求項 4】

プロセッサとメモリとストレージ装置とを有する計算機システムであって、

前記計算機システムは、

所定の処理を実行して前記処理の内容を含む固定長のログを生成し、前記ストレージ装置に前記ログを格納し、前記所定の処理の対象となるデータのバックアップデータを生成して前記ストレージ装置に格納する実行部と、

前記ログを格納するログ領域の終端を決定して前記ストレージ装置のログファイルに書き込み、リカバリの際には前記バックアップデータと前記ログファイルを読み込んで、前記ログを前記バックアップデータに適用するリカバ리를、前記ログファイルに書き込まれたログ領域の終端まで行うログ管理部と、

を有し、

前記実行部は、

当該ログが正常であるか否かを示すログヘッダ識別子を当該ログに設定し、

前記ログ管理部は、

前記ログヘッダ識別子を読み込んで当該ログヘッダ識別子が正常であるか否かを判定し、前記ログヘッダ識別子が不正である場合には、予め設定されたサイズを読み込んで次のログを探索し、前記ログヘッダ識別子が正常な場合には、前記ログを前記バックアップデータに適用してリカバリを行うことを特徴とする計算機システム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、リレーショナルデータベースのマネジメントシステムにおけるトランザクション処理の永続化を目的としたログ出力の同時実行性能の向上に関する。

【背景技術】

【0002】

デバイスの発達によりサーバに搭載されるCPUに集積されるコア数や、メモリ容量が増大している。これによりリレーショナルデータベースのテーブルやインデックスといった主要なデータをメモリ上に配置するインメモリデータベースが普及してきている。

【0003】

インメモリデータベースにおいても、システムの障害時にデータベースの状態を復元するためには、ハードディスクやフラッシュメモリなどから構成されるストレージ装置に変更履歴をログとして出力しておく必要がある。

【0004】

従来のハードディスクドライブ等をログを記録する永続デバイスとして想定した計算機システムにおいては、一回のI/O出力に対するレスポンスがトランザクション処理を行うサーバの処理時間に対して非常に大きかった。このため、トランザクションごとにログのI/O出力を行うのではなく、メモリ上に確保されたログバッファ領域にログを格納した後、ログマネージャスレッドが、ログバッファ領域のログをまとめて永続デバイスに時系列的かつシーケンシャルに出力していた。なお、永続デバイスは、永続的にデータを保持する不揮発性記憶装置である。

【0005】

複数のトランザクションで1つのログファイルを共有することを効率的に行う方法として、ログファイルを固定サイズのスロットに切り分け、複数のスレッド間でこれらのスロットをログマネージャで順に予約し、ログを書き込んでいく方法が特許文献1に開示されている。

【先行技術文献】

10

20

30

40

50

【特許文献】

【0006】

【特許文献1】米国特許出願公開第2014/0208083号明細書

【発明の概要】

【発明が解決しようとする課題】

【0007】

昨今、半導体製造技術の進歩によって1つのサーバに搭載可能なコア数が増加し、複数のスレッドで並列的にトランザクションを処理するため、前述のログマネージャを介して行っていたログ出力完了の待ち合わせ処理がボトルネックになってきていた。

【0008】

一方、複数のスレッド間でのI/O処理に伴うリソースの競合に関しては、NVM(Non Volatile Memory) Expressと呼ばれる規格において、スレッドごとにI/Oキュー等を独立しておき、I/O処理がスレッドごとに独立して実行させる技術が確立されている。

【0009】

これにより、各スレッドは他のスレッドと競合することなくI/O出力を行えるようになった。しかしながら、スレッドごとに割り当てられたI/Oキューなどのリソースを活用すると、以下の問題が生じる。

【0010】

第一はログを格納する記憶装置であるログデバイスの歯抜けの発生である。トランザクション処理を行う各スレッドが、それぞれ予約したログデバイスの領域に対して独立して割り当てられたI/Oキューを利用してI/O出力を行うと、スレッドが独立して発行したI/O処理の完了順序は必ずしも保証されない。

【0011】

したがって、ある時点においてサーバに障害が発生すると、ログデバイス上に本来は連続して配置されているはずのログのうち、I/O処理が完了したスレッドのログだけが永続化され、I/O処理が完了しなかったスレッドが書き込むはずであった領域は書き込みされる前の状態(空白)のまま残ってしまう。以下では書き込みが行われなかったブランク領域のことを歯抜け領域(toothless area)と呼ぶ。

【0012】

第二は、障害回復時の処理時間の増大である。従来例に示したログマネージャがまとめてログをI/O出力していたときには、ログマネージャが複数のトランザクションのログを連続したアドレスに記録されるように整形してから一度のI/O出力で書き込みを行っていた。このため、上記の歯抜け領域は発生せず、1つのログの後ろの領域にログが記録されていない領域が最初に出現して以降の領域は、ログが書かれていない領域だった。

【0013】

しかしながら、各スレッドが独立してログを書き出す場合では、上述のように歯抜け領域が発生する場合がある。したがって、ログをもとにデータベースを障害前の状態に復旧する障害回復処理において、ログデバイスを走査する際に、一旦ログが書かれていない領域(歯抜け領域)が出現した場合でも、走査をつづけた先に別のログが記録されている可能性がある。そこで、歯抜け領域が出現しても、ログデバイス全体を走査する必要があり、回復時間を大幅に増大させてしまう、という問題があった。

【0014】

そこで本発明は、上記問題点に鑑みてなされたもので、複数のログをストレージ装置へ書き込む計算機システムで、リカバリ時にログの読み込みを高速化することを目的とする。

【課題を解決するための手段】

【0015】

本発明は、プロセッサとメモリとストレージ装置とを有する計算機システムで、前記プロセッサが所定の処理を実行し、前記ストレージ装置に前記処理の内容を含むログを格納

10

20

30

40

50

するログの管理方法であって、前記プロセッサが、前記所定の処理の内容を含むログを生成する第1のステップと、前記プロセッサが、前記ログを前記ストレージ装置のログファイルに書き込む第2のステップと、前記プロセッサが、前記ログを格納するログ領域の終端を決定して前記ストレージ装置のログファイルに書き込む第3のステップと、前記プロセッサが、前記所定の処理の対象となるデータのバックアップデータを生成して、前記ストレージ装置に格納する第4のステップと、前記プロセッサが、前記バックアップデータと前記ログファイルを読み込んで、前記ログを前記バックアップデータに適用して、前記ログファイルに書き込まれたログ領域の終端までリカバリを実行する第5のステップと、を含み、前記第2のステップは、当該ログが正常であるか否かを示すログヘッダ識別子を当該ログに設定するステップと、当該ログのサイズを示す情報を当該ログに設定するステップと、を含み、前記第5のステップは、前記ログヘッダ識別子を読み込んで当該ログヘッダ識別子が正常であるか否かを判定するステップと、前記ログヘッダ識別子が不正である場合には、前記サイズを読み込んで次のログを探索するステップと、前記ログヘッダ識別子が正常なログの場合には、前記ログを前記バックアップデータに適用してリカバリを行うステップと、を含む。

10

【発明の効果】

【0016】

本発明によれば、ログ領域の終端がログファイルに記載されているので、障害回復等のリカバリ時にストレージ装置を走査する領域の範囲が特定され、ストレージ装置の走査時間が短縮され、結果として計算機システムの障害回復時間を短縮することができる。

20

【図面の簡単な説明】

【0017】

【図1】本発明の第1の実施例を示し、計算機システムの一部を示すブロック図である。

【図2】本発明の第1の実施例を示し、データベースサーバの一部を示す機能ブロック図である。

【図3】本発明の第1の実施例を示し、ログ管理部の構成の一部を示すブロック図である。

【図4】本発明の第1の実施例を示し、トランザクション処理の一部を示すフローチャートである。

【図5】本発明の第1の実施例を示し、ログ出力処理のフローチャートである。

30

【図6】本発明の第1の実施例を示し、ログのデータ構造の一部を示す図である。

【図7】本発明の第1の実施例を示し、リカバリ処理の一部を示すフローチャートである。

【図8】本発明の第1の実施例を示し、可変長のログでブランク領域の次のログを探索する例を示す図である。

【図9】本発明の第2の実施例を示し、固定長のログの構造の一部を示す図である。

【発明を実施するための形態】

【0018】

以下、本発明の一実施形態について添付図面を用いて説明する。

【実施例1】

40

【0019】

以下、図面を参照して、一実施例を説明する。

【0020】

なお、以下の説明では、同種の要素の参照符号には同一の親番号を含んでいる。同種の要素を区別して説明する場合には、個々の要素を識別する参照符号（例えば、アルファベット）を使用し（例えばスレッド110A、110B、...）、同種の要素を区別しないで説明する場合には、要素の参照符号のうちの親番号のみ使用することがある（例えばスレッド110）。

【0021】

また、以下の説明では、「プログラム」を主語として処理を説明する場合があるが、プ

50

プログラムは、プロセッサによって実行されることで、予め定められた処理を、適宜に記憶資源（例えばメインメモリ416）及び/又は通信インターフェイスデバイス等を用いながら行うため、処理の主語がプロセッサとされてもよい。プログラムを主語として説明された処理は、プロセッサ或いはプロセッサを有する装置（例えばデータベースサーバ）が行う処理としてもよい。また、プロセッサは、処理の一部又は全部を行うハードウェア回路を含んでもよい。プログラムは、プログラムソースから各コントローラにインストールされてもよい。プログラムソースは、例えば、プログラム配布計算機又は記憶メディアであってもよい。

【0022】

図1は、本発明の計算機システムの構成の一例を示すブロック図である。データベースサーバ401に外部ストレージ装置402が、例えば通信ネットワーク403を介して接続されている。

10

【0023】

データベースサーバ401は、計算機であって、例えば、パーソナルコンピュータや、ワークステーション又はメインフレームであってよく、もしくは、これらの計算機において仮想化プログラムによって構成された仮想的な計算機であってよい。

【0024】

データベースサーバ401は、I/Oアダプタ413と、メインメモリ416と、記憶デバイス415及びそれらに接続されたプロセッサ414を有する。プロセッサ414は、例えば、マルチコアのマイクロプロセッサ又は、マイクロプロセッサと専用ハードウェア回路を含んだモジュールで構成してもよい。プロセッサ414は、メインメモリ416にロードされたコンピュータプログラムを実行する。プロセッサ414により実行されるコンピュータプログラムは、例えば、オペレーティングシステム(OS)117及びデータベース管理システム(Data Base Management System: 以下DBMS)412である。

20

【0025】

メインメモリ416は、例えば、揮発性のDRAM(Dynamic Random Access Memory)等であり、プロセッサ414によって実行されるプログラムと、プログラムが使用するデータを一時的に記憶する。なお、メインメモリ416には、不揮発性の半導体メモリを採用してもよい。

【0026】

30

記憶デバイス415は、不揮発性の記憶媒体を含み、例えば、HDD(Hard Disk Drive)又はSSD(Solid State Drive)である。記憶デバイス415は、プログラム及びプログラムが使用するデータを格納してよい。I/Oアダプタ413は、通信ネットワーク403とデータベースサーバ401を接続する。

【0027】

外部ストレージ装置402は、複数の記憶デバイスを含む記憶デバイス群443を有する装置であり、例えば、ディスクアレイ装置である。なお、外部ストレージ装置402は、複数の記憶デバイスに代えて、単一の記憶デバイスであってもよい。

【0028】

外部ストレージ装置402は、複数のログを格納したログファイル301を記憶する。外部ストレージ装置402は、データベースサーバ401からログのI/O要求を受け付ける。外部ストレージ装置402は、当該I/O要求に従いデータ(例えばログ)の読み書きを行い、読み書きの結果をデータベースサーバ401に回答する。なお、記憶デバイス群443が有する記憶デバイスは、不揮発性の記憶媒体を有するデバイスであって、例えば、HDD又はSSDである。記憶デバイス群443は、RAID(Redundant Array of Independent Disks)グループで構成してもよく、所定のRAIDレベルでデータを記憶してもよい。記憶デバイス群443の記憶空間に基づく論理的な記憶デバイス(例えば、論理ユニット、論理ボリューム、ファイルシステムボリューム)が、データベースサーバ401に提供され、当該論理的な記憶デバイス上にログファイル301が格納されてもよい。なお、本実施例において、ログファイル301は、ログを格納するログ格納領

40

50

域の一例である。

【0029】

外部ストレージ装置402は、記憶デバイス群443に加えて、I/Oアダプタ441及びこれらに接続されたストレージコントローラ442を有する。I/Oアダプタ441は、外部ストレージ装置402を通信ネットワーク403に接続し、通信ネットワーク403を介して、データベースサーバ401と接続される。通信ネットワーク403を介した通信プロトコルとしては、例えば、ファイバチャネル(FC)や、SCSI(Small Computer System Interface)、又は、TCP/IP(Transmission Control Protocol/Internet Protocol)が採用されてよい。例えば、ファイバチャネルもしくはSCSIが採用される場合、I/Oアダプタ441(及び413)は、ホストバスアダプタと呼ばれることがある。

10

【0030】

ストレージコントローラ442は、例えば、メモリ及びプロセッサを含み、データベースサーバ401からのI/O要求に応じて、ログファイル301を格納した記憶デバイス群443との間でデータの読出し、もしくは、書込みを行う。

【0031】

図2は、実施例に係るデータベースサーバ401の機能ブロック図である。

【0032】

本実施例1のDBMS412は、例えば、インメモリデータベースである。DBMS412は、テーブル112と、インデックス113とを、メインメモリ416上に配置する。

20

【0033】

さらにDBMS412はロック機構116を含んでよい。ロック機構116は2以上のスレッド110A~110Cが競合することを避けるためである。ロック機構116は、テーブル112やインデックス113をロックするモジュールである。ロック機構116は、ロックの取得済か否かを表す情報を含むことができる。例えば、ロック機構116は、ロック取得済なら値「1」でよくロック未取得なら値「0」でよい。

【0034】

DBMS412は、ログバッファ114とログ管理部115を有する。ログバッファ114は、テーブル112やインデックス113への更新履歴を含んだログを一時的に格納する。ログ管理部115は、ログファイル301およびログファイル301へのログの書き出しを管理する。なお、ログ管理部115は、ログファイル301とバックアップデータを読み込んで、ログをバックアップデータに適用してテーブル112を復元するリカバリ処理部125を含むことができる。

30

【0035】

DBMS412は、クエリ発行元からクエリを受信し、受信したクエリの実行において、1又は複数のトランザクションを実行する。具体的には、DBMS412は、クエリ受付部421、クエリ実行プラン生成部422及びクエリ実行部424を有する。

【0036】

クエリ受付部421は、クエリ発行元が発行するクエリを受け付ける。クエリは、例えば、構造化問合せ言語(SQL、Structured Query Language)によって記述される。1つのクエリで複数のトランザクションが記述されていてもよいし、複数のクエリで複数のトランザクションが記述されてもよい。

40

【0037】

また、クエリ発行元は、DBMS412の内部のコンピュータプログラムであってよいし、DBMS412外部のコンピュータプログラムであってよい。例えば、外部のコンピュータプログラムは、データベースサーバ401内で実行されるコンピュータプログラム(例えばアプリケーションプログラム)であってよいし、データベースサーバ401に接続されたクライアント計算機等の装置で実行されるコンピュータプログラム(例えば、アプリケーションプログラム)であってよい。

50

【 0 0 3 8 】

クエリ実行プラン生成部 4 2 2 は、クエリ受付部 4 2 1 が受け付けたクエリから、当該クエリを実行するために必要な 1 以上のデータベースオペレーションを含むクエリ実行プランを生成する。

【 0 0 3 9 】

クエリ実行プランは、例えば、1 以上のデータベースオペレーションと、データベースオペレーションの実行順序の関係を含む情報であり、クエリ実行プラン情報 4 2 3 として格納される。クエリ実行プラン情報 4 2 3 は、データベースオペレーションをノード、データベースオペレーションの実行順序の関係をエッジとする木構造で表されることがある。

10

【 0 0 4 0 】

クエリ実行部 4 2 4 は、クエリ実行プラン生成部 4 2 2 が生成したクエリ実行プランに従って、クエリ受付部 4 2 1 が受け付けたクエリを実行し、クエリの実行結果をクエリ発行元に応答する。

【 0 0 4 1 】

この際、クエリ実行部 4 2 4 は、データベースオペレーションの実行に必要なデータの読み出し要求（参照要求）を発行し、当該読み出し要求に従いテーブル 1 1 2 からデータを読み出す。クエリ実行部 4 2 4 は、読み出したデータを使用して、クエリに応じたデータベースオペレーションを実行してデータを算出し、読み出し元レコードのデータを算出後のデータに更新する書き込み要求を発行する。

20

【 0 0 4 2 】

クエリ実行部 4 2 4 は、データベースオペレーションを、1 以上のスレッド 1 1 0 A ~ 1 1 0 C を実行することにより行う。なお、DBMS 4 1 2 において、複数のスレッド 1 1 0 A ~ 1 1 0 C が並行して実行される。このため、プロセッサ 4 1 4 は、複数のコアを有する。複数のコアは、1 又は複数の CPU に存在する。スレッド 1 1 0 は、タスクと呼ばれてもよい。スレッド 1 1 0 の実装としては、例えば、OS 1 1 7 が実現するプロセスやカーネルスレッド等のほか、ライブラリ等が実現するユーザスレッドを用いてよい。1 つのスレッド 1 1 0 により、1 以上のデータベースオペレーションに対応した 1 つのトランザクションが実行されてよい。以下、クエリ実行部 4 2 4 がスレッド 1 1 0 を実行することにより行われる処理の主語を、スレッド 1 1 0、とすることがある。

30

【 0 0 4 3 】

クエリ実行部 4 2 4（スレッド 1 1 0）は、トランザクションの実行と、トランザクションの実行結果（または処理内容）を含むログを生成する。そして、クエリ実行部 4 2 4 の各スレッド 1 1 0 は、外部ストレージ装置 4 0 2 内のログファイル 3 0 1 にログを書き込むために、外部ストレージ装置 4 0 2 に対する I/O 要求を OS 1 1 7 に発行する。OS 1 1 7 は、当該 I/O 要求を受け付け、外部ストレージ装置 4 0 2 へ I/O 要求を発行する。

【 0 0 4 4 】

I/O アダプタ 4 1 3 には、複数の I/O キュー 2 0 1（2 0 1 A ~ 2 0 1 C）が設定される。トランザクションの処理において、スレッド 1 1 0 が、ログの書き込みのための I/O 要求を外部ストレージ装置 4 0 2 へ発行するが、I/O キュー 2 0 1 には、当該 I/O 要求が格納される。具体的には、I/O 要求は、OS 1 1 7 により I/O キュー 2 0 1 に格納される。

40

【 0 0 4 5 】

外部ストレージ装置 4 0 2 が、ログファイル 3 0 1 を記憶する。ログファイル 3 0 1 に、I/O 要求の書き込み対象のログが記録される。

【 0 0 4 6 】

本実施例 1 では、クエリ実行部のスレッド 1 1 0 と、I/O キュー 2 0 1 が、1 : 1 対応している。つまり、スレッド 1 1 0 A ~ 1 1 0 C 毎に、1 つの I/O キュー 2 0 1 A ~ 2 0 1 C が設定される。具体的には、スレッド 1 1 0 A に I/O キュー 2 0 1 A が関連付

50

けられている。例えば、スレッド110Aが、テーブル112のレコードを更新したことを表すログのI/O要求を、ログファイル301に対して発行するようになっている。発行されたI/O要求は、ログバッファ114を経由して、OS117に送られる。OS117が、ログファイル301に対するI/O要求を受けて、当該I/O要求を、スレッド110Aに対応するI/Oキュー201Aに格納する。I/Oキュー201Aに格納されたI/O要求は、OS117によりI/Oキュー201から外部ストレージ装置402に送られる。外部ストレージ装置402は、当該I/O要求の書込み対象データであるログを、ログファイル301に書き込む。

【0047】

図2に示すDBMS412の構成は一例に過ぎない。例えば、或る構成要素は複数の構成要素に分割されていてもよく、複数の構成要素が1つの構成要素に統合されていてもよい。

【0048】

図3は、ログ管理部115の構成の一例を示すブロック図である。ログ管理部115は、ロック機構121と、ログファイルアドレス122と、ログ領域終端アドレス123と、ログ領域追加フラグ124と、リカバリ処理部125とを有する。

【0049】

ロック機構121も、図2に示したロック機構116と同様に、ログ管理部115のロックの取得済か否かを表すデータでよい。例えば、ロック機構121は、ロック取得済なら値「1」でよくロック未取得なら値「0」でよい。ログファイルアドレス122、ログ領域終端アドレス123、ログ領域追加フラグについて書き込みまたは読み出しを行うときには、ロック機構121のロックを取得していなければならない。

【0050】

ログファイルアドレス122は、ログファイル301におけるログの書込み先アドレスである。ログファイルアドレス122が表すアドレス(値)は、ログファイル301にログを書き出す度に出力するログのサイズ分だけ加算される。なお、ログファイルアドレス122や後述のログ領域終端アドレス123は、外部ストレージ装置402上でのログの格納領域の終端を示す値であり、例えばLBA(Logical Block Address)等を採用することができる。

【0051】

ログ領域終端アドレス123はログファイルアドレス122の上限値であり、当該上限値を超えた場所にログを書くことはできない。また、リカバリ処理においては、ログ領域終端アドレス123にリカバリ処理におけるログファイルの走査範囲の上限がセットされる。リカバリ処理において、ログ領域終端アドレス123を超えて走査を行わない。

【0052】

ログ領域追加フラグ124は、トランザクションを処理するスレッドが、ログ出力する領域で領域の追加処理を実行中の場合には、本フラグがセットされる。例えば、本フラグは追加中なら「1」、そうでない場合には「0」でよい。

【0053】

リカバリ処理部125は、ログファイル301のログを適用してテーブル112を復旧するリカバリ処理を実行する。リカバリ処理部125は、図示しない管理装置などから所定の指令を受け付けたときに処理を開始する。

【0054】

図4は、トランザクション処理のフローチャートである。なお、以下の説明では、1つのトランザクションを例に取り説明する。スレッド110Aが、トランザクションAを開始すると、トランザクションAに対応した指示(クエリ中の指示)に基づいて、参照及び更新セットの生成を行う(S301)。

【0055】

参照及び更新セットは、レコードの参照(テーブル112の読み出し要求)とレコードの更新(テーブル112、インデックス113の書込み要求)とのセットである。参照及び

10

20

30

40

50

更新セットは、テーブル112と、インデックス113を更新するための要求セットであるが、ステップS301の時点では、テーブル112、インデックス113の変更は行われず、トランザクションAに対応したローカルメモリ領域（メインメモリ461上に確保された領域（図示せず））に参照及び更新セットが保持される。

【0056】

次に、スレッド110Aが、コミット判定を行う（S302）。コミット判定は、例えば、トランザクションAが参照及び更新セットに基づいて、テーブル112、インデックス113に対して行う変更が他のトランザクションとの整合性を保っているか否かをデータベースのアイソレーションレベル（またはトランザクション分離レベル）に応じて行われる。

10

【0057】

コミット判定がNG（処理失敗など）の場合（S303：No）、スレッド110Aは、アボート処理を行い（S307）、アボート完了通知を出力して、トランザクションを終了する。

【0058】

コミット判定がOK（処理完了など）の場合（S303：Yes）、スレッド110Aは、ログ出力処理を実行する（S304）。ログ出力処理は、後述するように、所定の処理（トランザクション）が完了する度に処理の内容を含むログをログファイル301へ書き込む処理である。

【0059】

20

次に、スレッド110Aは、参照及び更新セットに基づいてテーブル112、インデックス113をそれぞれ更新し（S305）、コミット完了通知を出して（S306）、トランザクションを終了する。

【0060】

上記処理により、スレッド110Aは、トランザクション処理が完了すればコミット完了通知を出力し、トランザクション処理が失敗すればアボート完了通知を出力する。

【0061】

図5は、ログ出力処理（図4のS304）の一例を示すフローチャートである。まず、トランザクションを実行しているスレッド110Aは、ログ管理部115のロック機構121を取得する（S501）。

30

【0062】

次に、スレッド110Aは、ログ領域終端アドレス123とログファイルアドレス122の差が一定値以下であり、さらにログ領域追加フラグ124がセットされていないことを判定する（S502）。ここで、一定値は、ログ領域終端アドレス123の更新に要する時間内にログ領域を使い切らない十分な容量であり、予め設定されてDBMS102で保持しているものとする。

【0063】

ステップS502の判定結果がYesの場合、スレッド110Aは、ログ領域を拡張する準備を行う。つまり、スレッド110Aは、ログ領域追加フラグ124をセットし（S503）、ログバッファ114に生成した当該トランザクションのログに、ログ領域拡張ログ（またはログ領域拡張情報）を追加する（S504）。

40

【0064】

一方、ステップS502の判定結果がNoの場合には、スレッド110Aは、ログファイルアドレス122を取得し、ログバッファ114に生成したログ（すなわち、書き込み予定のログ）のサイズをログファイルアドレスに加算し（S505）、ログ管理部115のロック機構121を解放する（S506）。

【0065】

スレッド110Aは、ログバッファ114に用意したログの書き込み要求（ログ管理部115から取得したログファイルアドレス122を指定した書き込み要求）を発行する（S507）。スレッド110Aは、外部ストレージ装置402からI/Oアダプタ413を

50

介して書き込み完了通知を受信した場合に書き込み処理を完了する (S 5 0 8)。

【 0 0 6 6 】

上記の処理により、スレッド 1 1 0 A によってログファイル 3 0 1 には新たなログが書き込まれ、ログファイル 3 0 1 のログファイルアドレス 1 2 2 が更新される。

【 0 0 6 7 】

書き込みが完了すると、スレッド 1 1 0 A は、当該トランザクションのログにログ領域拡張ログを追加したか否かを判定する (S 5 0 9)。このログ領域拡張ログは、ステップ S 5 0 4 で追加される情報である。

【 0 0 6 8 】

ステップ S 5 0 9 の判定結果が、 Y e s の場合、スレッド 1 1 0 A は、ログ領域の拡張を行う。まず、スレッド 1 1 0 A は、ログ管理部 1 1 5 からロック機構 1 2 1 を取得する (S 5 1 0)。

【 0 0 6 9 】

次に、スレッド 1 1 0 A は、上記ステップ S 5 0 4 でログファイル 3 0 1 に記載したログ領域拡張情報に対応する値を加算してログ領域終端アドレス 1 2 3 にセットする (S 5 1 1)。すなわち、スレッド 1 1 0 A は、ログ領域の終端アドレスに予め設定したサイズを加算してログ領域を拡張する。

【 0 0 7 0 】

そして、スレッド 1 1 0 A は、ログファイル 3 0 1 のログ領域終端アドレス 1 2 3 を更新する。さらに、スレッド 1 1 0 A は、ログ領域追加フラグ 1 2 4 をクリアする (S 5 1 3)。上記ステップ S 5 0 9 ~ S 5 1 3 の処理により、ログファイル 3 0 1 のログ領域が拡張される。

【 0 0 7 1 】

一方、ステップ S 5 0 9 において、判定結果が N o の場合には、スレッド 1 1 0 A は、そのままログ出力処理を終了する。

【 0 0 7 2 】

上記処理により、スレッド 1 1 0 A は、ログをログファイル 3 0 1 へ書き込む際に、ログ領域が不足する恐れがある場合にはログ領域を拡張することができる。

【 0 0 7 3 】

図 6 は、ログのデータ構造の一例を示す図である。1 つのトランザクションに対して、少なくとも 1 つのログ 3 0 が生成される。

【 0 0 7 4 】

ログファイル 3 0 1 に格納されるログ 3 0 は、ログヘッダ識別子 3 3 とログサイズ 3 4 から構成されるログヘッダ 3 1 と、データベースへの変更履歴およびログ領域拡張ログを格納したログ本体 3 2 とから構成される。なお、ログサイズ 3 4 は、ログ 3 0 が可変長の場合にログ 3 0 のサイズを示す値である。

【 0 0 7 5 】

ログヘッダ識別子 3 3 は、ログヘッダ 3 1 の先頭に格納されており、値を所定の方法で検証することにより、当該アドレスからログ 3 0 が正常に開始していることが判定できるものである。例えばもっとも容易な実装方法として、ログヘッダ 3 1 の開始アドレスの値 (またはアドレスのハッシュ値) が格納されている場合、当該アドレスから有効なログヘッダ 3 1 が開始していると判定することができる。なお、偶然に記録されていたビット列が、このような判定で正常と誤判定されることがないように、例えば判定に使用するビット長を長く、アドレスとハッシュ値を併記するなどしてもよい。

【 0 0 7 6 】

リカバリ処理部 1 2 5 は、ログヘッダ 3 1 からログヘッダ識別子 3 3 を読み込んで、検証することにより当該ログ 3 0 が有効であるか否かを判定することができる。

【 0 0 7 7 】

なお、ログファイル 3 0 1 は、図示のログ 3 0 を複数格納するファイルである。ログファイル 3 0 1 には、上記ログ 3 0 に加えてログファイルアドレス 1 2 2 とログ領域終端ア

10

20

30

40

50

ドレス123が格納される。また、ログファイルアドレス122とログ領域終端アドレス123は、ログファイル301内の所定の領域（例えば、先頭）に格納される。

【0078】

図7は、リカバリ処理のフローチャートである。リカバリ処理では、DBMS102が外部ストレージ装置402へバックアップしたある時点のデータベース（テーブル112）と、バックアップを生成した時点からの差分を逐次的に記録したログ30とから最新のデータベース（テーブル112）を復元する。

【0079】

本実施例1では、DBMS102が所定のタイミングとなる度に外部ストレージ装置402へデータベースのバックアップを行う。また、DBMS102は、テーブル112の更新、追加、削除を行う度にログ30を生成してログファイル301へ書き込む。そして、DBMS102は、所定の指令を受け付けるとリカバリ処理部125でリカバリ処理を開始する。なお、所定のタイミングは、所定時間の経過や、テーブル112が更新されたときなどである。

【0080】

リカバリ処理部125は、まず、外部ストレージ装置402へバックアップされたデータベース（テーブル112）をメインメモリ416上にロードする（S701）。次に、リカバリ処理部125は、ログファイル301に記録されたログ領域終端アドレス123を、ログ管理部115の該当領域にセットする（S702）。また、リカバリ処理部125は、ログファイルアドレス122にログファイル301の先頭アドレスなど、リカバリを開始するアドレスを設定する。なお、リカバリを開始するアドレスは、データベースサーバ401の管理者などが設定しても良く、バックアップデータの生成時刻に対応するログ30のアドレスなどを指定することができる。

【0081】

次に、リカバリ処理部125は、ログファイル301のログファイルアドレス122がログ領域終端アドレス123以下であるか否かを判定する（S703）。

【0082】

ステップS703の判定結果がNOの場合、リカバリ処理部125は、適用すべきログ30はすべてデータベースに適用されているため、リカバリ処理を終了する。

【0083】

一方、ステップS703の判定結果がYESの場合、リカバリ処理部125は、ログファイルアドレス122が指し示すアドレスから、ログヘッダ31を読み出し（S705）、ログヘッダ識別子33が正常であるか否かを判定する（S706）。

【0084】

ログヘッダ識別子33の正常化否かの判定は、上述のようにログヘッダ31の開始アドレスと、ログヘッダ識別子33の値が一致した場合などで、リカバリ処理部125は、ログヘッダ31が正常に記録されていると判定することができる。なお、本実施例1では、ログヘッダ31が正常であれば、ログ本体32も正常に記録されていると判定する。

【0085】

一方、ステップS706の判定結果が正常でない（NO）場合、リカバリ処理部125は、ログファイルアドレス122に予め設定した値を加算して、ステップS703に戻って上記処理を繰り返す。

【0086】

ここで、ログファイルアドレス122に加算する予め設定した値とは、例えば、ログ30のサイズを指している。例えばログサイズが常に32バイトの倍数の固定長である場合には、ログファイルアドレス122に32バイト加算する。これは、正常なログを失った場合に、次の有効なログ30を探索するため、次のログ30の開始アドレスの可能性のあるアドレスを順に走査するためである。

【0087】

一方、ログ30のサイズが可変長の場合、ログファイルアドレス122に加算する所定

10

20

30

40

50

数は、ログヘッダ識別子 3 3 のサイズとすることができる。ログヘッダ識別子 3 3 は固定長で指定され、例えば、4 バイトで記述される。この場合、図 8 で示すように、ログヘッダ識別子 3 3 のサイズをログファイルアドレス 1 2 2 に順次加算することで、リカバリ処理部 1 2 5 は、ブランク領域が存在する場合でも、次に出現するログヘッダ識別子 3 3 (ログヘッダ 3 1) を探索することができる。なお、図 8 は、ログ本体 3 2 が可変長の場合に、ブランク領域の次のログヘッダ 3 1 を探索する例を示す図である。

【 0 0 8 8 】

ステップ S 7 0 6 の判定結果が正常 (Y E S) である場合、リカバリ処理部 1 2 5 は、ログヘッダ 3 1 に格納されたログサイズ 3 4 を読み出し (S 7 0 7)、ログファイルアドレス 1 2 2 に、当該ログサイズ 3 4 を加算する (S 7 0 8)。さらに、リカバリ処理部 1 2 5 は、当該加算されたログファイルアドレス 1 2 2 までのログ本体 3 2 を読み込んで、ログ本体 3 2 の内容をバックアップデータに適用してデータベース (テーブル 1 1 2) を復旧する (S 7 0 9)。その後、リカバリ処理部 1 2 5 は、ログファイルアドレス 1 2 2 を次のログ 3 0 に設定してからステップ S 7 0 3 に戻ってログ領域終端アドレス 1 2 3 まで上記処理を繰り返す。

【 0 0 8 9 】

以上の処理により、ログファイル 3 0 1 中のログ 3 0 に歯抜け領域 (ブランク領域) が存在する場合でも、ログヘッダ 3 1 単位で探索することで、次のログ 3 0 を取得することができる。

【 0 0 9 0 】

本実施例 1 によれば、データベースサーバ 4 0 1 で稼働する DBMS 1 0 2 のスレッド 1 1 0 A ~ 1 1 0 C が、それぞれの I / O キュー 2 0 1 A ~ 2 0 1 C へログを並列的に格納する。その後、I / O アダプタ 4 1 3 からネットワーク 4 0 3 を介して外部ストレージ装置 4 0 2 のログファイル 3 0 1 へログを書き込む。スレッド 1 1 0 A ~ 1 1 0 C は、ログファイル 3 0 1 で自身のログ 3 0 を書き込む領域のアドレスよりも小さいアドレスの領域 (直前のログ 3 0) へのログ書き出しが完了しなくても、自身のログ領域へログ 3 0 を書き出すことができる。

【 0 0 9 1 】

これにより、前記従来例に示したログマネージャスレッドによるトランザクション間のログ書き出し処理の待ち合わせ処理時間が削減される。したがって、トランザクションあたりのプロセッサ処理時間が削減され、計算機システム全体の処理性能が向上する。

【 0 0 9 2 】

そして、DBMS 1 0 2 のスレッド 1 1 0 は、ログ 3 0 を格納する外部ストレージ装置 4 0 2 のログファイル 3 0 1 に、障害回復時 (リカバリ時) にログファイル 3 0 1 を走査する領域の範囲 (ログ領域終端アドレス 1 2 3) を記録する。

【 0 0 9 3 】

これにより、リカバリ処理の際には、ログファイル 3 0 1 の走査時間を短縮し、障害回復時間を短縮することができる。各スレッド 1 1 0 は、ログ領域終端アドレス 1 2 3 とログファイルアドレス 1 2 2 の差が所定値以下になると、新たなログ領域を確保するためにログ領域終端アドレス 1 2 3 を所定量だけ増大させる。すなわち、ログ 3 0 の追加に応じて、徐々にログ領域を拡張することで、リカバリ処理の際の走査範囲を抑制することができるのである。

【 0 0 9 4 】

したがって、ログファイル 3 0 1 の領域は、ログ 3 0 の追加に応じて徐々に増大するので、リカバリ処理の際に探索するブランク領域の範囲を低減でき、障害回復までの時間を短縮することができる。

【 0 0 9 5 】

また、ログ 3 0 が可変長の場合には、各ログ 3 0 に固定長のログヘッダ識別子 3 3 を設定する。データベースサーバ 4 0 1 の障害によって、1 以上のスレッド 1 1 0 でログ 3 0 の書き込みが完了しないブランク領域が発生しても、リカバリ処理の際には、ログヘッダ

10

20

30

40

50

識別子 3 3 サイズ単位で、ログファイルアドレス 1 2 2 を加算することで、リカバリ処理部 1 2 5 は次の有効なログヘッダ識別子 3 3 を検出することが可能となる。

【 0 0 9 6 】

なお、上記実施例 1 ではデータベースサーバ 4 0 1 の I / O アダプタ 4 1 3 にスレッド 1 1 0 毎の I / O キュー 2 0 1 を設けた例を示したが、これに限定されるものではなく、外部ストレージ装置 4 0 2 に I / O キュー 2 0 1 A ~ 2 0 1 C を設けても良い。この場合、例えば、I / O アダプタ 4 4 1 に I / O キュー 2 0 1 A ~ 2 0 1 C を設けることができる。

【 0 0 9 7 】

また、上記実施例 1 では、各スレッド 1 1 0 は、所定の条件を満たすとログ領域終端アドレス 1 2 3 を変更して、ログ領域を拡張する例を示したが、これに限定されるものではない。例えば、スレッド 1 1 0 がログ 3 0 を書き込む度に、書き込んだログ 3 0 のサイズをログ領域終端アドレス 1 2 3 に加算し、ログファイル 3 0 1 を更新しても良い。

10

【 0 0 9 8 】

すなわち、スレッド 1 1 0 が、ログ領域へログを書き込む度にログ領域の終端を決定してログ出力を行っても良いし、ログ領域を拡張する際にログ領域の終端を決定してログ出力を行っても良い。

【実施例 2】

【 0 0 9 9 】

図 9 は、第 2 の実施例を示し、固定長のログ 3 0 の構造の一例を示す図である。本実施例 2 では、前記実施例 1 のログ 3 0 のサイズを固定長としたもので、その他の構成は前記実施例 1 と同様である。

20

【 0 1 0 0 】

固定長のログ 3 0 は、前記実施例 1 と同様のログヘッダ識別子 3 3 と、ログ本体 3 2 から構成される。

【 0 1 0 1 】

リカバリ処理の際、リカバリ処理部 1 2 5 は、ひとつのログ 3 0 の読み込みとログ適用が完了すると、ログファイルアドレス 1 2 2 に所定値を加算して、次のログ 3 0 の開始アドレスを算出する。リカバリ処理部 1 2 5 は、次のログ領域のログヘッダ識別子 3 3 を読み込んで、有効か否かを判定する。

30

【 0 1 0 2 】

ログヘッダ識別子 3 3 が有効でない場合、リカバリ処理部 1 2 5 は、ログファイルアドレス 1 2 2 に所定値を加算して、次のログ 3 0 の開始アドレスを算出して上記処理をログ領域終端アドレス 1 2 3 まで返す。

【 0 1 0 3 】

ログ 3 0 のサイズを固定長とすることで、リカバリ処理の際に有効なログ 3 0 を探索する処理は可変長の場合に比して高速になる。可変長の場合は、図 8 で示したように、ブランク領域内でログヘッダ識別子 3 3 を探索する必要がある。これに対して、固定長の場合では、図 9 で示すように図中ログ 4 の次のログ領域がブランク領域なので、リカバリ処理部 1 2 5 は、現在のログファイルアドレス 1 2 2 に所定値を加算することで、次のログ 6

40

【 0 1 0 4 】

以上のように、本実施例 2 では、ログ 3 0 のサイズを固定長とすることで、リカバリ処理の際にはブランク領域であってもログ領域単位で処理を行うことが可能となって、障害回復までの時間を短縮することが可能となるのである。

【 0 1 0 5 】

<まとめ>

なお、本発明は上記した実施例に限定されるものではなく、様々な変形例が含まれる。例えば、上記した実施例は本発明を分かりやすく説明するために詳細に記載したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。また、ある実

50

施例の構成の一部を他の実施例の構成に置き換えることが可能であり、また、ある実施例の構成に他の実施例の構成を加えることも可能である。また、各実施例の構成の一部について、他の構成の追加、削除、又は置換のいずれもが、単独で、又は組み合わせでも適用可能である。

【0106】

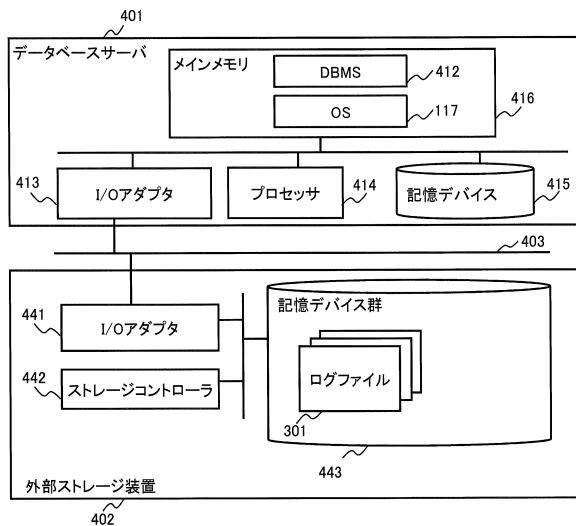
また、上記の各構成、機能、処理部、及び処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。また、上記の各構成、及び機能等は、プロセッサがそれぞれの機能を実現するプログラムを解釈し、実行することによりソフトウェアで実現してもよい。各機能を実現するプログラム、テーブル、ファイル等の情報は、メモリや、ハードディスク、SSD (Solid State Drive) 等の記録装置、または、ICカード、SDカード、DVD等の記録媒体に置くことができる。

10

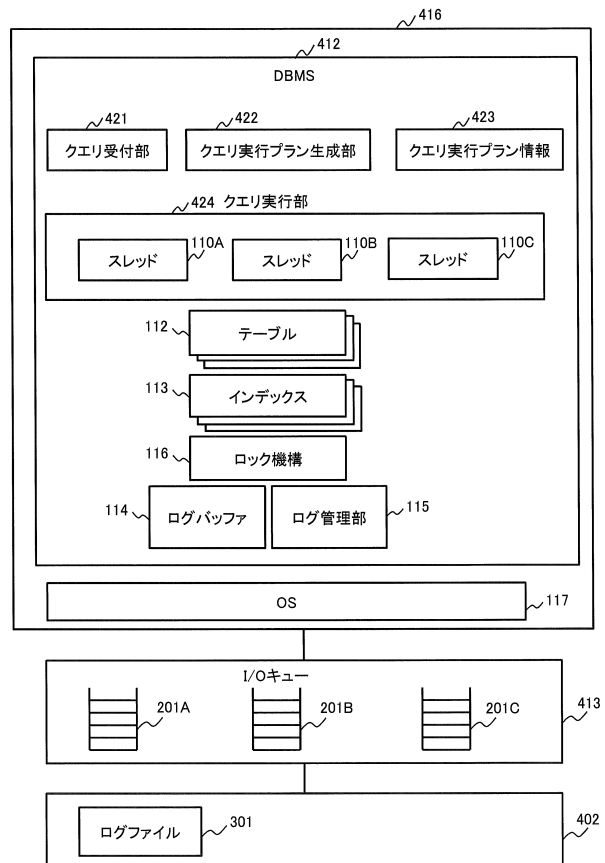
【0107】

また、制御線や情報線は説明上必要と考えられるものを示しており、製品上必ずしも全ての制御線や情報線を示しているとは限らない。実際には殆ど全ての構成が相互に接続されていると考えてもよい。

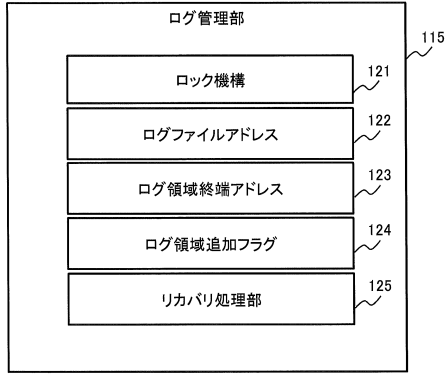
【図1】



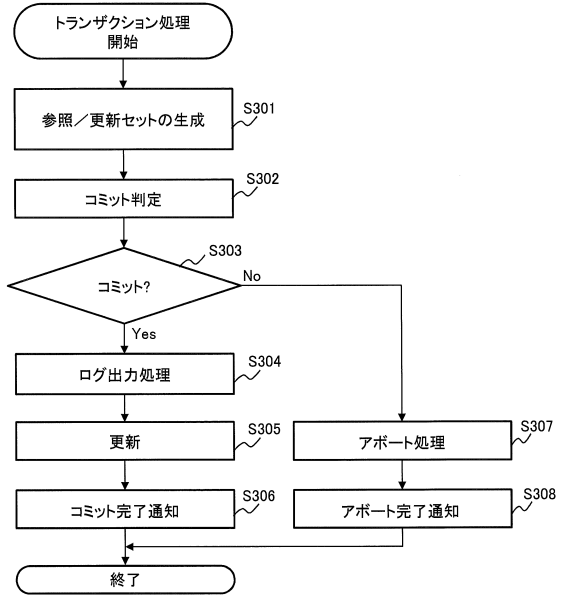
【図2】



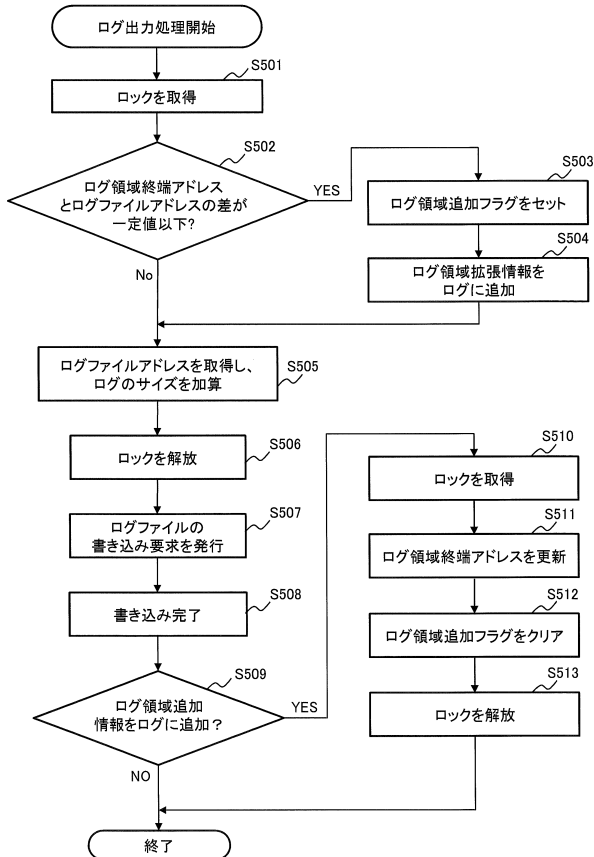
【図3】



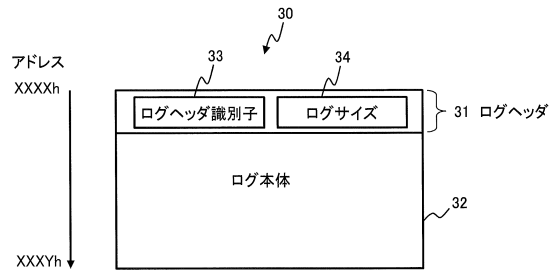
【図4】



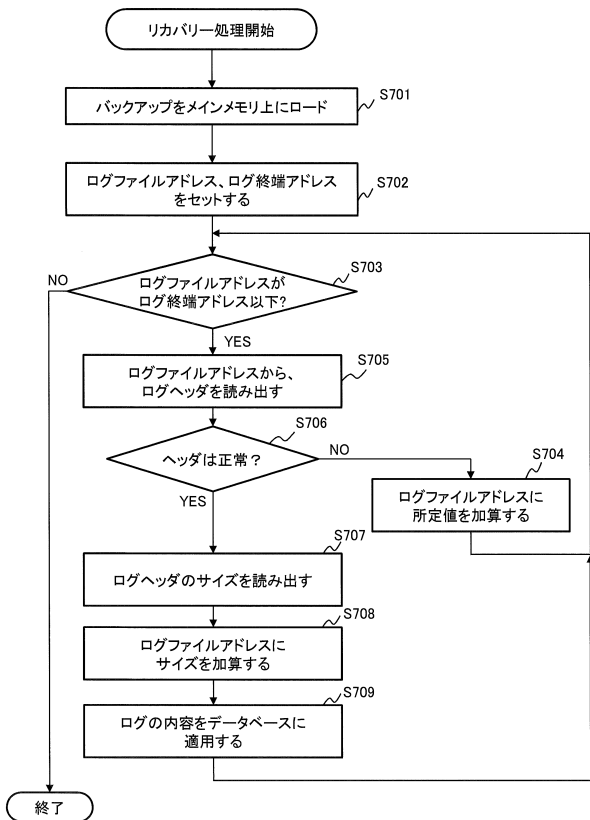
【図5】



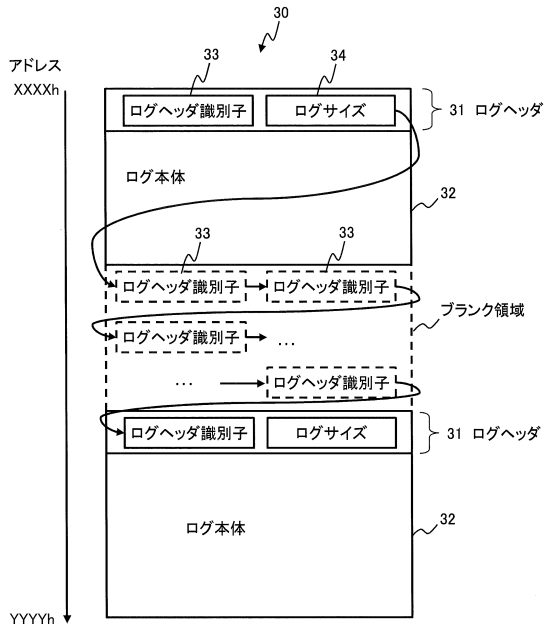
【図6】



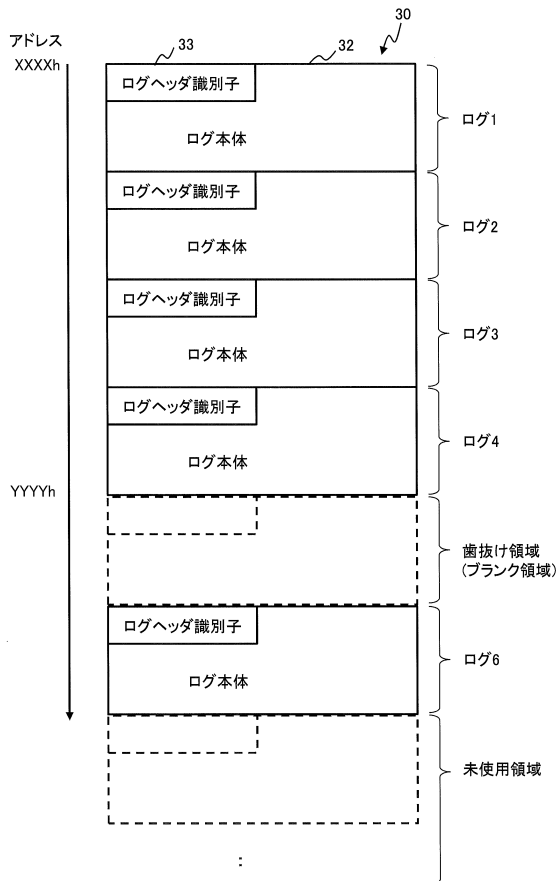
【図7】



【図8】



【図9】



フロントページの続き

審査官 原 秀人

- (56)参考文献 特開平07-248949(JP,A)
特開2006-268531(JP,A)
特開平10-247157(JP,A)
ジム グレイ, トランザクション処理 [下], 日本, 日経BP社, 2001年10月29日,
第1版, pp. 593(3)-616(26)

- (58)調査した分野(Int.Cl., DB名)
G06F 12/00
G06F 11/34