

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4126843号  
(P4126843)

(45) 発行日 平成20年7月30日(2008.7.30)

(24) 登録日 平成20年5月23日(2008.5.23)

(51) Int. Cl.	F I
<b>G06F 12/00 (2006.01)</b>	G06F 12/00 505
<b>G06F 17/30 (2006.01)</b>	G06F 12/00 518A
	G06F 17/30 414A

請求項の数 2 (全 16 頁)

(21) 出願番号	特願2000-101211 (P2000-101211)	(73) 特許権者	000005108
(22) 出願日	平成12年3月31日 (2000.3.31)		株式会社日立製作所
(65) 公開番号	特開2001-282599 (P2001-282599A)		東京都千代田区丸の内一丁目6番6号
(43) 公開日	平成13年10月12日 (2001.10.12)	(74) 代理人	100100310
審査請求日	平成16年8月24日 (2004.8.24)		弁理士 井上 学
		(72) 発明者	原 憲宏
			神奈川県川崎市幸区鹿島田890番地 株式会社日立製作所 ビジネスソリューション開発本部内
		審査官	上嶋 裕樹
		(56) 参考文献	特開平04-112240 (JP, A)
			特開平08-087511 (JP, A)
			特開昭61-184635 (JP, A)
			最終頁に続く

(54) 【発明の名称】 データ管理方法および装置並びにデータ管理プログラムを格納した記録媒体

(57) 【特許請求の範囲】

【請求項1】

一つのキーが、該キーを含むデータに関する識別情報をデータポイントとして管理し、キー値に基づいて上記キーに関連するデータポイントを見つけるためにB木インデクスを用いたデータ管理装置において、

上記一つのキーと同一のキーを含む複数のデータに関する複数のデータポイントを格納するデータ構造であり、そのキーを含むインデクスエントリが格納されるリーフノードとは異なるデータ構造であり、そのデータ構造はそのキーを含むインデクスエントリからポイントされ、前記同一キーを含む複数のデータに関する複数のデータポイントへのアクセスの際に排他を取得することは無い重複データ構造を有する記憶手段と、

あるキーに関連するデータポイントを追加もしくは削除する際に、そのキーに関連する複数のデータポイントを格納する前記重複データ構造に対し、そのデータポイントを追加もしくは削除するための情報としてテーブルデータ関連情報と前記重複データ構造へのポイントと前記キーを含む第1のログレコードを取得する手段と、

リーフノードに格納されたそのキーを含むインデクスエントリを更新する手段と、

前記重複データ構造へのそのデータポイントの追加もしくは削除が行われたことを示す情報を含む第2のログレコードを取得する手段と、

前記重複データ構造に、そのデータポイントを追加もしくは削除する手段と、

上記データポイントの追加もしくは削除を伴う処理結果を取り消す際に、

上記取得したログレコードを受け取り、受け取ったログレコードが前記第2のログレコ

ードであった場合には、その前記第2のログレコードを保持し、受け取ったログレコードが前記第1のログレコードであった場合には、第2のログレコードを保持しているかを確認し、第2のログレコードを保持している場合には、前記重複データ構造へのそのデータポインタの追加もしくは削除が完了したと判断する手段と、

上記判断において受け取ったログレコードが前記第1のログレコードであった場合にもかかわらず第2のログレコードを保持していないことを確認することにより前記重複データ構造へのそのポインタの追加もしくは削除が完了していないと判断した場合、第1のログレコードに含まれる前記重複データ構造へのポインタが示す前記重複データ構造をアクセスし、第1のログレコードに含まれる前記テーブルデータ情報を追加もしくは削除する手段を有することを特徴とするデータ管理装置。

10

【請求項2】

一つのキーが、該キーを含むデータに関する識別情報をデータポインタとして管理し、キー値に基づいて上記キーに関連するデータポインタを見つけるためにB木インデクスを用いたデータ管理プログラムを記録した計算機読み取り可能な記録媒体において、

上記一つのキーと同一のキーを含む複数のデータに関する複数のデータポインタを格納するデータ構造であり、

そのキーを含むインデクスエントリが格納されるリーフノードとは異なるデータ構造であり、そのデータ構造はそのキーを含むインデクスエントリからポイントされ、前記同一キーを含む複数のデータに関する複数のデータポインタへのアクセスの際に排他を取得することは無い重複データ構造を有し、

20

あるキーに関連するデータポインタを追加もしくは削除する際に、

そのキーに関連する複数のデータポインタを格納する前記重複データ構造に対し、そのデータポインタを追加もしくは削除するための情報としてテーブルデータ関連情報と前記重複データ構造へのポインタと前記キーを含む第1のログレコードを取得する機能と、

リーフノードに格納されたそのキーを含むインデクスエントリを更新する機能と、

前記重複データ構造へのそのデータポインタの追加もしくは削除が行われたことを示す情報を含む第2のログレコードを取得する機能と、

前記重複データ構造に、そのデータポインタを追加もしくは削除する機能とを有し、

上記データポインタの追加もしくは削除を伴う処理結果を取り消す際に、

30

上記取得したログレコードを受け取り、受け取ったログレコードが前記第2のログレコードであった場合には、その前記第2のログレコードを保持し、受け取ったログレコードが前記第1のログレコードであった場合には、第2のログレコードを保持しているかを確認し、第2のログレコードを保持している場合には、前記重複データ構造へのそのデータポインタの追加もしくは削除が完了したと判断する機能と、

上記判断において受け取ったログレコードが前記第1のログレコードであった場合にもかかわらず第2のログレコードを保持していないことを確認することにより前記重複データ構造へのそのポインタの追加もしくは削除が完了していないと判断した場合、第1のログレコードに含まれる前記重複データ構造へのポインタが示す前記重複データ構造をアクセスし、第1のログレコードに含まれる前記テーブルデータ情報を追加もしくは削除する機能とを、コンピュータへ実行させるデータ管理プログラムを記憶した計算機読み取り可能な記憶媒体。

40

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明はデータ管理技術に関し、あるインデクスのキーに対して複数のデータを管理する場合に適用して有効な技術に関するものである。

【0002】

【従来の技術】

近年インターネットを利用した業務システムの急速な普及に伴い、システムを支えるデータ

50

ベース管理システム(DBMS)の適用分野も拡大している。さらにそれに伴いDBMSの扱うデータ量も年々増加している。そのシステム中には、ワークフローにおける「業務推移状態」など、値域の狭いデータに対してDBMSの最も代表的な高速検索手段であるB木インデクスを利用するようなアプリケーションも多く、B木インデクスの安定した性能の提供が重要となっている。

【0003】

B木インデクスに関しては、例えば、文献(Jim Gray and Andreas Reuter, Transaction Processing: Concepts and Techniques, Morgan Kaufmann Publishers, 1993)に、基本的にはデータ構造、複数処理からの同時アクセス方法、および回復に関する基本方式について開示している。

10

【0004】

実システムとしての高速なB木インデクスを実現するためには、次の課題を解消する必要がある。

【0005】

(i) 複数処理による同時実行性の向上

レコードを含むデータベーステーブルに対して複数のトランザクションが同時にアクセスしてくる場合、問題が起こる。具体的には、あるトランザクションが1つのレコードを更新しようとしたとき、同時に他のトランザクションが同一レコードにアクセスしようとするとき、競合状況が発生する。競合問題の解決法の1つとして、レコードまたはB木インデクスのある部分(B木インデクス全体、ノード、インデクスエントリ、キー値など)に対するロッキング方式(排他的アクセス方式)がある。ロッキング方式は、データにアクセスする前に強制的にトランザクションにロックを取得させ。このとき他のトランザクションによってロックがかけられている場合、そのロックが衝突するために、必要なロックを取得できないことがある。ロッキング方式は自トランザクションの更新結果や参照結果を確実に保証してくれるが、ロックを取得中は他のトランザクションのアクセスをロック解除まで待たすので、システムの多重度(同時実行性)を高めスループットを上げるためには、1つのトランザクションが同時に取得するロックの数やロックによる影響範囲(粒度)を最小限に抑えることが重要である。

20

【0006】

(ii) 回復処理制御

実システムにおいては、トランザクションの中断、システムダウン、媒体障害など種々の障害に関して、トランザクションの原子性(Atomicity)を保証し、かつB木インデクス内のデータの整合性を保証する必要がある。そのためのログレコード取得方式およびログレコードを使用した回復処理制御が重要となる。ログレコードに取得方式には、インデクスを構成するノードの更新前・更新後の物理イメージをログレコードの情報をして取得する物理ログ方式などがある。物理イメージを用いた回復では、更新を行ったトランザクションの完了まで、他トランザクションはそのページへのアクセスが制限されることになる。すわち、回復処理制御は排他制御と密接に関係し、同時実行性にも影響を与える。

30

【0007】

図11に、従来の代表的なB木インデクスの構造を示す。

40

【0008】

B木インデクスは、1つのルートノードを頂点に、ルートノードから多数のレベルにわたりノードが枝分かれている。枝分かれたノードの内末端のノードをしばしばリーフノードと呼ぶ。またルートノードを含むリーフノード以外のノードをしばしば上位ノードと呼ぶ。各ノード内は以下に示す情報から構成される複数のインデクスエントリを持つ。リーフノード内のインデクスエントリは、データベース中のレコード(テーブルデータ)を示すポインタと、そのレコード(テーブルデータ)の特徴の1つを表すキー値から構成される。また上位ノード内のインデクスエントリは、次の下位レベルにあるノード(子ノード)を示すポインタと、その子ノードから枝分かれて最終的にリーフノードで管理されているキー値の範囲を示す1キー値から構成される。上位インデクスエントリ内のキー値は、

50

B木インデクスへのアクセス・プログラムがルートノードから目的とするレコードへのポインタを含むインデクスエントリが格納管理されているリーフノードへと辿っていく際の道しるべ(判定要素)の役割を担う。通常B木インデクスの1ノードは、データベース上では、アクセス単位であるページによって実現される。

【0009】

【発明が解決しようとする課題】

B木インデクスに対する排他制御方式の一つにインデクスの排他資源の粒度を「キー」とするキー値排他方式もしくはキーレンジ排他方式がある。これらインデクスのキーに対する排他制御を用いた場合、同一キーを持つデータ(行)に対する検索・更新処理はインデクスキー排他によりシリアライズされ、同時実行性が低下する。現在急速な普及をみせているWEB環境下におけるトランザクションは、従来のOLTPのトランザクションの長さに較べ長く、キーに対する排他制御方式では高スループットのサービスの提供は困難である。

10

【0010】

また、図11に示すように多数の同一キー有するB木インデクスに関して、以下のような問題もある。図11のB木インデクスでは、あるキーに関連するテーブルデータ情報を1つのリーフエントリ16で管理している。リーフエントリ16は、キー14、そのキー値に関連するテーブルデータ情報18(データレコード識別子)、そのキー値に関連するテーブルデータに関する情報18の個数17(重複数)を有する。テーブルデータ情報18は重複数分リスト形式にて昇順に配置されている。このような格納方式では、あるキーに対する重複数が多くなると、リーフエントリ自体の長さがリーフノード12に収まりきれなくなる。結果的に、1つのキーに対するテーブルデータ情報を複数のリーフエントリとして、複数のノードに分けて格納しなければならない。図11の例では、キー「28」のエントリがリーフノードN4、N5、N6と3つのノードに分かれて格納されている。この構造では、範囲検索に対応するためのリーフノードの水平方向のポインタを用いたスキャンにおいて、キー「28」のエントリがネックとなる。またさらに、ルートノード10および上位ノード13に格納されている上位エントリ13に含まれるキー14は、「28」に加えテーブルデータ情報でキーとしなければならず、ルートノード10および上位ノード13の数を増大する要因となっている。結果的に、重複キーは、重複キーそのものに対するアクセス処理だけでなく、B木インデクス全体のアクセス性能の低下および格納効率の低下の要因となっている。

20

30

【0011】

GrayはB木構造外へのデータ構造にてテーブルデータ情報を管理する方法を開示しているが、それに対する同時実行制御および回復制御方法に関しては開示されていない。また、物理ログによる回復処理方式について開示しているが、先に示したような同一ページに対する同時実行を実現することができない。このような問題は、同一キーに対するインデクスアクセスを行う場合、一般的に発生する。

【0012】

本発明の目的は上記問題を改善し、複数のアクセスが行われる環境下において、同一キーに対するインデクスアクセスを好適に行うことが可能な技術を提供することにある。

40

【0013】

【課題を解決するための手段】

前記課題を以下の手段により改善する。

【0014】

複数ユーザから同時にアクセスされる環境下で、多数の同一キーを持つB木インデクスを管理する方法において、ある一つのキーに関する複数のデータポインタを格納するデータ構造であり、そのキーを含むインデクスエントリが格納されるリーフノードとは異なるデータ構造であり、そのデータ構造はそのキーを含むインデクスエントリからポイントされ、複数処理による同時アクセスが可能であるデータ構造を有し、あるキーに関連するデータポインタを追加もしくは削除する際に、そのキーに関連する複数のデータポインタを格

50

納する前記データ構造に対し、そのデータポイントを追加もしくは削除するための情報を含む第1のログレコードを取得し、リーフノードに格納されたそのキーを含むインデクスエントリを更新し、前記データ構造へのそのデータポイントの追加もしくは削除が行われたことを示す情報を含む第2のログレコードを取得し、前記データ構造に、そのデータポイントを追加もしくは削除することにより上記課題を改善する。

【0015】

【発明の実施の形態】

複数ユーザからの同時アクセス環境下において、同一キーに対するインデクスアクセスを効率的に行うことが可能な一実施形態のデータベース処理システムについて説明する。

【0016】

まず、本発明の概念を図1を用いて簡単に説明する。

【0017】

本実施形態のデータベース管理システムにおけるB木インデクス1は、図1に示すようにキーを用いて効率的にテーブルデータに関連する情報を取得するためのB木構造101と、キーごとにそのキーに関連付けられた多数のテーブルデータ情報を管理するデータ構造102(以降「重複キーデータ構造」と呼ぶ)から構成される。

【0018】

B木構造101は、ルートノード10、中間ノード11、リーフノード12により構成され、1つのルートノードを頂点に、ルートノードから多数のレベルにわたりノードが枝分かれしている構造をもつ。ルートノード10および中間ノード11には、キー14とそのキー値に関連する下位レベルへのノードを示すポインタ13とを有する上位エントリ13が格納される。それぞれのノード内において、上位エントリ13およびリーフエントリ16は、そのインデクスエントリ13、16が有するキー値の順に格納管理されている。本実施例では、図の左から右方向へキー値が昇順に格納されている。それぞれのノード内において、それらソートされたインデクスエントリをバイナリサーチすることにより、効率的かつ安定して必要なインデクスエントリにアクセスすることができる。一つの上位インデクスエントリがポイントする子ノードには、その上位インデクスエントリ内キー値よりも小さいキー値を持つインデクスエントリが格納されている。

【0019】

キー14に関連するテーブルデータが多数の場合のリーフエントリ20は、キー14、そのキー値に関連するテーブルデータに関する情報の個数17、及び、テーブルデータに関する情報を格納する重複キーデータ構造102へのポインタ19を有する。図1の例では、キー値「28」に対して、関連するテーブルデータ情報数すなわち重複数が「35,678」であることを示している。それら35,678個の情報は、ポインタ19「Nx1」にて指される重複データ構造102に格納管理されている。

【0020】

一方、キー値「28」以外のインデクスエントリのように図1の例では、重複数すなわちテーブルデータに関する情報数があまり多くない場合、重複キーデータ構造102を持たないリーフエントリ16として格納される。リーフエントリ16は、キー14、そのキー値に関連するテーブルデータに関する情報18(データレコード識別子)、そのキー値に関連するテーブルデータに関する情報18の個数17を有する。データレコードの追加もしくは更新によって、そのキーに関連するテーブルデータ情報の重複数がある敷居値を超えた際に、上記重複キーデータ構造102をポイントするリーフエントリ20に移行する。

【0021】

また、多数のテーブルデータに関連する情報を格納管理する重複キーデータ構造102は以下の特徴を持つ。

【0022】

(1)アクセスの際に排他を取得することなく、複数処理による同時アクセスが可能である。

10

20

30

40

50

## 【 0 0 2 3 】

(2)通常処理のアクセスを制限することなく取り消し処理が可能である。

## 【 0 0 2 4 】

ここで、矢印 1 0 3 はデータレコードの更新に伴う B 木インデクス 1 に対するキー値「 2 8 」、テーブルデータ関連情報「 P 1 8 」の追加更新処理の様子を示している。まず、この追加更新処理は、キー値「 2 8 」を用いて B 木構造 1 0 1 のルートノード 1 0 ( N 1 ) からサーチを開始し、中間ノード 1 1 ( N 2 ) を経由し、更新対象のキー値「 2 8 」を有するリーフエントリ 2 0 が格納されているリーフノード 1 2 ( N 5 ) に辿り着く。そして、リーフエントリ 2 0 を確定し、第 1 のログレコード 4 1 を取得し、エントリ内の重複数 1 7 をインクリメントする。第 1 のログレコード 4 1 には、重複キーデータ構造 1 0 2 に対する更新処理を続行するために必要な情報が含まれている。さらに、エントリ内の重複キーデータ構造 1 0 2 へのポインタ 1 9 ( N x 1 ) を取得し、重複キーデータ構造 1 0 2 へのアクセスを行う。そして、第 2 のログレコード 4 2 を取得し、重複キーデータ構造 1 0 2 に対し、テーブルデータ関連情報「 P 1 8 」を追加する。第 2 のログレコード 4 2 には、重複キーデータ構造 1 0 2 に対しテーブルデータ関連情報 1 8 の追加が完了したことを示す情報が含まれている。

10

## 【 0 0 2 5 】

図 1 の例を用いてテーブルデータ関連情報の追加に関し説明したが、テーブルデータ関連情報の削除に関しても同様である。リーフエントリ 2 0 をアクセスし、第 1 のログレコード 4 1 を取得し、エントリ内の重複数 1 7 をデクリメントする。そして、重複キーデータ構造 1 0 2 をアクセスし、第 2 のログレコードを取得した後に、重複キーデータ構造 1 0 2 からテーブルデータ関連情報 1 8 の削除を行う。

20

## 【 0 0 2 6 】

以上のように、リーフエントリの更新から重複キーデータ構造への変更順序が一定であり、インデクス 1 を構成するデータリソースへの排他制御を必要としないことから、複数処理の同時実行が可能である。また、一連の変更処理結果に対し取り消し要因が発生した場合、もしくはリーフエントリに対する更新直後に処理の中断要因が発生した場合でも、前述の第 1 のログレコード 4 1 および第 2 のログレコード 4 2 を使用することにより、同時実行性を失うことなしに取り消し処理を実施することが可能である。

## 【 0 0 2 7 】

具体的には、取り消し処理時、第 1 のログレコード 4 1 と第 2 のログレコードがペアで取得されている場合、一連の変更処理は完了していると判断し、図 1 の例では、キー値「 2 8 」、テーブルデータ関連情報「 P 1 8 」の削除処理を行う。取り消し処理時、第 1 のログレコード 4 1 しか取得されていない場合は、リーフエントリ 2 0 に対する変更処理しか行われていないので、第 1 のログレコードに含まれる情報を用いて重複キーデータ構造 2 0 1 に対するテーブルデータ関連情報「 P 1 8 」の追加処理を行った後、改めてキー値「 2 8 」、テーブルデータ関連情報「 P 1 8 」の削除処理を行う。取り消し処理においても B 木インデクス 1 に対する変更の順序は常に片方向であるため、他処理との同時実行性を損なうことはない。以上のように、取り消し処理を交えても高い同時実行性を提供することが可能である。

30

40

## 【 0 0 2 8 】

さらに、本重複キーデータ構造 2 0 1 をリーフエントリ 2 0 からポイントし、リーフノード 1 2 の外に格納管理することから、リーフエントリ 2 0 が格納されているリーフノード 1 2 内には、すべてのテーブルデータ関連情報 1 8 をリーフエントリ 1 6 の有する方式に較べ、より多くの他のキーを有するリーフエントリの格納が可能となる。そしてリーフエントリを格納するリーフノード 1 2 の必要数を抑え、最終的に上位ノード 1 1 の必要数も抑え、B 木インデクス全体の格納容量削減にもなる。多くのテーブルデータ情報を持たないキーに対するアクセスに関しても、B 木構造の容量が抑えられるため、B 木構造を辿るすべての検索処理および変更処理に対し安定した性能を提供することが可能である。

## 【 0 0 2 9 】

50

次に図2に本実施形態のデータベース管理システムの概略構成を示す。ユーザが作成したアプリケーションプログラム6と、問い合わせやリソース管理などのデータベースシステム全体の管理を行うデータベース管理システム2がある。上記のデータベース管理システム2は、論理処理部21、物理処理部22、システム制御23と、データベースアクセス対象となるデータを永続的にあるいは一時的に格納するデータベース3、そしてシステムログ4を有する。また、上記データベース管理システム2はネットワークなどを介して他のシステムと接続されている。

【0030】

上記論理処理部21は、問合せの構文解析・意味解析を行う問合せ解析211と、適切な処理手順を生成する最適化処理212と、処理手順に対応したコードの生成を行うコード生成213を具備している。また、上記生成されたコードに基づき、コードを解釈しデータベース処理実行の指示を物理処理部22に対して行うコード解釈実行214を具備している。

10

【0031】

上記物理処理部22は、データベース3に格納されているテーブルデータ5に対し、データの格納、削除、更新、検索処理をデータベースバッファ24を介して行うテーブルデータ管理部221を具備する。また、データベース3に格納されているインデクス1を管理するB木インデクス管理部25を具備している。インデクス1は、キーを用いて効率的に関連情報にアクセスするためのB木構造と、キーごとにそのキーに関連付けられた多数のテーブルデータ情報を管理する重複キーデータ構造から成る。

20

【0032】

B木インデクス管理部25は、テーブルデータ5の更新に伴い、それに関連するB木インデクス1に対し、B木構造に対する変更を行うB木構造アクセス処理26の機能、および重複データデータ構造に対する変更を行う重複キーデータ構造アクセス処理27の機能を有する。また、インデクス1をアクセスすることにより、検索条件であるキーから関連するテーブルデータを高速に検索するための機能をB木構造アクセス処理26と重複キーデータ構造アクセス処理27内に有する。

【0033】

上記物理処理部22は、また、システムで共用するリソースの排他制御223を具備している。本実施例では、データに対する整合性は、テーブルデータ5に対する排他制御(排他リソース：テーブルデータ格納ページID, テーブルデータIDなど)のみで実現し、B木インデクス1のリソース(インデクスID, インデクスページID, キー値など)に対する排他制御は行わない。

30

【0034】

システムログ4は、データベース3へのデータ挿入、更新、削除などの更新履歴情報を蓄積する。その蓄積される情報には、B木インデクス1に対する変更に関する履歴情報も含まれる。そして、トランザクションの取り消しやシステムダウンなどの種々の状況において、B木インデクス1のデータ整合性を回復するために使用される。

【0035】

上記システム制御23は、ユーザからの問い合わせ入力、問い合わせ結果の返却、コマンド受付によるデータベース保守などを行う。また、前述のシステムログ4の管理を行い、上記物理処理部22と連携し、データベース3変更の際に履歴情報であるログレコードの取得26を行う。そして、システム制御23は、読み込んだログレコードがB木インデクスに対する変更処理に関する場合に、物理処理部22のB木インデクス管理部25にそのログレコードを渡し、取り消し処理を依頼する。依頼されたB木インデクス管理部25は、ログレコード内の情報を基に変更取り消し処理を行う。

40

【0036】

ユーザは、アプリケーションプログラム6を介し、データベース管理システム2を用いて、データベース3内のテーブルデータ5を構成するデータレコードを探索、アクセスし、変更することができる。データベース管理システム2はデータレコードへの効率のよいア

50

クセスを実現するために、B木インデクス1を利用する。

【0037】

データベース3内のB木インデクス1あるいはテーブルデータ5は、アクセス単位であるページ単位に、データベースバッファ24に読み込んでくる(GETしてくる)ことによりアクセスが可能になる。バッファ上のページを参照(検索)した後、バッファをRELEASEすることにより使用していたバッファに対する利用終了を宣言し、他処理からのそのバッファの利用を許可する。また、バッファ上のページに対して更新を行った後、バッファをPUTすることにより当該バッファ上のページがデータベースに反映することを宣言し、使用していたバッファの他処理からの利用を許可する。そのバッファ上のページはすぐにはデータベースに反映されず、ある契機にデータベースに反映される。バッファにGETしたページに他のトランザクションが不当にアクセスしないようにバッファに対してラッチをかける。ラッチは一種のロックであるが通常のロックよりもはるかに取得期間が短く、獲得と解除を安価に行うことができる。本実施例では、インデクスノードを参照する際に共用ラッチを、更新する際に排他ラッチをかける。排他ラッチに関する処理はシリアライズされ、共用ラッチ間ではバッファに対する同時参照を可能にする。そして、バッファ解放の際に獲得中ラッチを解除し、排他ラッチを獲得しようとしている他の処理のアクセスを許可する。

10

【0038】

図3は本実施形態のコンピュータシステムのハードウェア構成の一例を示す図である。コンピュータシステム3000は、CPU3002、主記憶装置3001、磁気ディスク装置等の外部記憶装置3003及び多数の端末3004で構成される。主記憶装置3001上には、図2を用いて先に説明したデータベース管理システム2が置かれ、外部記憶装置3003上にはデータベース管理システム3が管理するテーブルデータ5とインデクス1を含むデータベース3が格納される。また、システムログ4も外部記憶装置3003上に格納される。さらに、データベース管理システム2を実現するプログラム3100も外部記憶装置3003上に格納される。

20

【0039】

図8は本実施形態のインデクスのログレコードの一例を示す図である。

【0040】

図8の41に示す様に、図1にて説明したリーフエントリ変更時に取得される第1のログレコード41は、その変更処理が属するトランザクションの情報などを含むログレコードヘッダ410、変更を示す操作コード411、キー412、変更対象であるテーブルデータ情報413、そして変更されたリーフエントリが指す重複キーデータ構造へのポインタ414により構成される。操作コード411には、追加処理や削除処理を示すコードが設定される。図8の例ではその設定コードは追加処理を示す「INSERT」である。ログレコードヘッダ410には変更対象のB木インデクスのインデクスIDやそのインデクスが格納されるエリア情報も含まれる。キー412に設定される値は変更対象であるリーフエントリ20内のキー14の値と同値である。図8の例ではその値は「28」である。

30

【0041】

重複キーデータ構造へのポインタ414は、リーフエントリ変更は終了したが重複キーデータ構造変更が完了していない状況にて変更結果取り消し要因が発生した場合に、重複キーデータ構造に対する変更処理を続行するための情報である。

40

【0042】

また図8の42に示す様に、図1にて説明したリーフエントリ変更時に取得される第2のログレコードは、その変更処理が属するトランザクションの情報などを含むログレコードヘッダ420、変更を示す操作コード421、そして変更対象であるテーブルデータ情報422により構成される。411と同様に操作コード421には、追加処理や削除処理を示すコードが設定される。図8の例ではその設定コードは追加処理を示す「INSERT」である。また、テーブルデータ情報422には、第1のログレコード41の413と同じ値が設定されることになる。第2のログレコード42は、先に図1を用いて説明したように、

50



重複キーデータ構造に対する変更処理が終了しているかを示すマーカの役割を担う。

【 0 0 4 3 】

図 4 は本実施形態の B 木インデクス管理部 2 5 の処理手順を示すフローチャートである。図 4 では B 木インデクスに対するテーブルデータ情報追加の場合の B 木構造アクセス処理 2 6 の処理内容を表している。

【 0 0 4 4 】

まず、ステップ 4 0 1 において、キーを用いて、B 木構造をルートノードから上位ノードへと辿り、リーフノードを確定後、そのリーフノードをデータベースバッファに読み込みバッファ固定(GET)を行う。次にステップ 4 0 2 において、リーフノード内に格納されているリーフエントリをサーチし、変更対象であるリーフエントリの確定を行う。そして、ステップ 4 0 3 においてそのリーフエントリのキーに関連するテーブルデータ情報数(重複数)をインクリメントする。ここで、ステップ 4 0 4 にて、リーフエントリ内に含まれる重複キーデータ構造を指すポインタを取得する。さらに、図 8 を用いて説明した第 1 のログレコードを作成し(ステップ 4 0 5)、作成した第 1 のログレコードを取得する(ステップ 4 0 6)。最後にリーフノードをデータベースに書き出すように予約し、バッファ固定を解除(PUT)する(ステップ 4 0 7)。以上でリーフエントリ変更処理を終了し(ステップ 4 0 8)、引き続き重複キーデータ構造に対する変更処理に進む。

【 0 0 4 5 】

本フローチャートでは、リーフエントリ変更(4 0 4)、ログレコード取得(4 0 5)の順に処理が行われているが、実際のデータベースへのリーフノードに対する変更結果の反映は、W A L (Write-ahead log) プロトコルに従ってログレコードのシステムログへの書き出しが実施された後に行われる。

【 0 0 4 6 】

図 5 は本実施形態の B 木インデクス管理部 2 5 の処理手順を示すフローチャートである。図 5 では B 木インデクスに対するテーブルデータ情報追加の場合の重複キーデータ構造アクセス処理 2 7 の処理内容を表している。

【 0 0 4 7 】

重複キーデータ構造に対する変更処理は、図 4 のフローチャートで説明したリーフエントリに対する変更処理に引き続き行われる。

【 0 0 4 8 】

まず、ステップ 5 0 1 において、先に図 4 のステップ 4 0 4 にて取得したポインタを用いて重複キーデータ構造へアクセスし、ステップ 5 0 2 において、テーブルデータ情報を追加する位置をサーチする。そして、図 8 を用いて説明した第 2 のログレコードを作成し(ステップ 5 0 3)、作成した第 2 のログレコードを取得する(ステップ 5 0 4)。ここで、重複キーデータ構造に追加領域があるかを判定する(ステップ 5 0 5)。追加領域が無い場合、ステップ 5 0 6 に進み、新規領域の割り当てを行い、新規領域追加に伴う重複キーデータ構造の変更を行う(ステップ 5 0 7)。具体的には重複キーデータ構造が複数のページから構成される場合は、新規ページ割り当てが行われる。その後ステップ 5 0 8 へ進む。また、ステップ 5 0 5 の判定結果が領域有りの場合、ステップ 5 0 8 へ進み、重複キーデータ構造へテーブルデータ情報を追加し、重複キーデータ構造に対する変更処理を終了する(ステップ 5 0 9)。

【 0 0 4 9 】

図 6 は本実施形態の B 木インデクス管理部 2 5 の処理手順を示すフローチャートである。図 6 では B 木インデクスに対するテーブルデータ情報削除の場合の B 木構造アクセス処理 2 6 の処理内容を表している。

【 0 0 5 0 】

まず、ステップ 6 0 1 において、キーを用いて、B 木構造をルートノードから上位ノードへと辿り、リーフノードを確定後、そのリーフノードをデータベースバッファに読み込みバッファ固定(GET)を行う。次にステップ 6 0 2 において、リーフノード内に格納されているリーフエントリをサーチし、変更対象であるリーフエントリの確定を行う。ここで、

10

20

30

40

50

ステップ603において、リーフエントリ内に含まれる重複キーデータ構造を指すポインタを取得する。そして、ステップ604において、そのリーフエントリのキーに関連するテーブルデータ情報数(重複数)が、1かどうかを判定する。重複数が1の場合、ステップ605に進みリーフエントリの削除を行い、ステップ607に進む。ステップ604の判定結果が重複数が1より大きい場合、ステップ606に進み、そのリーフエントリのキーに関連するテーブルデータ情報数(重複数)をデクリメントする。

【0051】

さらに、図8を用いて説明した第1のログレコードを作成し(ステップ607)、作成した第1のログレコードを取得する(ステップ608)。最後にリーフノードをデータベースに書き出すように予約し、バッファ固定を解除(PUT)する(ステップ609)。以上でリーフエントリ変更処理を終了し(ステップ610)、引き続き重複キーデータ構造に対する変更処理に進む。

10

【0052】

図7は本実施形態のB木インデクス管理部25の処理手順を示すフローチャートである。図7ではB木インデクスに対するテーブルデータ情報削除の場合の重複キーデータ構造アクセス処理27の処理内容を表している。

【0053】

重複キーデータ構造に対する変更処理は、図6のフローチャートで説明したリーフエントリに対する変更処理に引き続き行われる。

【0054】

まず、ステップ701において、先に図4のステップ603にて取得したポインタを用いて重複キーデータ構造へアクセスし、ステップ702において、削除対象テーブルデータ情報の存在位置をサーチする。そして、図8を用いて説明した第2のログレコードを作成し(ステップ703)、作成した第2のログレコードを取得する(ステップ704)。

20

【0055】

次に、ステップ705において、重複キーデータ構造からのテーブルデータ情報削除を行う。ここで、テーブルデータ情報の削除によって、回収可能な重複キーデータ構造を構成する領域が発生したかを判定する(ステップ706)。判定の結果、回収可能領域が発生した場合、ステップ707へ進み領域の回収を行い、領域回収に伴う重複キーデータ構造の変更を行う(ステップ708)。具体的には重複キーデータ構造が複数のページから構成される場合は、ページの回収が行われる。回収されたページは他の処理における新規ページ割り当てにより再利用されることとなる。ステップ708後ステップ709にて重複キーデータ構造に対する変更処理を終了する。また、ステップ706における判定結果が、回収可能領域発生なしの場合、そのままステップ709にて重複キーデータ構造に対する変更処理を終了する。

30

【0056】

図9は本実施形態のインデクス管理部のインデクス更新結果取り消し処理の処理手順を示すフローチャートである。

【0057】

まず、ステップ901において、システム制御からB木インデクス変更に関するログレコードを受け取る。ステップ902において、処理すべきログレコードが存在するかを判定する。処理すべきログレコードが存在しない場合、ステップ913にてインデクス更新結果取り消し処理を終了する。処理すべきログレコードが存在する、すなわち受け取りログレコードが存在する場合、ステップ903へ進みログレコードの種別を判定する。ここでまずログレコードが第2のログレコードかどうかを判定する(ステップ903)。第2のログレコードである場合、ステップ904へ進み第2のログレコードを保持する。そして、ステップ901へ戻り次のログレコードの処理を行う。ここで次のログレコードは、第1のログレコードのはずである。ステップ903の判定において、第1のログレコードではない場合、さらにステップ905にて、ログレコードが第1のログレコードかどうかを判定する。第1のログレコードでない場合、ステップ906にてその他のログレコードに

40

50

対応する更新結果取り消し処理を行い、ステップ901へ戻り、次のログレコードの処理を行う。第1のログレコードである場合、以下の処理を行う。まず第2のログレコードが保持されているかを判定し(ステップ907)、保持されている場合には、リーフエントリおよび重複キーデータ構造に対する一連の変更処理が完了していると判断し、保持されている第2のログレコードを廃棄し(ステップ908)、ステップ910へ進む。第2のログレコードが保持されていない場合、リーフエントリの変更処理は終わっているが対となる重複キーデータ構造に対する変更処理が完了していないと判断し、ステップ909にて重複キーデータ構造の対する更新処理を続行する。その後ステップ910へ進む。ステップ910では、処理中のログレコードが通常処理にて取得されたものであるかどうかを判定する。通常処理にて取得されたログレコードである場合、関連する変更結果を取り消す必要がある。ステップ911にてその取り消し処理を行う。また、通常処理時取得ではない、すなわち取り消し処理時における取得の場合、取り消し処理はステップ909を含めすでに完了していると判断し、一連のB木インデクス更新結果の取り消し処理を終了し、ステップ901へ戻り、次のログレコードの処理を行う。

10

**【0058】**

図10は本実施形態の重複キーデータ構造102の構成の一例を示す図である。

**【0059】**

図10の例では、重複キーデータ構造102は、実際にテーブルデータ情報を格納するページ1002と、そのページ1002を指すポインタを有する管理エントリ1003を格納するページ1001より構成される。

20

**【0060】**

ページ1001に格納される管理エントリ1003は、ページ1002を指すポインタ1005、およびそのページ1002内に格納されるテーブルデータ情報の最大値1004を有する。ページ1001内の管理エントリ1003は、テーブルデータ情報の最大値1004の順にソートされている。

**【0061】**

また図10に示すように、各ページ1001は、そのページ1001内に格納される管理エントリ1003内の最大情報1004よりも大きな最大情報を持つ管理エントリが格納されるページ1001へのポインタを有する。リーフエントリ20からは、最も小さいテーブルデータ情報を格納するページ1002をポイントする管理エントリが格納されてるページ1001がポイントされる。

30

**【0062】**

以上のデータ構造により、図5および図7にて説明した重複キーデータ構造内のサーチを効率的に行うことが可能である。

**【0063】**

また、ページ1001もしくはページ1002の新規割り当て・回収は、データベースバッファへのGET・PUTにより容易に実現可能である。新規割り当て・回収の際の第2のログレコードは、図8の例で示した構成情報の他にページ1001もしくはページ1002のページ識別子を付け加えることにより、取り消し処理における変更続行処理が容易に可能である。

40

**【0064】**

以上示したフローチャートの処理は、図3で例として示したコンピュータシステムにおけるプログラムとして実行される。しかし、そのプログラムは図3の例の様にコンピュータシステムに物理的に直接接続される外部記憶装置に格納されるものと限定はしない。ハードディスク装置、フロッピーディスク装置等のコンピュータで読み書きできる記憶媒体に格納することができる。また、ネットワークを介して図3のコンピュータシステムとは別のコンピュータシステムに接続される外部記憶装置に格納することもできる。

**【0065】**

このようにすることにより、同一キーを持つテーブルデータ情報が、アクセスの際に排他を取得することなく複数処理による同時アクセスが可能であり、且つ通常処理のアクセス

50

を制限することなく取り消し処理が可能である重複キーデータ構造で管理されるため、同一キーに対して同時にアクセスする複数の処理に対し、取り消し処理をも交えた高い同時実行性を提供することが可能である。

【0066】

その重複キーデータ構造はリーフエントリからポイントされ、リーフノードの外に格納管理されることから、重複キーを有するリーフエントリが格納されているリーフノード内により多くの他のキーを有するリーフエントリの格納が可能となる。その結果、B木インデクスの格納容量が抑えられ、安定したアクセス性能を提供することが可能になる。

【0067】

以上のように、多数の重複キーを有する場合でも同時実行性に優れ、格納効率のよいB木インデクスを用いたデータ管理方法を提供することができる。

10

【0068】

本B木インデクスを用いたデータ管理方法は、ワークフローにおける「業務推移状態」など、値域の狭いデータに対してB木インデクスを利用するようなアプリケーションに対して非常に有効である。そのアプリケーションが利用するB木インデクスでは、値域が狭いため一つのキーに関連するデータ重複数が膨大であり、さらに「業務推移状態」などの変化によってB木インデクスに関して頻繁に更新が発生するためである。

【0069】

以上、B木について説明してきたが、本発明は、複数のアクセスが行われる環境下において、同一キーに対するインデクスアクセスが発生する場合に同様に適用可能であり、同一キーに対するインデクスアクセスを好適に行うことが可能なデータ管理方法および装置を提供することが可能である。

20

【0070】

【発明の効果】

本発明によれば、複数のアクセスが行われる環境下において、同一キーに対するインデクスアクセスを好適に行うことが可能なデータ管理方法および装置を提供することができる。

【図面の簡単な説明】

【図1】本発明の概念図である。

【図2】本実施形態のB木インデクスを有するデータベース処理システムの機能ブロックを示す図である。

30

【図3】本実施形態のコンピュータシステムのハードウェア構成の一例を示す図である。

【図4】本実施形態のB木インデクス管理部25のインデクス更新処理の処理手順を示すフローチャートである。

【図5】本実施形態のB木インデクス管理部25のインデクス更新処理の処理手順を示すフローチャートである。

【図6】本実施形態のB木インデクス管理部25のインデクス更新処理の処理手順を示すフローチャートである。

【図7】本実施形態のB木インデクス管理部25のインデクス更新処理の処理手順を示すフローチャートである。

40

【図8】本実施形態のB木インデクスのログレコードの一例を示す図である。

【図9】本実施形態のB木インデクス管理部のインデクス更新結果取消し処理(回復処理手順)の処理手順を示すフローチャートである。

【図10】本実施形態の重複キーデータ構造の構成の一例を示す図である。

【図11】従来の多くの重複キーを有するB木インデクスの構成を示す図である。

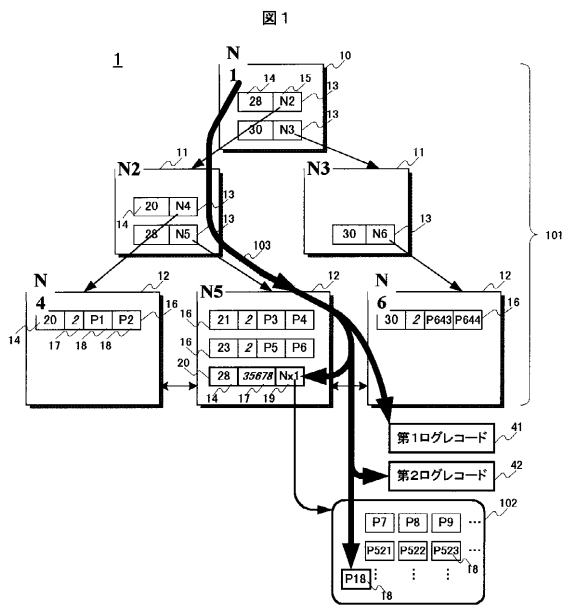
【符号の説明】

1 ... B木インデクス、2 ... データベース管理システム、3 ... データベース、4 ... システムログ、5 ... テーブルデータ、6 ... アプリケーションプログラム、21 ... 論理処理部、22 ... 物理処理部、23 ... システム制御、25 ... B木インデクス管理部、26 ... B木構造アクセス処理、27 ... 重複キーデータ構造アクセス処理、101 ... B木構造、102 ... 重複キ

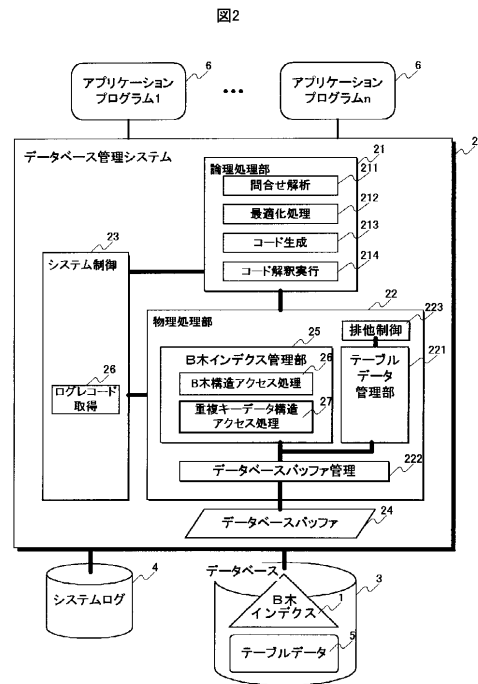
50

-データ構造、4 1...第1ログレコード、4 2...第2ログレコード。

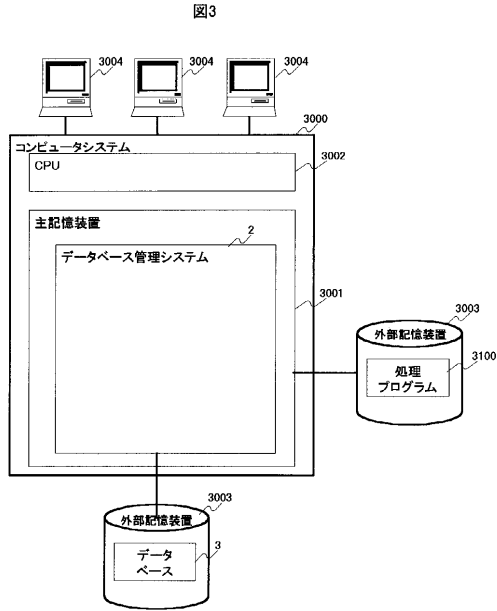
【図1】



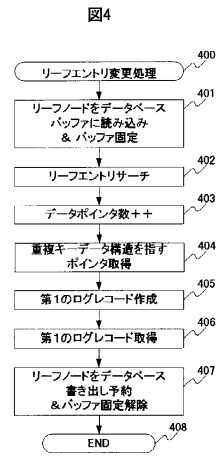
【図2】



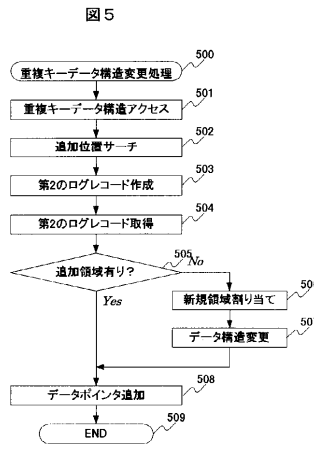
【 図 3 】



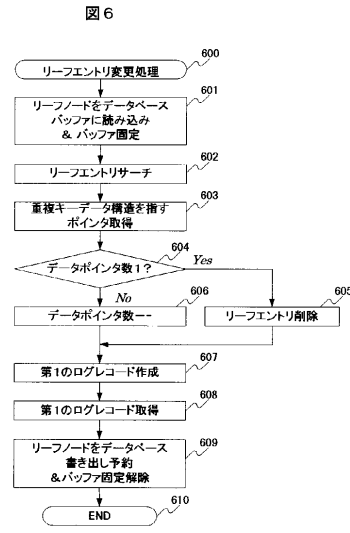
【 図 4 】



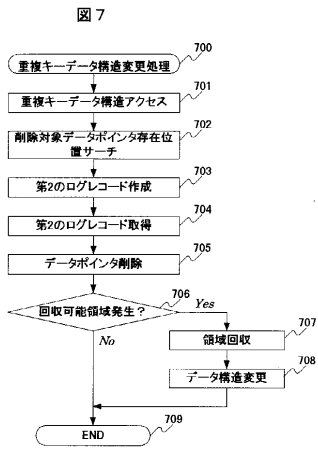
【 図 5 】



【 図 6 】

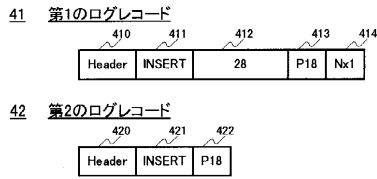


【 図 7 】

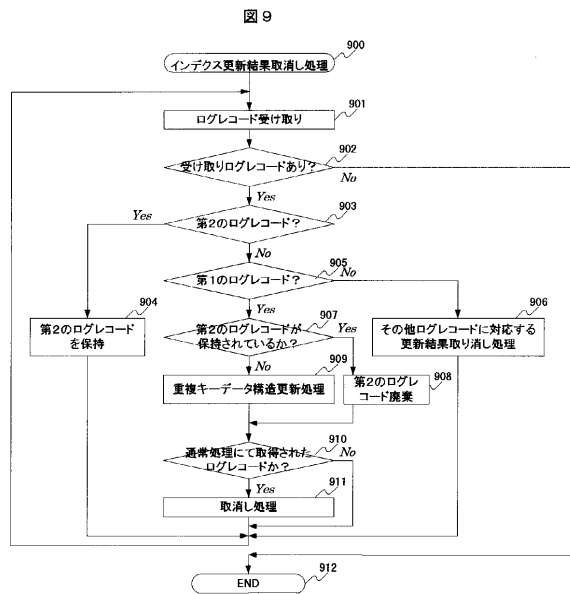


【 図 8 】

図 8

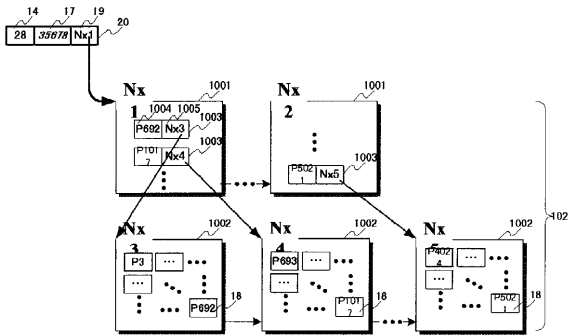


【 図 9 】



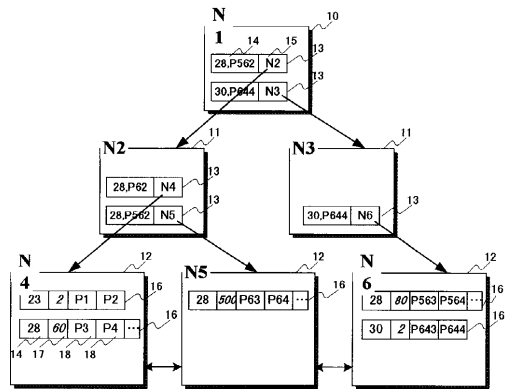
【 図 10 】

図 10



【 図 11 】

図 11



フロントページの続き

(58)調査した分野(Int.Cl. , D B名)

G06F 12/00

G06F 17/30

JSTPlus(JDreamII)