

ロックフリー索引のトライ木化による高速化に関する研究

井戸 佑^{1,a)} 杉浦 健人^{1,b)} 石川 佳治^{1,c)} 陸 可鏡^{1,d)}

概要：代表的な索引構造である B+木は様々なデータを格納可能な汎用性の高い索引である．一方で，扱うキーに制限を加え，索引構造を最適化させることで性能を向上させる研究も行われてきた．本研究では，著者らの研究室で開発しているロックフリー B+木 (B^c 木) に対し同様の拡張および性能改善を行う．具体的には，扱うキーをバイナリ比較可能なものに制限することでトライ木の構造を適用し，空間利用効率の向上およびそれに伴う検索性能の向上を図る．

1. はじめに

Mass 木は B^+ 木にトライ木構造を適応させることで，キャッシュ効率を改善させた． B^c -forest では B^c 木に対し，同様の改善を図る．

Mass 木の観点では B^+ 木から B^c 木になることにより，ロックフリーによる書き込み性能の改善を図る．また B^c 木の観点では，整数型や文字列型などのバイナリ比較可能なキーに対し，キャッシュ効率の改善を図る．

2. 関連研究

関連の深い索引構造として，同時実行制御においてロックを取得しない B^c 木，および B^+ 木にトライ木の構造を組み合わせた Mass 木について紹介する．

2.1 B^c 木

B^c 木はマッピングテーブル，ノード内バッファという構造上の特徴と CAS 命令を用いたロックフリー索引である． B^c 木の概形を図 1 に示す．

2.1.1 データ構造の概観

B^+ 木と同様に， B^c 木は索引層およびデータ層によって構成される．索引層のノード（中間ノード）は分割キーと子ノードへのポインタの組を格納し，木の下方への検索を補助する．構造は B^{link} 木に則っており，各ノードが同じ階層の右兄弟への参照リンクを持つ．ノード間の繋がりはマッピングテーブルにより仮想化する．各ノードは自身の子ノードや兄弟ノードへのポインタを直接持つ代わりに

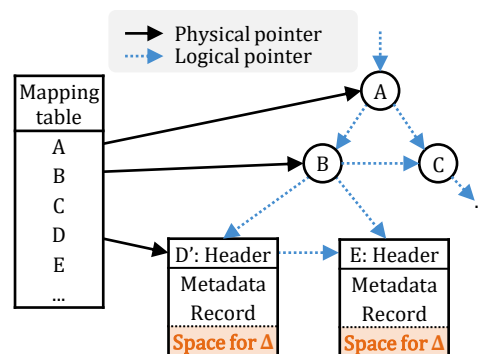


図 1 B^c 木の概形

マッピングテーブル上の ID (logical page ID, LPID) を持つ．各ノードへの参照はマッピングテーブルを用いた間接参照を採用し，マッピングテーブル内の物理ポインタを差し替えることでそのノードへの参照を一括で変更する．

各ノードの領域は不変領域と可変領域（ノード内バッファ）に分けられる．不変領域はノードヘッダおよびソート済みのレコードを格納する．ヘッダは不変領域の情報を管理し，構造変更時のみその値が変更される．可変領域はステータスワードの格納と差分レコードを挿入するための書き込みバッファの役割を果たす．ステータスワードは可変領域の情報を管理し，ノードの現在の状態や残容量などを管理する．

2.1.2 レコード操作の概観

ステータスワードを CAS 命令で更新することによって，ロックフリーな書き込みを実現している．各ノードへ構造変更操作を行う際は，構造変更後のノードから構造変更前のノードへ物理リンクを張り，古いノードへの参照を可能にする．

¹ 名古屋大学大学院情報学研究科
Graduate School of Informatics, Nagoya University

a) ido@db.is.i.nagoya-u.ac.jp

b) sugiura@i.nagoya-u.ac.jp

c) ishikawa@i.nagoya-u.ac.jp

d) lu@db.is.i.nagoya-u.ac.jp

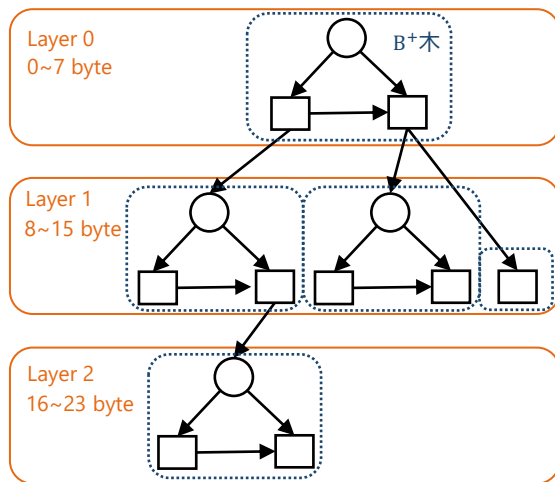


図2 Mass木の概形

2.2 Mass木

Mass木はB⁺木を基本単位とした階層構造やレコードメタデータの削除により、キャッシュ効率を改善した索引構造である。Mass木の概形を図2に示す。

Mass木は複数のB⁺木とlayer構造から構成される。Layer 0はキーの先頭0~7 byteで構成されるB⁺木である。先頭8 byteで一意性が確保できる場合には、Layer 0で完結する。先頭8 byteで一意性が確保出来ない場合、Layer 1(キーの8~15 byteで構成されるB⁺木)を作成し、Layer 0からLayer 1への物理リンクを張る。同様に、複数のB⁺木やLayerを作成し、トライ木に似た構造を持つのがMass木の特徴である。Mass木は整数型や文字列型など、分割が可能なキーに限定することで上記に示す階層化(共通部分の集約)を可能にしている。

また、Mass木は固定長キーおよび固定長ペイロードに特化したノードレイアウトを利用している。B^c木のような可変長キーおよび可変長ペイロードに対応する索引構造では、各レコードに対応する固定長レコードメタデータを利用することでノード内のレコードの配置等を管理している。Mass木では、8 byte分割によりキーが8 byteで固定される。更に、ペイロードに関してもポインタの活用やインラインを固定長に限定することで、メタデータの利用を回避している。

以上のトライ木構造の利用やレコードメタデータの削除により、Mass木はキャッシュ効率を改善している。

3. B^c-forestの構造

B^c-forestはBc木をバイナリ比較可能なキーに最適化した索引構造である。Mass木のように8 byte単位でキーを分割、階層分けし、各階層でのレコード管理にはBc木を利用する。つまり、各Bc木の中間ノードでは固定長の部分キーおよび子ノードへのポインタのみを管理することとなり、レコードメタデータの除外によるキャッシュ効率の

改善が可能となる。一方で、葉ノードではposting listを用いて共通する部分キーを持つレコードを管理し、少数のレコードのみからなる下層の生成を抑制する。

3.1 中間ノードにおけるレコードメタデータの除外

B^c-forestでは中間ノードはキーを8 byteで分割しているため、固定長キーとして扱うことが出来る。また、ペイロードは子ノードへのポインタであるため固定長である。この特性を利用し、B^c-forest内のB^c木における中間ノード内のレコードメタデータの除外を行う。これにより、索引層における探索性能の改善及びキャッシュ効率を改善を図る。

3.2 葉ノードにおけるposting listの導入

Mass木においては、空間利用効率が問題となる。図3は末尾3 byteが異なる19 byteの2キーを格納したMass木の索引構造である。先頭8 byteが共通するため、layer 1を作成する。同様に8~16 byte目が共通するため、layer 2を作成する。Mass木では本来、1つのノードに最大16個のレコードを格納することが出来る。しかしlayer 1では、1つのレコードしか格納されていない状態でlayer 2を作成している。layer 1にのみ注目すると、空間利用率は約6.25%しかない。

B^c-forestでは空間利用率の改善として、posting listを導入する。posting listでは、1つのキーに対して複数のペイロードを対応付けることが出来る。図4はposting listを導入した際の索引構造を示したものである。layer 1においてposting list作成することで、layer 2の無駄な階層化と空間利用率の悪化を防ぐ。

4. B^c-forestのノード操作

4.1 書き込み

4.2 読み取り

5. B^c-forestの構造変更操作

5.1 統合

5.2 分割

5.3 新層作成

6. おわりに

謝辞 本研究はJSPS 科研費JP20K19804, JP21H03555, JP22H03594, JP22H03903の助成、および国立研究開発法人新エネルギー・産業技術総合開発機構(NEDO)の委託業務(JPNP16007)の結果得られたものである。

参考文献

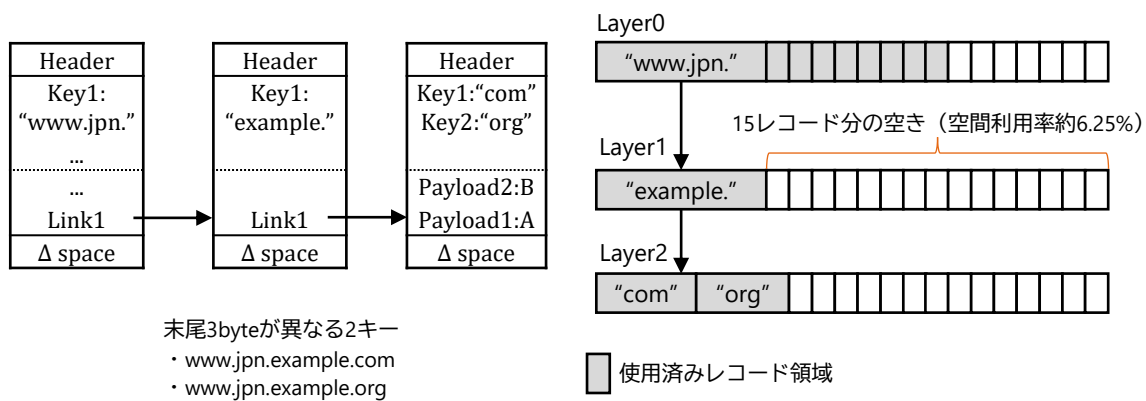


図3 Mass 木の空間利用率

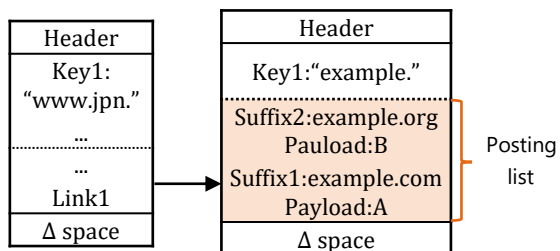


図4 posting list の導入