VAST木:木構造索引の圧縮を用いたベクトル命令による 大規模データ探索の高速化

山室 健^{1,a)} 鬼塚 真^{2,b)}

受付日 2014年12月21日, 採録日 2015年4月7日

概要:大規模なデータに対して空間コストが低く, CPUの実行効率が高い on-memory の木構造索引 (VAST 木)を提案する.2分木や B+木に代表される索引は基本的でかつ重要なデータ構造として幅広く活用され ている.しかしデータの大規模化に対して索引サイズの肥大化や CPUの実行効率低下の問題が顕在化する.特に CPUの実行効率の改善には,並列化・分岐処理・キャッシュ構造など様々な技術的な課題を考慮 する必要がある.VAST 木ではデータ並列度が高い比較処理を行うために分岐ノードには不可逆な圧縮手 法を,葉ノードには CPU に最適化された可逆な圧縮手法をそれぞれ適用する.この設計により索引全体の データ削減を行いながら実行効率の改善も同時に実現する.データ数が 2³⁰の分岐ノードのサイズは既存 手法に対して 95%以上削減,また索引全体のサイズで 47~84%削減した.データ数が 2²⁸の探索では,ス ループット性能は 2 分木に対して 6 倍,既存手法で最も効率的な FAST に対して 1.24 倍の性能向上を確 認した.探索処理中のメモリ帯域の消費量も 2 分木に対して 94.7%, FAST に対して 40.5%削減している.

キーワード:木構造探索, SIMD 命令, 圧縮

Vast-Tree: A Vector-advanced and Compressed Structure for Massive Data Tree Traversal

TAKESHI YAMAMURO^{1,a)} MAKOTO ONIZUKA^{2,b)}

Received: December 21, 2014, Accepted: April 7, 2015

Abstract: We propose a compact and efficient on-memory index structure, called VAST-Tree, for massive data sets. Several indexing techniques such as binary trees and B+trees are widely-used in the database community. Unfortunately, we find that these techniques suffer two major shortcomings when applied to massive data sets; for one thing, their indices are so large that they could overflow regular main memory. For another, they suffer from a variety of penalties (e.g., conditional branches and cache misses). VAST-Tree applies lossy compression to keys in branch nodes and processor-friendly compression to keys in leaf nodes in order to reduce the size of indices. This lossy compression can also improve data parallelism for key comparisons in branch nodes. Compared to current alternatives, VAST-Tree can compress branch nodes by more than 95% and the overall index size by $47 \sim 84\%$ for 2^{30} keys. With 2^{28} keys, VAST-Tree has roughly 6.0-times and 1.24-times throughput and saves the memory consumption by more than 94.7% and 40.5% as compared to binary trees and state-of-the-art FAST, respectively.

Keywords: Tree traversal, SIMD, compression

 NTT ソフトウェアイノベーションセンタ
 NTT Software Innovation Center, Musashino, Tokyo 180– 8585, Japan

² 大阪大学大学院情報科学研究科 Graduate School of Information Science and Technology, Osaka University, Suita, Osaka 565–0871, Japan

^{a)} yamamuro.takeshi@lab.ntt.co.jp

^{b)} onizuka@ist.osaka-u.ac.jp

1. 研究背景

データ集合から指定された条件に合致した部分集合を高 速に探索するため2分木やB+木に代表される索引は基本 的でかつ重要なデータ構造として幅広く活用されている. これらのデータ構造は索引対象のデータ数が多くなった場

- **表1**2分木とB+木の索引サイズ(GiB)比較.括弧内は分岐ノー ドのサイズを示す
- Table 1
 Total sizes (GiB) of binary trees and B+trees. The values in parentheses show only the sizes of branch nodes.



- 図1 索引対象のデータ数が 2²²~2²⁸ における Xeon5260 を用いた 実行時間内訳(命令実行時間,ストール時間,分岐ペナルティ) 比較.括弧内は各実行における IPC を表す
- Fig. 1 Ratios of execution time and counts of instructions for exact-match scans on a Xeon X5260 processor. The values in parenthesis indicate IPC in the processor.

合に索引サイズの肥大化と CPU の実行効率(CPU の1サ イクルあたりの実行命令数)の低下の問題が顕在化する. 表1にデータ数が2²⁸~2³²に対する2分木とB+木の索 引サイズを示す.2分木は広く知られた一般的な方法で実 装を行い,B+木は PostgreSQL v9.0.1の索引構造を用い た.データ数が多くなると索引サイズが10~100 GiB 程度 になり,現在市販されているサーバのメモリサイズと比較 しても非常に大きいことが分かる.索引がメモリサイズを 超えると極端な性能劣化が発生するため,索引のサイズ削 減は重要な課題の1つといえる.

CPU を考慮した設計を行っていない手法の実行効率は 一般的に低く [1],木構造索引に関しても既存研究で分岐 処理やキャッシュ・TLB^{*1}ミスによるストール時間が探索 処理の実行効率を下げていることが報告 [9] されている. 図1にXeon5260を用いたデータ数が 2²²~2²⁸の2分木 の完全一致探索の実行時間内訳を示す.分岐処理によるペ ナルティとストール時間が内訳の約 60%~80%を占めてお り,結果的に実行効率を示す IPC (Instructions Per Cycle) は 0.3 前後であることが分かる.

木構造索引における CPU の実行効率(命令あたりに必要な CPU サイクル数)の課題を取り扱った最新の提案は Kim らによる FAST [12], [13] である. FAST は on-memory で 更新がない処理を前提に, CPU を考慮したデータ構造設 計と探索処理を行うことで分岐処理によるペナルティや キャッシュミスを改善している.また1つの命令で複数 のデータを処理可能(SIMD: Single Instruction Multiple Data)なベクトル命令を活用することで比較処理の高速化 も実現している.しかし,FASTを規模の大きなデータに 適用する場合に以下のような課題が存在する.

- 効率化のためにキャッシュラインやページに対する内 部構造のアラインメント調整をするが、データ規模が 大きくなった場合に多量のパディング領域による索引 サイズの肥大化が発生する。
- SIMD 命令で分岐ノード内の比較キーの同時処理を行うが,最大のデータ並列数が3に制限されている.
- 葉ノードの圧縮が取り扱われていない.

本研究の目的は大規模なデータに対して空間コストが低 く,探索における比較処理の並列数を改善した圧縮木構造 索引である VAST 木を提案する. VAST 木では FAST と 同様に、2次記憶装置を考慮しない on-memory の索引構造 を前提としている. VAST 木は分岐ノードに SIMD 命令で のデータ並列度を改善する不可逆な圧縮を, 葉ノードに復 元処理の際の分岐処理(if-then 命令)の除去やキャッシュ を考慮したデータ配置による CPU 最適化を行った可逆圧 縮をそれぞれ適用する.この設計により索引全体のデータ 削減を行いながら同時に実行効率の改善を実現する.分岐 ノードを複数の領域に分割して,各領域ごとに異なるパラ メータの不可逆圧縮手法を適用する.木構造における分岐 ノードは上部に比べて下部が大きくなるため,下部に位置 する領域ほど圧縮率が高くなるように設計することで、分 岐ノード全体のサイズが小さくなるように構成を行う. ま た圧縮後の比較キーの bit 長を8や16にあわせることで SIMD 命令で同時比較が可能なデータ並列数を改善する. 上記の不可逆な圧縮手法はデータ削減と実行効率改善を同 時に実現する一方, 圧縮前後で正しい順序関係を保持しな い場合がある.具体的には圧縮する前に大小関係のある値 を圧縮後に同値と判定する誤った比較(不正な比較処理) が発生する.結果的にこの不正な比較処理で、本来探索さ れるべき正しい探索位置とは異なる誤った探索位置(探索 のズレ)が返される.この不正な比較処理によって発生す る正しい探索位置からの探索のズレの量を Δe と定義する. VAST 木では分岐ノードの処理後に得られる葉ノード上の 位置から順読み込みを行い,探索のズレの訂正処理に行う. 本研究の貢献は以下である.

- 大規模なデータに対して空間コストが低く,探索における比較処理の並列数を改善した VAST 木を提案する.本手法は木構造の各ノードに適した圧縮手法を適用することで索引サイズを削減しながら,同時にデータ並列度の高い SIMD 比較処理やアライメントの調整によって実行効率も改善する.
- 人工データと Twitter Public Timeline (2010年5月~ 2011年4月)の実データを用いて、VAST木と既存

^{*1} Translation Look-aside Buffer の略,よく参照する論理ページ と物理ページの対応関係をキャッシュする CPU 内の機構である。

手法との圧縮率や性能比較を行った.データ数が 2³⁰ の VAST 木の分岐ノードのサイズは既存手法に対し て 95%以上削減,また索引全体のサイズで 47~84%削 減した.データ数が 2²⁸の探索では,スループット性 能は 2 分木に対して 6 倍,既存手法で最も効率的な FAST に対して 1.24 倍の性能向上を確認した.

- CPU コア数の増加に対して線形的に性能改善を行う ために重要なメモリ帯域の消費量を、データ数が 2²⁸ の探索処理において 2 分木と FAST に対してそれぞれ 94.7%、40.5%だけ削減した。
- VAST 木で発生する探索のズレをモデル化することで Δeの分布を分析する.またこのモデルを活用するこ とで最悪値での訂正処理を 7.1~18.5%改善した.

次の2章では FAST で用いられている SIMD 命令を活 用した分岐ノードの比較処理を概説する.3章で VAST 木 の各ノードに適用する圧縮手法,分岐ノードの不可逆圧縮 を活用した SIMD 命令によるデータ並列度の改善, Δe の 訂正処理など設計に関する詳述を行い,擬似コードをあわ せて示す.4章では VAST 木のプロトタイプを用いた評価 結果を示し,5章で探索中に発生する Δe のモデル化や探 索で発生するメモリ帯域の消費量変化を考察する.最後の 6章と7章で関連研究と結論を述べる.

2. 前提知識:SIMD 命令を用いた木構造探索

FAST における SIMD 命令を用いた木構造探索の概略を 図 2 に示す.まず索引内の比較キーを図中の三角形で示 す部分木にまとめる.この部分木を SIMD ブロックと呼 び,1つの SIMD 命令で同時に処理される最小単位である. FAST では128 bit の SIMD 用レジスタ*2と32 bit のキー を前提に比較処理の同時実行を提案している.探索中の分 岐処理を取り除き CPU 実行効率を高めるために,SIMD 命令の比較結果と移動オフセットの対応表を用いて分岐処 理を積和演算に置き換えて制御依存をデータ依存に変換 する.

図2は以下のように、幅優先順でメモリ上に配置される.

 $[34, 78, 91], [2, 11, 23], [35, 39, 49], [80, 87, 88], \dots$

括弧の整数列が SIMD ブロックを表し,下部線が図 2 中 の太線の SIMD ブロックと対応している. '79' を探索対象 の入力クエリとした場合,まず初めに '79' を SIMD レジス タ A に (79, 79, 79, *) としてロードをする ('*' は任意の 値を表している).次に一番左の SIMD ブロックを SIMD レジスタ B に (34, 78, 91, *) としてロードを行い, SIMD レジスタ A と B の比較処理を行うことで比較結果 (1, 1, 0, *) を SIMD レジスタ C に出力する. SIMD レジスタ A 内の値が大きければ 1 を,そうでなければ 0 を返す.その



- 図 2 SIMD 命令による木構造探索の例.レジスタに SIMD ブロッ ク内の比較キー列を読み込み,同時比較を行う.事前計算した 対応表から,次に比較するべき SIMD ブロックへ移動する
- Fig. 2 An example of tree traversal with SIMD instructions. This instance shows that 128 bit SIMD registers can load four 32 bit integers, and three keys can be compared simultaneously.

後,比較結果と移動オフセット対応表を用いて次に処理するべき SIMD ブロックへ移動する.各 SIMD ブロックは
12B (3 個の 32 bit 整数の比較キー)で,図から比較結果
(1,1,0,*)に対応した値は3である.結果的に処理した
SIMD ブロックの 12B と,比較結果から算出したオフセット 36B を加算して処理中の位置から48B (12B + 12B × 3)だけ位置をずらせばよいことが分かる.

FAST では葉ノードまでの探索に必要な読み込みキャッ シュライン数を削減するために、図2で説明した2段の SIMD ブロックがキャッシュラインに収まるようにブロッ ク化する. SIMD ブロック内の比較キーの数を 3, SIMD ブロックを2段とした場合に比較キー集合の総サイズが (SIMD ブロックが 12B で、5 つの SIMD ブロックがある ため) 60B となり, 結果として一般的な CPU のキャッシュ ラインサイズ (64B) に収まる. この構成により, 2段の SIMD ブロックの探索に必要な読み込みキャッシュライン 数は必ず1となる.また上記と同様に,複数のキャッシュ ラインブロック集合が1つのページに収まるようにブロッ ク化している.しかしこのブロック化により,1回のSIMD 命令による比較キー数が3に制限されている点と、ブロッ ク化の際に発生するパディング(2段の SIMD ブロックを 64B のキャッシュラインに収めた場合の差分 4B) の影響 による索引サイズの肥大化が課題となる.

3. 提案手法: VAST 木

3.1 データ構造設計の概要

VAST 木では分岐ノードと葉ノードにそれぞれ異なる圧 縮手法を適用することで、索引全体のデータ削減を行いな がら同時に実行効率の改善を実現する.分岐ノードには SIMD 命令でのデータ並列度を改善する不可逆な圧縮を, 葉ノードには復元処理の際の分岐処理の除去やキャッシュ

^{*&}lt;sup>2</sup> 128 bit のデータ処理用の SIMD 命令として Intel の CPU では SSE 命令が, AMD の CPU では 3DNOW!が実装されている.



- 図3 VAST 木の概要.上部 P₃₂は FAST [12], [13] と同様の構造 を適用(非圧縮部)して、中部 P₁₆ と下部 P₈には不可逆な圧 縮手法を適用する(圧縮部).葉ノードに対しては可逆の圧縮 手法を適用することで木構造全体のデータサイズ削減を図る
- Fig. 3 An overview of VAST-Tree. The top layer (P_{32}) of VAST-Tree uses FAST techniques, and VAST-Tree compresses the middle and the bottom layers $(P_{16}$ and $P_8)$ of the trees. The heights of these layers are H_{32} , H_{16} and H_8 , respectively. Moreover, keys in leaf nodes (key1, key2, ..., keyN) are compressed by using loss-less compression.

を考慮したデータ配置で実行効率を改善した可逆な圧縮 をそれぞれ適用する.分岐ノード内の比較キーは 32 bit か 64 bit であることを想定して,これ以降は 32 bit を前提に 説明を行う.図3に VAST 木の全体構造を示す.図中の $H_{nbit}, SH_{nbit}, CLH_{nbit}, CH_{nbit}$ はそれぞれの領域の高さ を表しており, nbit はそれぞれの領域に含まれる比較キー のbit 長を示す.木の高さは 2分木相当であり,たとえば 図2の SIMD ブロックの高さは 2 である.木構造の分岐 ノードは上部に比べて下部が大きいため,下部に位置する 領域ほど圧縮率が高くなるように設計する.分岐ノードと 葉ノードの設計をそれぞれ 3.1.1 項と 3.1.2 項で概説する.

3.1.1 分岐ノードの設計概要

FAST と同様に VAST 木も図 3 中の SIMD ブロックを SIMD 命令で同時比較する. SIMD ブロック内の比較キー は下部の領域ほど bit 長が小さいことから,データ並列数 が 3 に限られている FAST とは異なり SIMD 命令で同時 比較可能な数を改善している. VAST 木は具体的に以下 3 つの領域*³で構成される.

- P₃₂, FAST と同様の構造を適用, (2^{SH₃₂} 1) 個の比 較キーを同時処理.
- P₁₆, 32 bit の比較キーを 16 bit に圧縮し, (2^{SH₁₆}-1) 個の比較キーを同時処理.
- P₈, 32 bit の比較キーを 8 bit に圧縮し, (2^{SH₈} 1) 個 の比較キーを同時処理.

 $P_{16} \ge P_8$ 領域の比較キーに "prefix truncation" と "suffix truncation" [2] を適用する. 上位 bit の連続する 0 と,

*3 64 bit の場合は $P_{64}/P_{32}/P_{16}/P_8$ の 4 つの領域で構成する.



図 4 圧縮ブロックと SIMD ブロックのアライメントの再構成 Fig. 4 Alignment of compression blocks and SIMD blocks.

下位 bit を除去することで非可逆な圧縮を行い,0ではな い有意な上位 bit のみを新たな比較キーとする.これによ りハッシュ手法などを用いて比較キーを圧縮する方法とは 異なり,上位 bit を用いたキーの比較を行えるようにして いる.しかし下位 bit を除去しているため,圧縮前後で正 しい順序関係が保持できず探索中に不正な(圧縮前に大小 関係のある値を,圧縮後に同値と誤って判定する)比較処 理が行われる場合がある.この不正な比較処理によって発 生する正しい探索位置からの探索のズレの量を Δe と定義 する.VAST 木では,正しい探索位置からのズレ Δe を探 索処理後に葉ノードを順読み込みすることで訂正を行う. 比較キーの非可逆圧縮に関しては 3.2.1 項で,不正な比較 処理と正しい探索位置からのズレ Δe の具体例と訂正処理 を 3.2.3 項で説明する.

図 3 中の各ブロックのキャッシュラインやページに対 するアライメントの調整は,性能改善のために重要である と先行研究で指摘されている [12]. *P*₃₂ の領域は FAST に 倣い,キャッシュラインブロック内の SIMD ブロックを キャッシュラインの先頭から連続位置に配置し,末尾の SIMD ブロックサイズ未満の領域はパディングで埋める. ページブロック内の領域に関しても同様の再構成を行う. FAST では索引全体に対して上記のアライメントの調整を 行うが,この再構成によって発生するパディング領域の影 響が索引の下部になるほど高くなるため索引サイズが肥大 化する.そこで VAST 木では再構成と索引サイズのトレー ドオフに着眼して,*P*₁₆ と *P*₈ の領域では図 4 に示すよう に各ブロックを SIMD 長のみでアライメントすることでパ ディングの影響を最小化するように設計する.図中の圧縮 ブロックのヘッダに関しては 3.2.1 項で説明を行う.

3.1.2 葉ノードの設計概要

葉ノード上のデータ(key1, key2,..., keyn) は昇順の列 と見なすことができるため, VAST 木では P4Delta [26] を 活用する.葉ノードの探索処理においては任意位置のデー タを高速に参照する必要があるため, k 個ずつのデータ列 を1つのチャンクとして圧縮を行うことで任意位置の要素 の復元を行いやすくしている.葉ノードの圧縮方法は 3.3 節で説明を行う.圧縮により単体キャッシュラインにより 多くの比較キーを格納できることから,探索のズレの訂正

			14bit			圧縮された 8bitの比較キー列
w_1 :	0	0000	0000 00	000 000	0000	0000 0000
w ₂ :	1938	0000	0000 01	11 1011	1111	0001 1110
W ₃ :	3923	0000	0000 11	11 0101	0011	0011 1110
w ₄ :	6281	0000	0001 10	000 1000	1001	0110 0010
 w ₈ :1	15864	0000]	0011 11	01 1111	1000	1111 0111
	8bit (nbit=8) 6bit (v _{rshift} =6)					

図 5 P₈内の8つの値 (w₁, w₂,...,w₈) における不可逆圧縮の例 Fig. 5 An example of the lossy compression on the eight values (w₁, w₂,...,w₈) of P₈.

処理において参照が必要なキャッシュライン数を少なくす ることが可能である.

3.2 分岐ノード構造における再構成

3.2.1 比較キーの不可逆圧縮

VAST 木では B+の分岐ノードの圧縮に用いられる "prefix truncation" と "suffix truncation" を適用する. これら は比較キー間で共通した先頭と末尾の bit 列を除去するこ とで圧縮を行う既存手法である. ここで v_1, v_2, \ldots, v_m (V) を圧縮ブロック内に含まれる昇順の比較キー列とする. 比 較キーの圧縮は以下のように行う.

- (1) V の各値を V 中の最小値 (v_{min}, 昇順であるため v_{min} = v₁ となる) と減算することで W を算出する.
 具体的には W = 0, v₂ - v_{min}, v₃ - v_{min}, ..., v_m - v_{min} となる.
- (2) Wの最大値(w_m)の左端の'1'から nbit 分の範囲に 対応した W内の各値の bit 列を取り出して,圧縮後の 比較キー列とする. nbit の値は P₁₆ = 16, P₈ = 8 で ある.

図5に P_8 の領域での比較キー圧縮の例を示す. v_{min} と 図中の v_{rshift} の値は圧縮ブロック内探索の際に利用する ため,圧縮ブロックのヘッダに記録(図4)する.圧縮ブ ロック内の探索に関しては次の3.2.2項で説明を行う.

nbit の値が高いほど圧縮率が高くなるが,一方で圧縮前後で順序関係を保持しない比較キーの出現確率が高くなる. 圧縮前後で順序関係を保持しない具体的な例は 3.2.3 項で 説明する.

3.2.2 圧縮ブロック内の探索処理

VAST 木における探索処理は圧縮ブロック内の SIMD ブ ロックを比較処理しながら,探索対象のキーを含む葉ノー ドまでのたどるべき圧縮ブロックを算出する.この項では 圧縮ブロック内の探索処理を,続く 3.2.3 項では探索のズ レの訂正処理をそれぞれ説明する.圧縮ブロック内の比較 キー列は不可逆な圧縮手法で変換されているため,そのま までは探索対象の入力クエリ key_q と比較することができ ない.そこで key_q を処理中の圧縮ブロック内の比較キー と同様の変換を行うことで,比較処理を行う.具体的な変 換処理は以下のように行う.

- (1) v_{min} と v_{rshift} を圧縮ブロックのヘッダから取り出す.
- (2) key_q から v_{min} の値を減算して v_{rshift} だけ右シフトす

ることで変換後の入力クエリ *key*[']_q を取得する. 上記のように変換することで取得した *key*[']_q と処理中の SIMD ブロックを 2 章で説明した SIMD 命令による探索処 理に従った方法で比較を行う. 3.4.2 項で *P*₈ の領域におけ る探索処理を行う擬似コードを示す.

3.2.3 **Δ***e* の訂正処理

まず 3.2.1 項の比較キーの不可逆圧縮によって変換前後 で順序関係を保持しない具体例を示す. P_8 の領域におい て比較キー "3220 (2進表記で<u>1100 1001</u> 0100)"が"201 (1100 1001)"に変換された場合を考える.ここで探索対象 の入力クエリが"3219 (2進表記で<u>1100 1001</u> 0101)"であ る場合,変換後の値が"201"となる.このとき圧縮前に大 小関係のある値を圧縮後に同値と判定する誤った比較(不 正な比較処理)が発生する.葉ノード上の正しい探索位置 を pos_q ,不正な比較処理によって発生する pos_q からのズ レの量を Δe と定義する.つまり VAST 木の分岐ノード内 の比較処理を完了した直後の探索位置は $key_{pos_q} + \Delta e$ で表 す.そのため訂正処理は Δe を 0 にして,目的のキーであ る key_{pos_q} を取得することである.3.4.3 項で訂正処理の擬 似コードを示す.

3.2.4 Δeの訂正処理を活用した索引サイズの削減

前の 3.2.3 項の訂正処理の機構を活用することで、VAST 木の分岐ノードをさらに圧縮する.具体的には、分岐ノー ドにおける最下部の SIMD ブロック(図 3 の P_8 領域の最 下端)をすべて取り除く.この再構成により、すべての探 索処理に対して最大 2^{SH_8-1} のズレが発生する.しかし以 下の理由から探索性能への影響は小さいことが予想される.

- *SH*₈の値が小さければ,訂正処理は単体キャッシュライン内の処理で完結できる可能性が高い.
- また葉ノードは圧縮されているため、単体キャッシュ ラインに含まれるデータ数は多い。

予備実験から性能劣化に対する圧縮効果が高いことが判 明したため,次の4章の評価実験ではすべてのパターンに 対して上記を適用する.

3.3 葉ノード構造における再構成

P4Delta [26] を葉ノード上の昇順データ ($key_1, key_2, ..., key_n$) に対して適用する方法を説明する. P4Delta は 差分値を圧縮する方法で, 差分値 d は d[1] = key[1], d[i] = key[i] - key[i - 1] ($1 < i \le n$) と定義する. こ の手法は圧縮する値を "coded" と "exceptions" の 2 つに 分類する. 大半の値 (e.g., 90%以上) が表現可能な bit 長 b を定め, b-bit で格納する値が "coded" である. 一方 2^b 以上の値は "exceptions" として非圧縮で格納する.

図6に葉ノードの圧縮の概要図を示す. 3.1.3項の葉ノー



- 図 6 葉ノードにおける可逆圧縮の概要.上部が圧縮前のキー列で、 下部が圧縮後を表す.任意位置のデータを参照できるように、 k 個ずつの列を1つのチャンクとして圧縮する
- Fig. 6 Compression of keys in leaf nodes. Upper array is a uncompressed list of keys, and lower array is a compressed one. VAST-Tree compresses k consecutive keys into a single chunk.

ドの設計概要で説明したように任意位置のデータを高速に 参照するため、図で示すように k 個ずつの列を1つのチャ ンクとして圧縮を行うことで任意位置の要素の復元を行い やすくしている.それぞれのチャンクは内部の最小値(昇 順であるため先頭の値)と k 個の圧縮されたキー("coded" と "exceptions")が含まれる.最小値はチャンクの先頭に 格納することで、訂正処理の際に不必要な復元処理を行わ ないようにしている.訂正処理中のチャンクの扱いに関し ては、3.4.3 項の訂正処理の擬似コードで説明を行う.

チャンクを適切ににアライメントすることで,訂正処理 中で必要なキャッシュライン数を改善することが可能にな る. k は 2 分法(Bisection Method)を用いて決定する. チャンクのアライメントのサイズとキーの bit 長をそれぞ れ nbit_{algn}, nbit_{key} として, k の決定は以下のように行う.

- (1) k_{left} と k_{right}の初期値をそれぞれ「nbit_{algn}/nbit_{key}」
 と nbit_{algn} とする.
- (2) [(k_{left} + k_{right})/2] (k_{middle}) 個のデータを圧縮した
 場合に、キャッシュラインに収まるかを確認する.
- (3) もし収まれば k_{right} に k_{middle} を, そうでなければ k_{left} に k_{middle} を設定する.
- (4) k_{left} と k_{right} が等しくなるまで(2) と(3) を繰り 返す.
- (5) k_{left} (もしくは k_{right})の値を k とする.

最小値を $[nbit_{algn}/nbit_{key}]$ として圧縮前後でサイズが 変わらない場合を,最大値を $nbit_{algn}$ としてキーの平均 bit 長が1になる場合を想定して k の有効範囲を設定する.

3.4 VAST 木の疑似コード

VAST 木を構築する擬似コードを Algorithm 1 に,分岐 ノード (P_8)を探索する擬似コードを Algorithm 2 に, Δe を訂正処理する擬似コードを Algorithm 3 にそれぞれ示す.

Algorithm 1 VAST 木を構築するための疑似コード

- 1: /* 2: * keys[]: 索引対象のキー列
- 3: * height: 変換済みの索引の高さ
- 4: * nbit: 比較キーの bit 長
- 5: * *H_{nbit}*: P_{nbit} における非圧縮部/圧縮部の高さ
- 6: * CH_{nbit} : P_{nbit} における圧縮ブロックの高さ
- 7: * SH_{nbit} : P_{nbit} における SIMD ブロックの高さ
- 8: * bytes_{output}: 構築された VAST 木
- 9: */
- 10: height = 0
- 11: // keys から P₃₂ 領域の FAST を構築
- 12: build_fast(keys, H_{32} , bytes_{output})
- 13: height = H_{32}
- 14: for *nbit* in {16, 8} do
- 15: for $i \leftarrow 1$ to H_{nbit}/CH_{nbit} do
 - 16: for $j \leftarrow 1$ to 2^{height} do
 - 17: // P_{nbit} における (i, j) 位置の部分木を構成する
 - 18: // 2^{CH_{nbit}} 個の比較キーを抽出
 - 19: $keys_{ext}[] = \text{extract_cb_keys}(keys, CH_{nbit}, i, j)$
 - 20: // 取り出した比較キー列から圧縮ブロックを構築
 - 21: $keys_{cmp}[] = lossy_compress(nbit, keys_{ext}[])$
- 22: $v_{min} = \text{extract}_{v_{min}}(keys_{ext}[])$
- 23: $v_{rshift} = \text{extract}_{v_{rshift}}(keys_{ext}[])$
- 24: build_cb($keys_{cmp}[], v_{min}, v_{rshift}, SH_{nbit}, bytes_{out})$
- 25: end for
- 26: height = height + CH_{nbit}
- 27: end for 28: end for
- _____

3.4.1 VAST 木の構築

VAST 木の構築(Algorithm 1)は、まず P_{32} の領域に おける比較キーの抽出と FAST のデータ構造の構築を行 う(12行目). VAST 木における $P_{16} \ge P_8$ の領域の構築は 15~27行目における内外の2つのループで構成される.外 部ループは H_{nbit}/CH_{nbit} 回実行され、内部ループは処理 中の高さにおける圧縮ブロック数(2^{height})回実行される. 1回の内部のループ内処理(17~24行目)で圧縮ブロック を1つ構築する.処理している圧縮ブロックに含まれる比 較キーを索引対象のキー列 keys[]から抽出(19行目)し、 圧縮ブロックを構築する(21~24行目).

3.4.2 分岐ノードの探索

 P_8 の領域の探索(Algorithm 2)は10~29行目における H_8/CH_8 回実行される外部ループと,21~26行目における CH_8/SH_8 回実行される内部ループで構成される.1回の 内部のループ内処理でSIMD ブロックを,1回の外部ルー プ内処理で圧縮ブロックをそれぞれ1つ探索する.圧縮ブ ロックの探索ではまずヘッダから最小値 v_{min} と右シフト 値 v_{rshift} を取り出し,探索対象の入力クエリ key_q に対し て内部の比較キーと同様の変換を行う(13~15行目).次 の内部ループ処理でのSIMD ブロック内の探索は,変換後 の入力クエリ key'_q と処理中のSIMD ブロックを比較処理 して,次に処理するべきブロックの位置を算出する(23~ 25行目).木の高さと累積データサイズの対応表 $size_{tree}$ [] と、比較結果と移動オフセットの対応表 *lookup*[]を用いる ことで、分岐処理を用いず SIMD 比較命令の結果から積和 演算のみで木構造の探索が可能になる。外部ループ処理が

Algorithm 2 P₈ 探索のための疑似コード

```
1: /*
2: * keyq: 探索対象の入力クエリ
3: * posc: 処理中の圧縮ブロックを示す位置情報
4: * poss: 処理中の SIMD ブロックを示す位置情報
5: * lookup[]: 比較結果と移動オフセットの対応表 (図 2)
6: * sizecb8: P8 における単体圧縮ブロックのサイズ
7: * size<sub>sb8</sub>: P<sub>8</sub> における単体 SIMD ブロックのサイズ
8: * size<sub>tree</sub>[]: 木の高さ(0~31)と累積データサイズの対応表
9: */
10: for i \leftarrow 1 to H_8/CH_8 do
11:
     // 処理中の圧縮ブロック内の比較キーと
12:
     // 同様の変換を key<sub>q</sub> に適用
13:
     v_{min} = \text{get}_{-}v_{min}(pos_c)
14:
     v_{rshift} = \text{get}_{-}v_{rshift}(pos_c)
15:
     key'_q = (key_q - v_{min}) >> v_{rshift}
     // 圧縮ブロック内の SIMD ブロック内の比較キーと
16:
17:
     // key' を比較することで,次に処理するべき
     // 圧縮ブロックの位置(pos<sub>c</sub>)を算出
18:
19:
     pos_s = 0
20:
     cur\_pos_c = pos_c
21:
     for j \leftarrow 1 to CH_8/SH_8 do
22:
        // pos_cの8つの比較キーとkey'_aをベクトル命令で比較
23:
        res = \text{compare}_simd8(key'_q, pos_c)
        pos_s = pos_s \times 2^{SH_8} + \text{lookup}[res]
24:
25:
       pos_c = cur_{-}pos_c + pos_s \times size_{sb} + size_{tree}[SH_8 + j]
26:
     end for
     pos_c = cpos \times 2^{CH_8} + spos
27:
28:
     pos_c = cpos \times size_{cb} + size_{tree}[H_{32} + H_{16} + i]
29: end for
30: // Δe を含んだ探索結果の位置情報を返す
31: return pos_c
```

Algorithm 3 Δe 訂正処理の疑似コード

```
1: /
 2: * keya: 探索対象の入力クエリ
 3: * pos_q: Algorithm 2 の結果 (\Delta e \ \epsilon含んだ位置情報)
 4: * k: 単体チャンク内のキー数
 5: */
 6: loop
     // 各チャンクの先頭値(チャンク内の最小値)を取得
 7:
 8:
     v_{min} = get_v_{min}(\lceil pos_q/k \rceil)
 9:
     pos_q = pos_q - k
10:
     if v_{min} < key_q then
11:
       break
12:
      end if
13: end loop
14: keys_{chunk}[] = decompress_{p4delta}(\lceil pos_q/k \rceil)
15: for i \leftarrow 1 to k do
16:
     if key_q == keys_{chunk}[i] then
        // 探索結果の位置情報 (pos_q - \Delta e) を返す
17:
18:
        return pos_q + i
19:
     end if
20: end for
21: // 未発見のエラーを返す
22: return -1
```

すべて完了した際の pos_c の値が、分岐ノード探索直後の Δe を含んだ葉ノード上の位置である.

3.4.3 **Δ**eの訂正処理

訂正処理 (Algorithm 3) では, 6~13 行目で正しい結果 を含んだチャンクを探索して, 14~22 行目で探索された チャンクを復元して正しい位置情報を返す. 前半の処理で は, チャンク先頭に含まれた値を抽出することで, その値 が入力値 key_q より小さくなるチャンクを探索する. 探索 されたチャンクを復元 (14 行目) して, 入力クエリに対 応したデータが含まれれば正しい位置情報 (16~19 行目) を, そうでなければエラー (22 行目) を返す.

4. 評価実験

評価実験では人工データ('Synthetic')と実データの両 方を用いることで異なるキー分布での評価を行う.人工 データでは $1/\lambda$ のパラメータで決められるポアソン分布に 従うデータを用いた.特に明示しない場合には $1/\lambda$ を 16 に設定して評価を行った.一方実データには Twitter Public Timeline(2010年5月~2011年4月)の 36,068,948 件(約 2^{25})の tweet を用いて,内部に含まれる *IDs* と *Timestamps* の属性を抽出して使用した. *IDs* はユーザを 示す識別子を,*Timestamps* は tweet を投稿した時刻をそ れぞれ表す整数値である.図7に実データの差分値 dの 分布を示す.図から明らかにように,大半のデータが重複 するような偏ったデータであることが分かる.

VAST 木の索引構造における各変数(*SH*₃₂, *CLH*₃₂, *PH*₃₂, *SH*₁₆, *CH*₁₆, *SH*₈, *CH*₈)は評価環境がキャッ シュラインサイズ 64B, ページサイズ 4 KiB (kernel-2.6.18-194の CentOS v5.5), SIMD レジスタ 128 bit であることを 前提に 3.1 節で説明したアライメントの調整を行うために (2, 4, 8, 3, 7, 4, 8)と設定した.本評価は Xeon5670(物 理 6 コアで最大メモリ帯域 31.8 GiB/s)のサーバを用いて 実施した. Xeon5670は 32 KiB の L1 キャッシュ, 256 KiB の L2 キャッシュ, 12 MiB の LLC (Last Level Cache)を



図 7 Twitter Public Timeline (2010 年 5 月~2011 年 4 月)の tweet 情報の IDs と Timestamps の差分値 d の分布

Fig. 7 Distributions of the first 10 d-gaps for realistic data used in our experiments. These data were collected from May, 2010 to Apr., 2011 from the Public Timeline of Twitter. 搭載している. CPU の実行時間内訳は oprofile v0.9.6 を用 いて取得を行った. 評価実施時に Xeon5670 上で oprofile をサポートしていなかったため,実行時間内訳の評価のみ Xeon5260 (物理2コアで最大メモリ帯域21.2 GiB/s)を用 いた. 評価用コードは C++ で記述し,gcc v4.1.2 の "-O3" でコンパイルして評価を行った.

4.1 VAST 木の評価

4.1.1 索引構造の圧縮性能

表 2 に索引の分岐ノードのサイズ比較結果を,表 3 に VAST 木の葉ノードの圧縮率をそれぞれ示す.分岐ノード のサイズはデータ分布に依存しないため,人工データのみ を使用した.表 2 から分かるとおり,すべての条件にお いて VAST 木の圧縮率は高く,データ数が 2^{30} で VAST 木は他の手法と比べて 95%以上の削減を実現している^{*4}. この分岐ノードの高い圧縮率は 3.1.1項で説明した "prefix truncation" と "suffix truncation" による比較キー圧縮と SIMD ブロックによるアライメント調整によるパディング除 去に加えて, 3.2.4項で説明した最下部の SIMD ブロックの 除去により実現している. VAST 木の変数 ($H_{32}/H_{16}/H_8$) の値は索引サイズと 3.2.3項で説明した Δe とのトレード オフを考慮して設定する必要がある.予備実験の結果から データ数が 2^{24} のときは $H_{32} = 8$ と $H_{16} = 6$ が, そうでな

- **表 2** 索引の分岐ノードサイズ (GiB) 比較. VAST 木の括弧内の 値は H₃₂ と H₁₆ をそれぞれ表している
- **Table 2** Total branch node sizes (GiB) of VAST-Tree, binary
trees, and FAST. The values in parentheses show H_{32}
and H_{16} . The branch nodes of VAST-Tree are obviously highly compressed by using the lossy compression.

データ数	2^{24}	2^{26}	2^{28}	2^{30}
VAST 木 (0, 6)	0.00449	0.00449	0.130	0.130
VAST 木 (8, 0)	0.00225	0.0186	0.0186	0.519
VAST 木 (8, 6)	0.00449	0.00449	0.130	0.130
VAST 木 (8, 12)	0.00248	0.00337	0.00337	0.251
FAST	0.252	1.25	1.25	64.3
2 分木	0.156	0.625	2.50	10.0

表 3 VAST 木の葉ノードの圧縮率. 括弧内の値は 3.3 節で説明し た手法で決められた k の値

Table 3 Compression ratios of lea	eaf nodes (VAST-Tree)	
-----------------------------------	-----------------------	--

アライメントサイズ	64B	128B	256B
$1/\lambda = 16$.142(113)	.133(240)	.129(497)
$1/\lambda = 64$.225(71)	.219(151)	.206(311)
IDs	.219(71)	.211(151)	.199(321)
Timestamps	.500(32)	.320(100)	.285(311)

*4 表 2 の値は近似しているため H₃₂ = 0 と H₃₂ = 6 のパターン で圧縮率が同値になっているが実際は多少異なる. 索引の上部 (H₃₂)は下部に対して比重が小さいため,構造変化によるデータ サイズの影響が小さく非常に近い値になっている. い場合は $H_{32} = 8 \ge H_{16} = 12$ が Δe の値が低く探索性能 が高い結果となったため、以降の評価実験ではこれらの値 を使用した.

表3にVAST木の葉ノードの圧縮率とチャンク内に含まれる比較キー数kを示す.葉ノードの圧縮率はデータ分布に依存するため、人工データと実データの両方を用いて評価を行った.偏った分布(図7)の場合、またチャンクのアライメントのサイズを大きくした場合に圧縮率が改善することが分かる.表2の結果とあわせて、VAST木ではデータ数が2³⁰の場合に索引全体のサイズを47~84%削減している.チャンクをキャッシュラインにアライメントした場合に探索性能が最も高くなったため、以降の評価実験ではこの値を使用した.

索引構築時間に関して Kim らによる先行研究 [12] では 64 M 個のデータに対する索引構築時間が 0.1 秒以下とある が,VAST 木では分岐ノード内の比較キー圧縮や,葉ノー ド圧縮の 2 分法を用いた最適化により同数の人工データに 対する索引構築で 100x 以上の構築時間を要する.そのた め FAST が想定している短い構築時間を活用した索引デー タ更新のための再構築を行うことができず,データの更新 が少なく参照量が非常に高い場合にのみ VAST 木の適用は 向いている.

4.1.2 VAST 木の性能評価

図 8 に索引対象のデータ数が 2²⁵ における完全一致探 索の実行時間内訳を示す.2分木と FAST は性能に対する データ分布の影響がないことから,データ数が 2²⁵の人工 データ(Synthetic)を評価に用いた. VAST 木は葉ノード 圧縮の有無で2パターンの評価を行った.図から VAST 木 のすべてのパターンで2分木に対してストール時間と分岐 ペナルティが 72.8%と 50%程度まで改善され、結果として IPC が改善されていることが分かる.葉ノードを圧縮しな い VAST 木 ("VAST-Tree w/o P4Delta") では, 3.3.1 項 で説明した分岐ノード内の比較処理のデータ並列数改善と アライメントの再構成によって,実データを用いた評価で IPC が FAST と比較して改善されている.人工データを 用いた評価では FAST に対して IPC が低くなっているが、 これは △e の訂正処理による読み込みキャッシュライン数 増大の影響だと考えられる. 葉ノードを圧縮した VAST 木 ("VAST-Tree w P4Delta") ではすべてのパターンにおいて Δe の訂正処理の改善, CPU 実行効率の高い P4Delta の復 元処理によって IPC の値が改善されていると考えられる.

図 9 に人工データを,表4 に実データをそれぞれ用い たスループット性能の評価結果を示す.葉ノードの圧縮を した VAST木("VAST-Tree w P4Delta")は索引サイズの 観点で利点はあるが,データ数が 2^{24} と 2^{28} の条件で IPC が高いにもかかわらず FAST に対して性能が劣ることが 分かる.これは Δe の訂正処理改善のための P4Deltaの適 用において,復元処理に必要な実行命令数が増大(図 8)



図 8 索引対象のデータ数が 2²⁵ における完全一致探索の実行時間内訳(命令実行時間,ストー ル時間,分岐ペナルティ)比較. 括弧内は各実行における IPC を表す

Fig. 8 Ratios of execution time and counts of instructions with the exact-match scans of three techniques: binary trees, FAST, and VAST-Tree. The total number of keys is 2²⁵, and the values in parenthesis indicate IPC in each technique.



- 図 9 人工データを用いた完全一致探索のスループット性能比較. E縮を用いた索引のサイズ削減を実現しながら VAST 木は FAST とほぼ同等か,それ以上の性能を実現している
- Fig. 9 Throughput with exact-match scans as determined from synthetic data sets.
- **表 4** Twitter Public Timeline の tweet 情報を用いた完全一致探 索のスループット性能比較 (×10⁶)
- Table 4Throughput $(\times 10^6)$ with exact-match scans by using
realistic data sets.

使用データ	IDs	Timestamps
2 分木	9.05	左と同値
FAST	34.3	左と同値
VAST \bigstar w P4Delta	67.7	31.4
VAST \bigstar w/o P4Delta	78.9	41.1

していることが原因だと考えらえる.一方で,葉ノード の圧縮を行わない VAST 木("VAST-Tree w/o P4Delta") は FAST に対して性能改善を実現している.データ数が 2²⁸の探索では,スループット性能は2分木に対して6倍, FAST に対して 1.24 倍である.表4の実データを用いた スループット性能も図9とほぼ同様の傾向が得られてい る.*IDs*のスループット性能が*Timestamps*よりも高い 理由は,*IDs*のデータの偏りの高さが影響していると考え られる(図7).この項の性能評価の結果により,葉ノード の圧縮の適用は IPC の改善になるがスループット性能の 改善にはならないことが判明した.そのため現実的な実装 においては VAST 木の索引サイズがサーバのメモリに対し



- 図 10 データ数が 2²⁸の探索処理での △e 分布. 括弧内の左値が平 均値を,右値が最悪値をそれぞれ示している. 99%以上の探 索における △e が 100 以下であることが分かる
- Fig. 10 Distributions of the errors caused by the use of our lossy compression. The evaluation used synthetic data with 2^{28} keys. The left values in parenthesis indicate the averaged amounts of errors, and the right values show the worst errors in each condition.



- 図 11 人工データの総数を変化させた場合の Δe 分布. 圧縮率が高 く不正な比較が行われる確率の高い H_8 の値が大きくなるこ とで、 Δe の平均値・最悪値が高くなっている
- Fig. 11 Distributions of the errors for the various numbers of keys in leaf nodes (synthetic data). The error amounts get bigger as the number of keys increases, because the maximum error amounts widen due to the increase of H_8 .

て十分余裕がある場合には性能向上の観点で葉ノードの圧 縮を行わず,そうではない場合には P4Delta を適用して索 引サイズを圧縮する選択が適している.

最後に探索処理中に発生する Δe の分布を図 10 と図 11 に示す. それぞれの図中の括弧に Δe の平均値と最悪値を



図 12 完全一致探索における平均メモリ帯域の消費量(B)比較.
 2²⁸個のデータ探索において2分木に対して94.7%, FAST
 に対して40.5%の削減を達成している

Fig. 12 Comparison of averaged memory bandwidth (B) for a single tree search consumed by all the techniques. VAST-Tree achieves the saving of memory bandwidth consumption as compared to binary trees and FAST because of the conditional branch-free and compressed structure.

あわせて示す. Δe が高くなると CPU の命令数やメモリ帯 域の消費量に影響がでることから非常に重要な指標である が,図 10 の結果から 99%以上の探索における Δe が 100 以下であることが分かる.表 3 の結果とあわせて訂正処理 に必要なキャッシュラインは 1~2 程度で探索性能へのペ ナルティは小さいことが分かる.実データ(図 10)やデー タ数を大きくした場合(図 11)に最悪値が高くなるが,平 均値が 100 以下程度に抑えられていれば図 9 と表 4 の結 果から探索処理の平均性能に関しては影響が小さいと判断 できる.

4.1.3 メモリ帯域の消費量評価

先行研究で CPU コア数の増加に対して線形的に性能改 善するには、メモリ消費量を抑えることが重要であると指 摘されている [12], [16], [20]. 図 12 に完全一致探索にお ける平均メモリ帯域の消費量を示す.2分木は分岐処理に おける投機的なメモリ参照の影響でメモリ帯域の消費量が 高く、FAST と VAST 木は分岐命令を含まないため2分木 に対して値が低い.結果的に2²⁸ 個のデータ探索で2分木 に対して 94.7%、FAST に対して 40.5%の削減を達成して いる.

5. 考察

前章の評価実験の結果から VAST 木の性能はデータ分布 に依存するため、5.1 節でデータ分布と Δe の関係を分析し てモデル化を行い、5.2 節では最悪値 Δe (図 10 と図 11) を考慮した訂正処理の改善方法を提案する.

5.1 *Δe*のモデル化

まず初めに分岐ノード内で発生する不正な比較処理の 回数を分析する. 図 13 で SIMD ブロックと葉ノード上 のキー列の関係を示す. SIMD ブロック内の比較キーを n_k とした場合,比較キー列は昇順であるため $n_k < n_{k+1}$



LANGERSTRATE ON NOT SERVICE STRATE OF NITE STRATE STRATE STRATE

図 13 SIMD ブロックと葉ノード上のキー列の分析. 図中の高さhの ST_A と ST_B の部分木下のキー数は 2^h である

Fig. 13 An analysis of the errors on the ordering of keys under a certain SIMD block. The number of keys under these sub-trees $(ST_A \text{ and } ST_B)$ is assumed to be 2^h .

の関係がある.SIMD ブロックまでの高さを h とした場 合,SIMD ブロック以下には $key_A \sim key_{A+2^h}$ に対応した 比較キーを持つ部分木 ST_A と $key_B \sim key_{B+2^h}$ に対応した 比較キーを持つ部分木 ST_A がそれぞれ存在する.ここで $n_k \leq key_q < n_{k+1}$ の順序関係を持つ探索対象の入力クエ リを key_q を考える.もし図中の SIMD ブロック内で不正 な比較処理が発生して $n_{k+1} \leq skey < n_{k+2}$ と判断されて しまった場合,探索処理は ST_A ではなく ST_B に移動す る. ST_B 内の比較キー集合は必ず key_q 以上の値になるた め,最終的に探索処理は key_B を返す.この分析により分 岐ノード内の不正な比較処理はたかだか 1 回しか発生しな いことが分かる.

上記の分析結果から幾何分布を用いて,最上部から連続 する複数の正しい比較処理と続く1回の不正な比較処理で Δe のモデル化を行う.高さhにおける不正な処理が発生 する確率と探索のズレ量,木の高さをそれぞ n_{ph} , de_h , Hとする.探索のズレの総量を表す確率変数Eは以下の ように計算する.

$$E = \frac{1}{H} \left(\sum_{h=1}^{H} de_h \times p_h \sum_{k=1}^{h} (1 - p_{k-1})^{k-1} \right)$$
(1)

実際に不正な比較処理は P_{32} ではなく P_{16} と P_8 の領域 で発生するため,式(1) は以下のようになる.

$$E = \frac{1}{H_{16} + H_8} \sum_{k=1}^{\frac{H_{16}}{CH_{16}} + \frac{H_8}{CH_8}} e_k$$
(2)

 e_k は高さkにおける各圧縮ブロックにおける探索のズレ 量を表し、以下のように計算する.

$$e_{k} = \sum_{l=1}^{\frac{CH_{16}}{SH_{6}} + \frac{CH_{8}}{SH_{8}}} de_{l} \times p_{k}' \sum_{n=1}^{l} \left(1 - p_{n-1}'\right)^{n-1}$$
(3)

del は高さlにおける SIMD 比較命令によって発生する 探索のズレ量を表し、以下のように計算する.

$$de_{l} = \begin{cases} 2^{H_{16} + H_{8} - l \times SH_{16}} & (l \leq \frac{H_{16}}{CH_{16}}) \\ 2^{H_{8} - (l - \frac{H_{16}}{CH_{16}}) \times SH_{8}} & otherwise \end{cases}$$
(4)

 p'_k は高さkにおける圧縮ブロック内における不正な比較処理の発生確率を表す.圧縮ブロック内の値は同じbit 長で圧縮されているため,不正な比較処理が発生する確率 は同じものとして, p'_k を以下のように計算する.

$$p_k' = \frac{c_k \times r_k}{qr},\tag{5}$$

ここでは qr が葉ノードの範囲 $(key_n - key_1)$, c_k が高 さ k における比較キーの総数, r_k が高さ k の比較キーが圧 縮処理で削られた範囲の量をそれぞれ表す. たとえば圧縮 処理によって下位 7 bit 削除されていれば r_k は 2^7 となる. さらに c_k は以下のように計算する.

$$c_{k} = \begin{cases} 2^{H_{32} + (k-1) \times CH_{16}} & (k \le \frac{H_{16}}{CH_{16}}) \\ 2^{H_{32} + H_{16} + (k - \frac{H_{16}}{CH_{16}} - 1) \times CH_{8}} & otherwise \end{cases}$$
(6)

 r_k は明らかに探索対象のデータ分布に依存する.ここで 入力データの差分値 dを表す確率変数を X とする.高さ kの (P_{nbit} の領域に存在する) 圧縮ブロックに対応した葉 ノード上の比較キー列における dの合計値が 2^{nbit} を超え た場合, r_k が0より大きな値になる.そのため dの合計値 を Y_k とした場合に r_k は以下のように計算する.

$$r_{k} = \begin{cases} 2^{\log_{2}Y_{k}-16} & (k \leq \frac{H_{16}}{CH_{16}} \text{ and } \log_{2}Y_{k} > 16) \\ 0 & (k \leq \frac{H_{16}}{CH_{16}} \text{ and } \log_{2}Y_{k} \leq 16) \\ 2^{\log_{2}Y_{k}-8} & (k > \frac{H_{16}}{CH_{16}} \text{ and } \log_{2}Y_{k} > 8) \\ 0 & (k > \frac{H_{16}}{CH_{16}} \text{ and } \log_{2}Y_{k} \leq 8), \end{cases}$$
(7)

$$Y_{k} = \begin{cases} \sum_{n=1}^{H_{16}+H_{8}+(1-k)\times CH_{16}} X_{n} & (k \leq \frac{H_{16}}{CH_{16}}) \\ \sum_{n=1}^{H_{8}+(1-k-\frac{H_{16}}{CH_{16}})\times CH_{8}} X_{n} & otherwise \end{cases}$$
(8)

図 14 に上記のモデルを用いた Δe の推定値と実値の比 較結果を示す.評価には4節と同様に $1/\lambda$ のパラメータで 決められるポアソン分布に従う人工データを用いた. λ の 値が小さいときと大きいときに多少の乖離はあるが,全体 の傾向は推定できていることが分かる.次の 5.2 節では上 記の結果を用いた最悪値 Δe の訂正処理の改善方法を説明 する.

5.2 訂正処理の改善

訂正処理を順読み込みではなく、 Δe の推定値と2分探 索を用いることで改善する方法を説明する. pos'_q を探索 直後の Δe を含んだ葉ノード上の位置で、正しい位置は pos_q とする. 図 15 で示すように pos'_q の位置から pos_q ま で $range_{\Delta e}$ の範囲の2分探索を行うことで訂正処理を行 う. そのため高い確率で pos_q が $range_{\Delta e}$ 内に含まれるよ うに $range_{\Delta e}$ を設定する必要がある. 前の5.1 節の分析か



図 14 Δe のモデルによる推定値と実値の比較 Fig. 14 Estimation of the errors incurred by VAST-Tree.



図 15 最悪値 △e を考慮した 2 分探索による訂正処理の概要

Fig. 15 An overview of our proposed optimization technique for binary searches so as to minimize the penalty of the worst errors.

表5 最悪値 Δe の平均探索時間比較 (CPU 内の rdtsc 値)

 Table 5
 Hardware-based timestamp counts of a single tree search with the worst error.

訂正手法	順読み込み	改善手法	改善率(%)
IDs	557285	523514	7.1
Timestamps	551297	449340	18.5

ら $range_{\Delta e}$ は以下のように計算する.

$$range_{\Delta e} = t \times \sqrt{V(E)}$$
 (9)

ここでtは変数, V(E)はEの分散をそれぞれ表す. pos_q が $range_{\Delta e}$ 内に含まれる確率はtに対する Chebyshev 不等式を用いて以下のように表される.

$$P(|E - Ex(E)| \ge range_{\Delta e}) \le \frac{1}{t^2}$$
(10)

ここで Ex(E) は E の期待値を表す.上記の結果より pos_q が $range_{\Delta e}$ 内に含まれる確率を考慮して t の値を設定 すればよいことが分かる.上記が確率が 99%以上になるよ うに値を設定して以降の実験を行った.表 5 に実データを 用いて最悪値 Δe の場合の平均探索時間比較を示す.CPU 内のハードウェアカウンタである rdtsc を用いて時間の取 得を行った.結果的に IDs で 7.1%, Timestamps18.5%の 時間がそれぞれ節約できていることが分かる.

6. 関連研究

先行研究 [3] においてデータベース処理のボトルネック が CPU/メモリであることが指摘されて以降,データベー ス技術において広く用いられる圧縮や結合などに代表され る重要な処理を近代的なハードウェア上で最適化するため の手法が数多く提案されている [5], [7], [8], [10], [19], [20]. 近年のハードウェア最適化に関わる研究動向に関しては文献 [15] が詳しい.ここでは特にハードウェア最適化された 索引技術と索引構造の圧縮技術の関連研究を概説する.

6.1 ハードウェア最適化された索引技術

T-tree [14] がメモリ上で最適化された索引構造として提 案された後,キャッシュ構造を考慮した B+木向けの最 適化手法が様々提案された [4], [9], [17], [18]. それらの手 法の主な着眼点は、 キャッシュミスを最小化するための B+木の分岐ノードサイズの検討であった.2つの先行研 究[4], [9] が分岐ノードをキャッシュラインより多少大きく 設定することで B+木の処理性能を改善できるという同様 の指摘を行った. これは分岐ノードをキャッシュラインよ り小さく設定すると、TLB ミスによるペナルティが顕在化 して性能を悪化させることが理由である. また異なる研究 アプローチとして、キャッシュミスによる遅延を隠ぺいす るための "query buffering" の手法も提案 [25] された.こ れは同じ分岐ノードを参照したクエリを一定時間保持して 同時に処理することで、ノードの参照回数を最小化する手 法である.結果的に参照するキャッシュライン数が減り, キャッシュミスによる遅延とメモリ帯域の消費量を改善す る.これらは主にキャッシュ構造に着眼した手法であり, 本研究で対象にしている分岐除去やアライメントの最適化 など他の CPU の実行効率に関わる分析はない.

SIMD 命令を用いた探索処理の効率化手法は 2000 年代 後半以降で提案された手法である [12], [13], [21]. その中 でも Kim らによる FAST は最も効率的で,キャッシュ構 造/アラインメント/分岐除去/データ並列などを用いるこ とで CPU の実行効率を改善している. しかし研究背景で 指摘したように,大規模なデータに対する探索処理に関し て FAST は,データ構造のアライメント調整による索引 サイズ肥大化などの技術的な課題がある. 本研究はハード ウェア最適化と索引サイズにおけるトレードオフに初めて 着眼しており,そのうえで圧縮による索引サイズの削減や SIMD 命令による比較処理のデータ並列数改善を提案した.

より近代的な手法として GPU や x86 互換の Intel[®] MIC (Many Integrated Core) などに代表されるメニーコア環 境における探索処理の改善手法も提案 [11], [22] されてお り, FAST においてもこれらの環境を用いた評価 [13] が行 われている.しかし現在のアーキテクチャ上の制約によ り,ホスト側のデータを GPU 側に転送する処理が必要に なり,その時間が処理全体の 15%から 90%を占めるという 先行研究の報告がある [6].ハードウェアを考慮した最適 化の有無で平均で 24x,最大で 53x の性能差が発生すると いう報告 [1] があり,大規模なデータ処理のおいてはます ますこれらの最適化が重要になることが考えられる.

6.2 索引構造の圧縮技術

索引構造の圧縮に関しては、データベース研究の黎明 期に提案された"prefix/suffix truncation"や"NULL suppression" [2] が有名であるが、これらの手法は索引サイズ削 減のみに着眼した提案である。本提案手法は"prefix/suffix truncation"を分岐ノードに適用することで、索引サイズ を削減しながら同時に CPU の実行効率を改善する点が異 なる.

葉ノードのデータは昇順の列と見なすことで近代的 な CPU に最適化された様々な圧縮手法が適用可能であ る.様々な既存手法の中から,ここでは特に効率的な P4Delta [26] と VSEncoding [23] を紹介する. 3.3 節で説明 した P4Delta は主に転置インデックスの要素技術として広 く用いられている [24].一方で VSEncoding は圧縮の際に 必要なパラメータを動的計画法で最適化する手法で,圧縮 率に優れた手法である.VAST 木の葉ノードにおいては, 探索と Δe の訂正処理で任意位置のデータを高速に参照す る必要があるため,圧縮率と復元性能を損なわずに上記の 要件を満たしやすい P4Delta を採用した.しかし,上記の 要件を満たすことができれば VAST 木と葉ノードに適用す る圧縮手法は独立しているため今後提案されるより高効率 な圧縮手法を容易に適用することが可能である.

7. 結論

本論文では大規模なデータに対して空間コストが低 く、CPUの実行効率が高いon-memoryの索引構造である VAST 木を提案した.VAST 木は木構造内の分岐ノードに SIMD 命令でのデータ並列数向上を可能とする不可逆な圧 縮手法を,葉ノードに CPU に最適化された可逆な圧縮手 法をそれぞれ適用することでデータ削減と実行効率改善を 同時に実現した.分岐ノードの不可逆圧縮は索引サイズの 大幅の削減に,葉ノードの圧縮は VAST 木の探索ズレ Δe の訂正処理の際の CPU 実行効率改善に寄与している.最 後に Δe をモデル化することで訂正処理をさらに改善する 方法を考察して提案を行った.データ構造における空間コ ストを抑えつつ,アルゴリズムの CPU 実行効率を高める ことは規模が増大傾向にあるデータ処理においては今後も 重要な要素技術である.

参考文献

- Kim, C. et al.: Closing the Ninja Performance Gap, Intel Labs Technical Report (2013).
- [2] Bayer, R. and Unterauer, K.: Prefix B-trees, ACM Trans. Database Systems, Vol.2, pp.11–26 (1977).
- [3] Boncz, P.A., Manegold, S. and Kersten, M.L.: Database Architecture Optimized for the New Bottleneck: Memory Access, *Proc. VLDB*, pp.54–65 (1999).
- [4] Chen, S., Gibbons, P.B. and Mowry, T.C.: Improving index performance through prefetching, *SIGMOD Record*, Vol.30, No.2, pp.235–246 (2001).

- [5] Chhugani, J. et al.: Efficient implementation of sorting on multi-core SIMD CPU architecture, *Proc. VLDB*, Vol.1, pp.1313–1324 (2008).
- [6] Fang, W. et al.: Database compression on graphics processors, *Proc. VLDB*, pp.670–680 (2010).
- [7] Govindaraju, N.K., Gray, J., Kumar, R. and Manocha, D.: GPUTeraSort: High Performance Graphics Coprocessor Sorting for Large Database Management, *Proc. SIGMOD*, pp.325–336 (2006).
- [8] Govindaraju, N.K., Lloyd, B., Wang, W., Lin, M. and Manocha, D.: Fast computation of database operations using graphics processors, *Proc. SIGMOD*, pp.215–226 (2004).
- Hankins, R.A. and Patel, J.M.: Efficit of node size on the performance of cache-conscious B+-trees, *Proc. SIG-METRICS* (2003).
- [10] He, B. and Yu, J.X.: High-throughput transaction executions on graphics processors, *Proc. VLDB Endowment*, Vol.4, pp.314–325 (2011).
- [11] Kaldewey, T. et al.: Parallel search on video cards, Proc. HotPar (2009).
- [12] Kim, C. et al.: FAST: Fast architecture sensitive tree search on modern CPUs and GPUs, *Proc. SIGMOD*, pp.339–350 (2010).
- [13] Kim, C. et al.: Designing Fast Architecture Sensitive Tree Search on Modern Multi-Core/Many-Core Processors, ACM Trans. Database Systems, Vol.9, No.4 (2011).
- [14] Lehman, T.J. and Carey, M.J.: A Study of Index Structures for Main Memory Database Management Systems, *Proc. VLDB*, pp.294–303 (1986).
- [15] Manegold, S., Kersten, M.L. and Boncz, P.: Database architecture evolution: Mammals flourished long before dinosaurs became extinct, *Proc. VLDB Endowment*, Vol.2, pp.1648–1653 (2009).
- [16] Matthew, R.: When Multicore Isn't Enough: Trends and the Future for Multi-Multicore Systems, *HPEC* (2008).
- [17] Rao, J. and Ross, K.A.: Cache Conscious Indexing for Decision-Support in Main Memory, *Proc. VLDB*, pp.78– 89 (1999).
- [18] Rao, J. and Ross, K.A.: Making B+-trees cache conscious in main memory, *Proc. SIGMOD*, pp.475–486 (2000).
- [19] Satish, N., Harris, M. and Garland, M.: Designing efficient sorting algorithms for manycore GPUs, *Proc. IPDPS*, pp.1–10 (2009).
- [20] Satish, N. et al.: Fast sort on CPUs and GPUs: A case for bandwidth oblivious SIMD sort, *Proc. SIGMOD*, pp.351–362 (2010).
- [21] Schlegel, B. et al.: k-ary search on modern processors, Proc. DaMoN, pp.52–60 (2009).
- [22] Sewall, J. et al.: PALM: Parallel Architecture-Friendly Latch-Free Modification to B+Trees on Many-Core Processors, *Proc. VLDB* (2011).
- [23] Silvestri, F. and Venturini, R.: VSEncoding: Efficient coding and fast decoding of integer lists via dynamic programming, *Proc. CIKM* (2010).
- [24] Yan, H., Ding, S. and Suel, T.: Compressing term positions in web indexes, *Proc. SIGIR*, pp.147–154 (2009).
- [25] Zhou, J. and Ross, K.A.: Buffering accesses to memoryresident index structures, *Proc. VLDB*, pp.405–416 (2003).
- [26] Zukowski, M., Heman, S., Nes, N. and Boncz, P.: Super-Scalar RAM-CPU Cache Compression, *Proc. ICDE*, pp.59–71 (2006).



山室 健

NTT ソフトウェアイノベーションセ ンタ研究員.2008年上智大学大学院 理工学研究科博士前期課程修了,修 士(工学).DBMSのコア技術,およ びハードウェア特性を考慮した探索/ 圧縮アルゴリズムの研究に従事.日本

データベース学会会員.



鬼塚 真 (正会員)

1991年東京工業大学工学部情報工学 科卒業.同年日本電信電話株式会社 入社.2000~2001年ワシントン大学 客員研究員,2010~2014年日本電信 電話株式会社特別研究員,2012~2014 年電気通信大学客員教授,現在,大阪

大学大学院情報科学研究科教授.博士(工学).大規模グ ラフデータの分散データ処理に関する研究開発に取り組 んでいる.2004年情報処理学会山下記念賞,2008年デー タベース学会上林奨励賞等受賞.電子情報通信学会,日本 データベース学会,ACM 各会員.

(担当編集委員 堀井 洋)