

CS 764: Topics in Database Management Systems Lecture 15: Adaptive Radix Tree

Xiangyao Yu 10/23/2025

Announcement

Midterm exam

- All papers before the exam are included
- Guest lectures are **not** included
- Nov. 4 (Tuesday) noon-Nov. 6 (Thursday) noon, central time
- The exam questions will be posted on Piazza (as a word document)
- Please use Piazza to ask questions privately

Suggested ways to submit your solutions:

- Directly type your solutions in this MS word document, or
- Print out the exam and submit a photocopy of your solutions, or
- Convert the exam into a pdf file and write your solutions on it

Today's Paper: B-tree Locking

The Adaptive Radix Tree: ARTful Indexing for Main-Memory Databases

Viktor Leis, Alfons Kemper, Thomas Neumann

Fakultät für Informatik Technische Universität München Boltzmannstrae 3, D-85748 Garching <lastname>@in.tum.de

Abstract-Main memory capacities have grown up to a point where most databases fit into RAM. For main-memory database systems, index structure performance is a critical hottleneck. Traditional in-memory data structures like balanced binary search trees are not efficient on modern hardware, because they do not optimally utilize on-CPU caches. Hash tables, also often used for main-memory indexes, are fast but only support point

To overcome these shortcomings, we present ART, an adaptive radix tree (trie) for efficient indexing in main memory. Its lookup performance surpasses highly tuned, read-only search trees, while supporting very efficient insertions and deletions as well. At the same time, ART is very space efficient and solves the problem of excessive worst-case space consumption, which plagues most radix trees, by adaptively choosing compact and efficient data structures for internal nodes. Even though ART's performance is comparable to hash tables, it maintains the data in sorted order, which enables additional operations like range scan and prefix lookup.

I. INTRODUCTION

transactional databases fit into RAM. When most data is use data level parallelism to perform multiple comparisons cached, traditional database systems are CPU bound because simultaneously with Singe Instruction Multiple Data (SIMD) they spend considerable effort to avoid disk accesses. This instructions. Additionally, FAST uses a data layout which has led to very intense research and commercial activities in avoids cache misses by optimally utilizing cache lines and main-memory database systems like H-Store/VoltDB [1], SAP the Translation Lookaside Buffer (TLB). While these opti-HANA [2], and HyPer [3]. These systems are optimized for mizations improve search performance, both data structures the new hardware landscape and are therefore much faster. Our cannot support incremental updates. For an OLTP database system HyPer, for example, compiles transactions to machine system which necessitates continuous insertions, updates, and code and gets rid of buffer management, locking, and latching deletions, an obvious solution is a differential file (delta) overhead. For OLTP workloads, the resulting execution plans mechanism, which, however, will result in additional costs. are often sequences of index operations. Therefore, index Hash tables are another popular main-memory data strucefficiency is the decisive performance factor.

an in-memory indexing structure. Unfortunately, the dramatic therefore much faster in main memory. Nevertheless, hash processor architecture changes have rendered T-trees, like all tables are less commonly used as database indexes. One reason traditional binary search trees, inefficient on modern hardware. is that hash tables scatter the keys randomly, and therefore only The reason is that the ever growing CPU cache sizes and support point queries. Another problem is that most hash tables the diverging main memory speed have made the underlying do not handle growth gracefully, but require expensive reorassumption of uniform memory access time obsolete. B^+ -tree ganization upon overflow with O(n) complexity. Therefore, variants like the cache sensitive B+-tree [5] have more cache- current systems face the unfortunate trade-off between fast friendly memory access patterns, but require more expensive hash tables that only allow point queries and fully-featured, update operations. Furthermore, the efficiency of both binary but relatively slow, search trees. Because the result of comparisons cannot be predicted easily, prefix tree, and digital search tree, is illustrated in Figure 1.

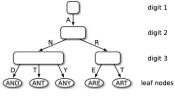


Fig. 1. Adaptively sized nodes in our radix tree

the long pipelines of modern CPUs stall, which causes additional latencies after every second comparison (on average).

These problems of traditional search trees were tackled by recent research on data structures specifically designed to be efficient on modern hardware architectures. The k-ary search After decades of rising main memory capacities, even large tree [6] and the Fast Architecture Sensitive Tree (FAST) [7]

ture. In contrast to search trees, which have $O(\log n)$ access More than 25 years ago, the T-tree [4] was proposed as time, hash tables have expected O(1) access time and are

and B+-trees suffers from another feature of modern CPUs: A third class of data structures, known as trie, radix tree,

Outline

B-tree vs. Trie

Adaptive Radix Tree

- Adaptive types
- Collapsing inner nodes
- Search and insert operations

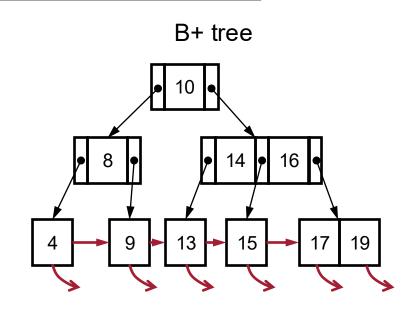
Evaluation

B+ Tree Revisit

Modern indexes fit in main memory

Keys are stored in each level of the tree

Must always traverse to the leaf node to check existence (e.g., cannot stop at an inner node)

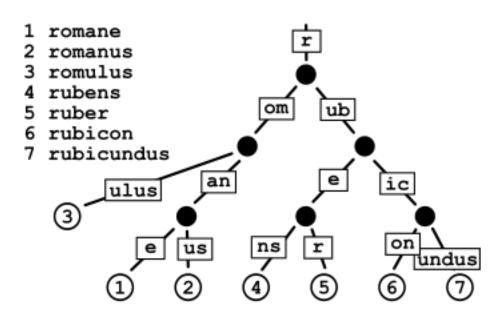


Trie (aka. digital tree or prefix tree)

Path to leaf node represents key of the leaf

Operation complexity is O(k) where k is the length of the key

Keys are most often strings and each node contains characters

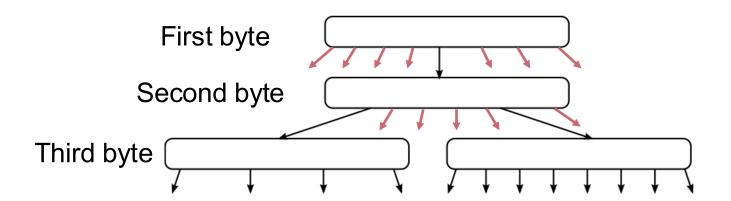


Source: https://en.wikipedia.org/wiki/Radix_tree

Static Radix Tree

Span (s): The number of bits within the key used to determine the next child

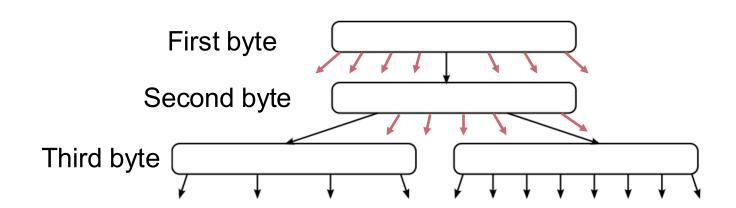
An inner node is an array of 2^s pointers



Static Radix Tree

Span (s): The number of bits within the key used to determine the next child

An inner node is an array of 2^s pointers



Example:

k = 32 bit keys

Consider the index space consumption to insert one key

Extra space $\approx \left[\frac{k}{s}\right]$ (# of levels) \times 2^s (node size per level)

Static Radix Tree

Span (s): The number of bits within the key used to determine the next child

An inner node is an array of 2^s pointers

Large span

- => reduced height
- => exponential tree size

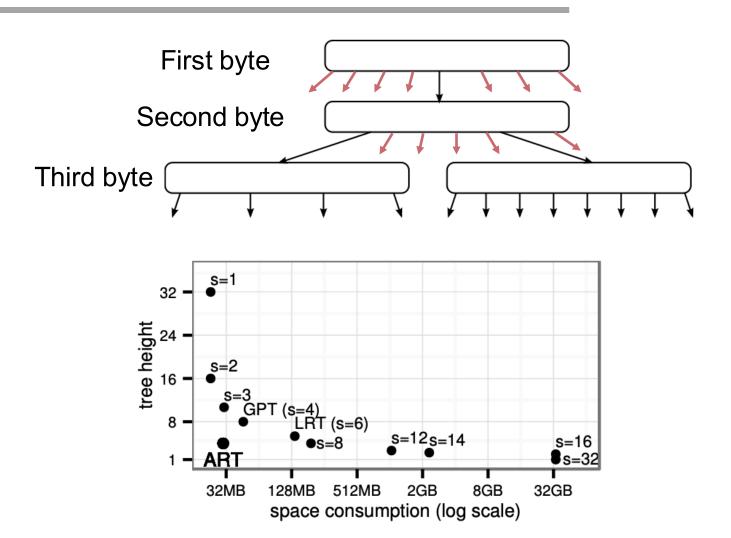
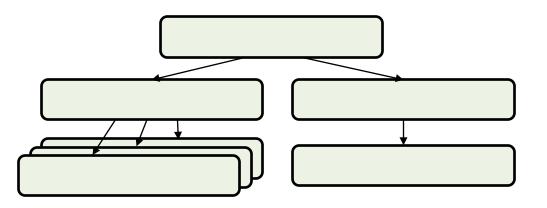
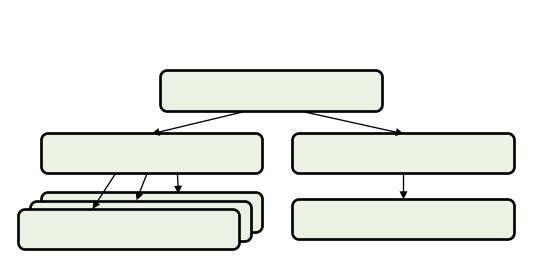


Fig. 3. Tree height and space consumption for different values of the span parameter s when storing 1M uniformly distributed 32 bit integers. Pointers are 8 byte long and nodes are expanded lazily.



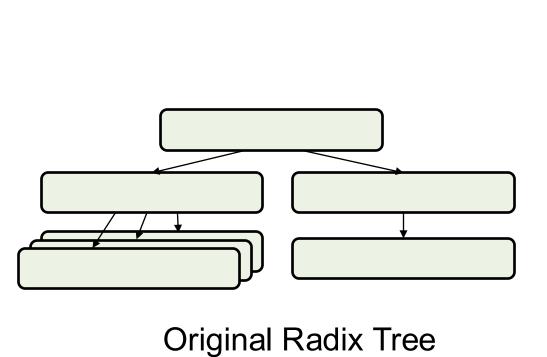
Original Radix Tree

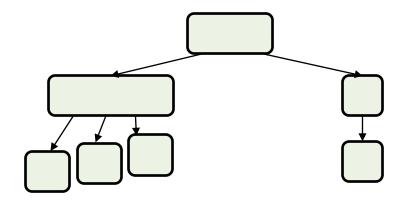


Original Radix Tree

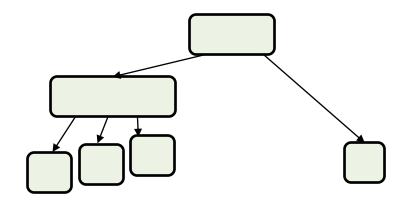
Optimization 1: adaptive node type

Key idea: Use a small node type when only a small number of children pointers exist

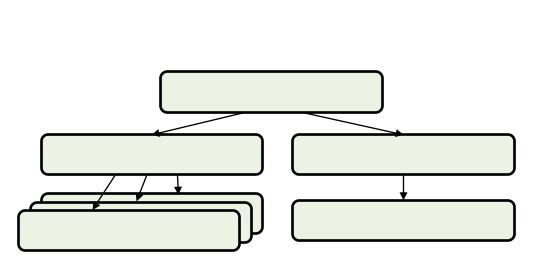




Optimization 1: adaptive node type

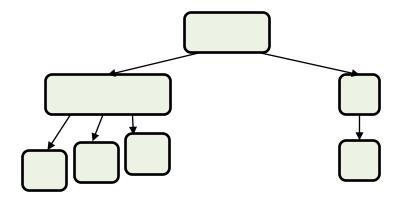


Optimization 2: collapsing inner nodes

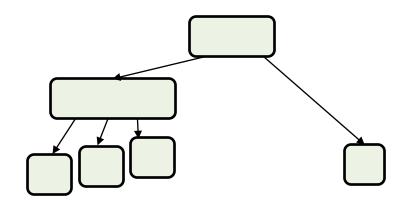


Original Radix Tree





Optimization 1: adaptive node type



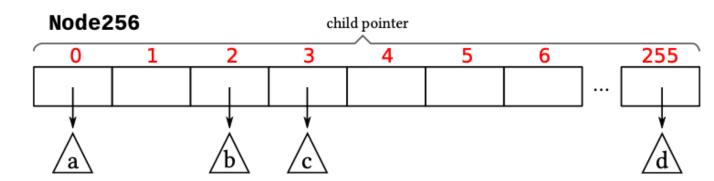
Optimization 2: collapsing inner nodes

Node4 and Node16

Node48

Node256

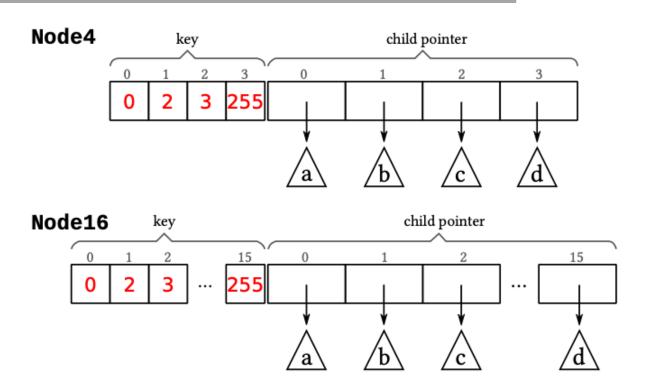
- 256 child pointers indexed with partial key byte directly
- (Same as original radix tree)
- Used for 49–256 entries



Node4 and Node16

- Store up to 4 (16) partial keys and the corresponding pointers
- Partial keys are sorted
- Use SIMD instructions to accelerate key search

Node48 Node256

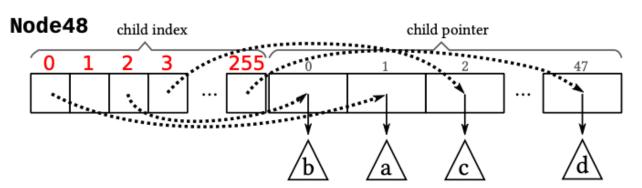


Node4 and Node16

Node48

- 256 entries indexed with partial key byte directly
- Each entry stores a one-byte index to a child pointer array
- Child pointer array contains 48 pointers to children nodes

Node256

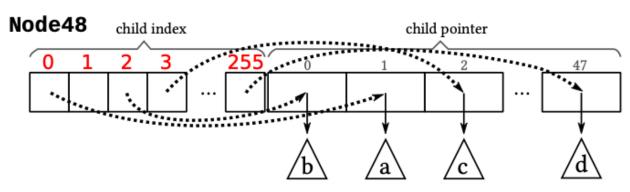


Node4 and Node16

Node48

- 256 entries indexed with partial key byte directly
- Each entry stores a one-byte index to a child pointer array
- Child pointer array contains 48 pointers to children nodes

Node256



Discussion Question

Q1: Is Node48 more space efficient compared to Node4/16 layout?

Q2: What is the key advantage of Node48 layout?

Lazy expansion: remove path to single leaf

- Inner nodes created only required to distinguish at least two leaf nodes
- In the example, root can directly point to leaf FOO, eliminating the two inner nodes
- Requires the key to be stored at the leaf or in the database

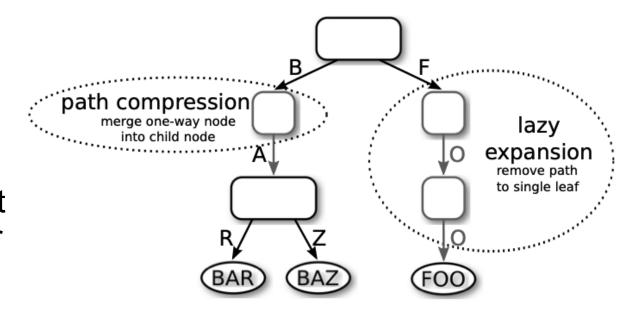


Fig. 6. Illustration of lazy expansion and path compression.

Lazy expansion: remove path to single leaf

Path compression: merge one-way node into child node

 Removes all inner nodes that have only a single child

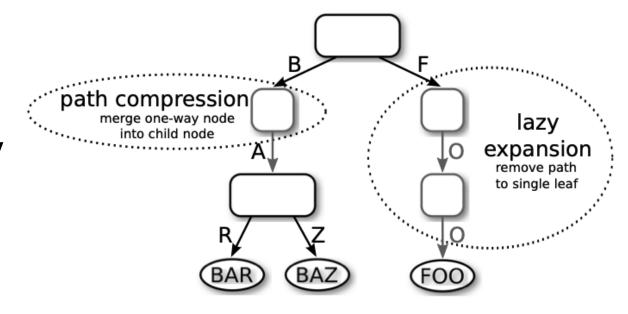


Fig. 6. Illustration of lazy expansion and path compression.

Lazy expansion: remove path to single leaf

Path compression: merge one-way node into child node

- Removes all inner nodes that have only a single child
- Pessimistic: child node stores the compressed partial key

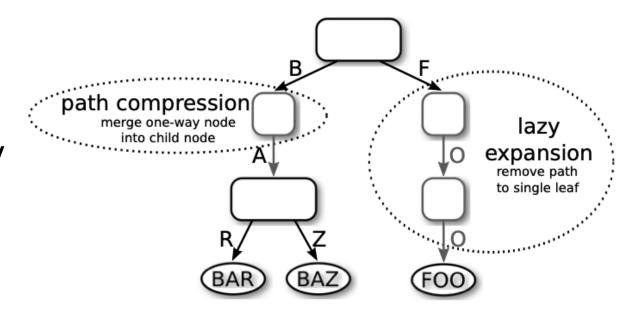


Fig. 6. Illustration of lazy expansion and path compression.

Lazy expansion: remove path to single leaf

Path compression: merge one-way node into child node

- Removes all inner nodes that have only a single child
- Pessimistic: child node stores the compressed partial key
- Optimistic: child node stores only the length of compressed partial key

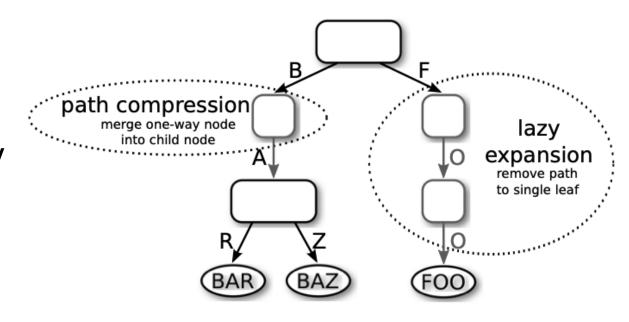


Fig. 6. Illustration of lazy expansion and path compression.

Lazy expansion: remove path to single leaf

Path compression: merge one-way node into child node

- Removes all inner nodes that have only a single child
- Pessimistic: child node stores the compressed partial key
- Optimistic: child node stores only the length of compressed partial key
- Hybrid: use constant vector to store partial key, switch to optimistic approach if the vector overflows

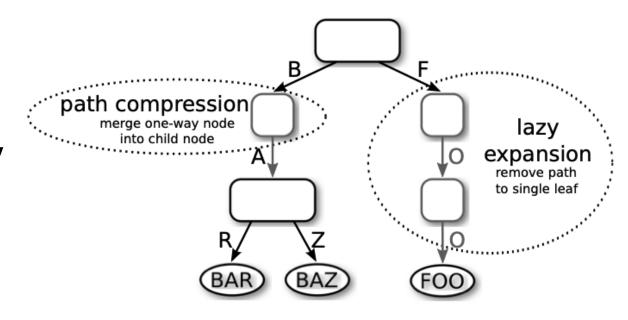


Fig. 6. Illustration of lazy expansion and path compression.

Search Algorithm

```
search (node, key, depth)

if node==NULL

return NULL

if isLeaf(node)

if leafMatches(node, key, depth)

return node

return NULL

if checkPrefix(node, key, depth)!=node.prefixLen

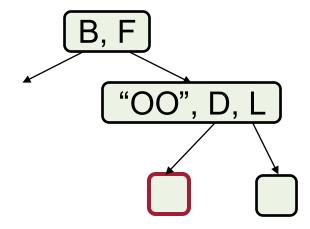
return NULL

depth=depth+node.prefixLen

next=findChild(node, key[depth])

return search(next, key, depth+1)
```

Fig. 7. Search algorithm.



Example: search for **FOOD**

ART requires at most 52 bytes of memory to index a key

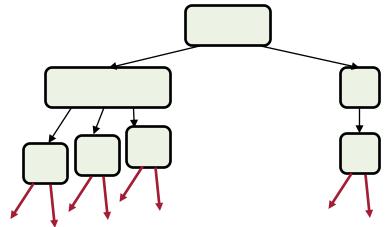
 Assume each child node has 52 byte budget, show that each node will have at least 52 bytes budget remain



ART requires at most 52 bytes of memory to index a key

 Assume each child node has 52 byte budget, show that each node will have at least 52 bytes budget remain

$$b(n) = \begin{cases} x, & \text{isLeaf}(n) \\ \left(\sum_{i \in c(n)} b(i)\right) - s(n), & \text{else.} \end{cases}$$

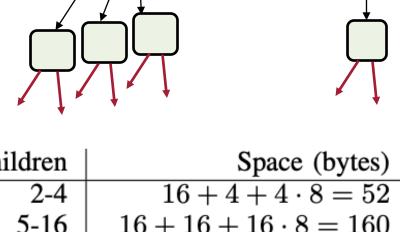


ART requires at most 52 bytes of memory to index a key

 Assume each child node has 52 byte budget, show that each node will have at least 52 bytes budget remain

$$b(n) = \begin{cases} x, & \text{isLeaf}(n) \\ \left(\sum_{i \in c(n)} b(i)\right) - s(n), & \text{else.} \end{cases}$$

Node4: $b(n) \le 52 * 2 - 52 = 52$



Type	Children	Space (bytes)
Node4	2-4	$16 + 4 + 4 \cdot 8 = 52$
Node16	5-16	$16 + 16 + 16 \cdot 8 = 160$
Node48	17-48	$16 + 256 + 48 \cdot 8 = 656$
Node256	49-256	$16 + 256 \cdot 8 = 2064$

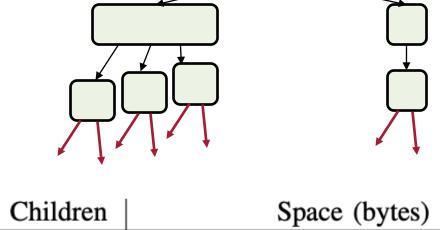
ART requires at most 52 bytes of memory to index a key

 Assume each child node has 52 byte budget, show that each node will have at least 52 bytes budget remain

$$b(n) = \begin{cases} x, & \text{isLeaf}(n) \\ \left(\sum_{i \in c(n)} b(i)\right) - s(n), & \text{else.} \end{cases}$$

Node4: $b(n) \le 52 * 2 - 52 = 52$

Node16: $b(n) \le 52 * 5 - 160 = 100$



Type	Children	Space (bytes)
Node4	2-4	$16 + 4 + 4 \cdot 8 = 52$
Node16	5-16	$16 + 16 + 16 \cdot 8 = 160$
Node48	17-48	$16 + 256 + 48 \cdot 8 = 656$
Node256	49-256	$16 + 256 \cdot 8 = 2064$

ART requires at most 52 bytes of memory to index a key

 Assume each child node has 52 byte budget, show that each node will have at least 52 bytes budget remain

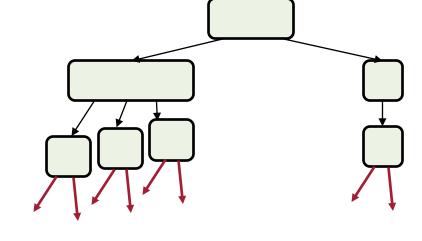
$$b(n) = \begin{cases} x, & \text{isLeaf}(n) \\ \left(\sum_{i \in c(n)} b(i)\right) - s(n), & \text{else.} \end{cases}$$

Node4: $b(n) \le 52 * 2 - 52 = 52$

Node16: $b(n) \le 52 * 5 - 160 = 100$

Node48: $b(n) \le 52 * 17 - 656 = 228$

Node256: $b(n) \le 52 * 49 - 2064 = 484$



Type	Children	Space (bytes)
Node4	2-4	$16 + 4 + 4 \cdot 8 = 52$
Node16	5-16	$16 + 16 + 16 \cdot 8 = 160$
Node48	17-48	$16 + 256 + 48 \cdot 8 = 656$
Node256	49-256	$16 + 256 \cdot 8 = 2064$

Discussion

Space consumption

- ART requires at most 52 bytes of memory to index a key
- Q: What if the key itself is larger than 52 bytes?

Discussion

Space consumption

- ART requires at most 52 bytes of memory to index a key
- Q: What if the key itself is larger than 52 bytes?

Binary comparable keys

- If only binary-comparable keys are used as keys of a radix tree, the data is stored in sorted order and all operations that rely on this order can be supported (e.g., range scan)
- For finite and totally ordered domains, always possible to transform values to binary-comparable keys

Evaluation—Single-Threaded Lookup

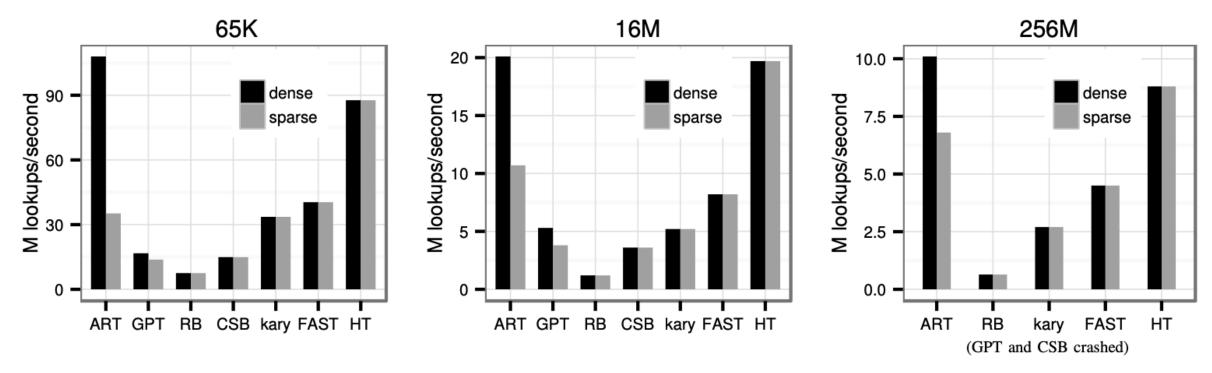


Fig. 10. Single-threaded lookup throughput in an index with 65K, 16M, and 256M keys.

Evaluation—Single-Threaded Insert

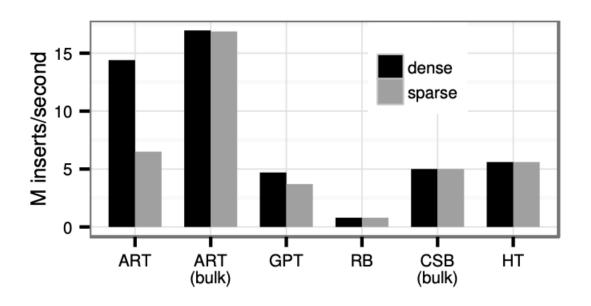
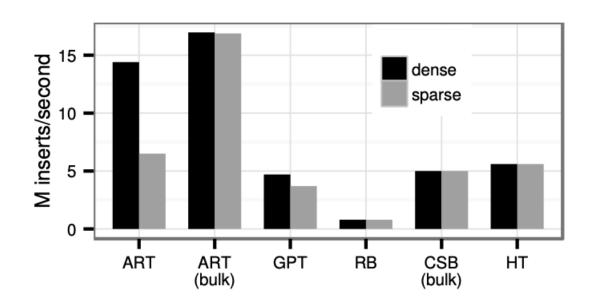


Fig. 14. Insertion of 16M keys into an empty index structure.

Evaluation—Single-Threaded Insert



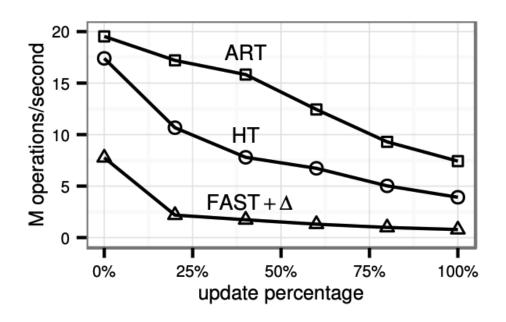


Fig. 14. Insertion of 16M keys into an empty index structure.

Fig. 15. Mix of lookups, insertions, and deletions (16M keys).

Evaluation – More Baselines

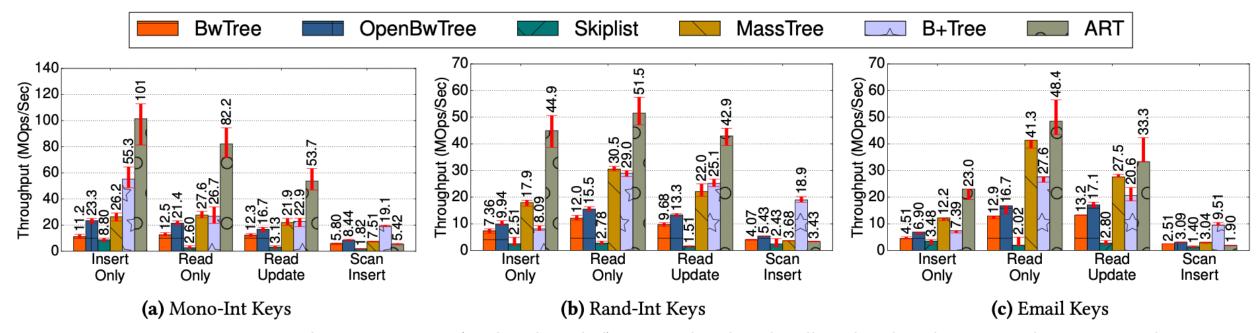
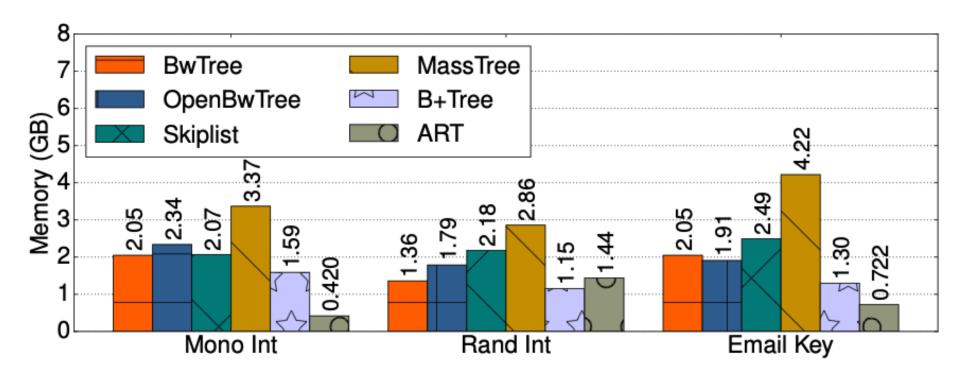


Figure 14: In-Memory Index Comparison (Multi-Threaded) – 20 worker threads. All worker threads are pinned to NUMA node 0.

^{*} Wang, Ziqi, et al. Building a bw-tree takes more than just buzz words. SIGMOD 2018

Evaluation – Memory Usage



(b) Multi-Threaded – Read/Update

^{*} Wang, Ziqi, et al. Building a bw-tree takes more than just buzz words. SIGMOD 2018

Q/A – Adaptive Radix Tree

Node type transition overhead for write-heavy workload?

Adaptive node selection threshold through learning/tuning?

How to make ART multi-threaded? Lock-free?

ART for persistent memory with durability?

Any real-world implementation?

– DuckDB: https://duckdb.org/2022/07/27/art-storage

Next Lecture

Jim Gray, et al., <u>Granularity of Locks and Degrees of Consistency in a Shared Data Base</u>. Modelling in Data Base Management Systems, 1976