

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

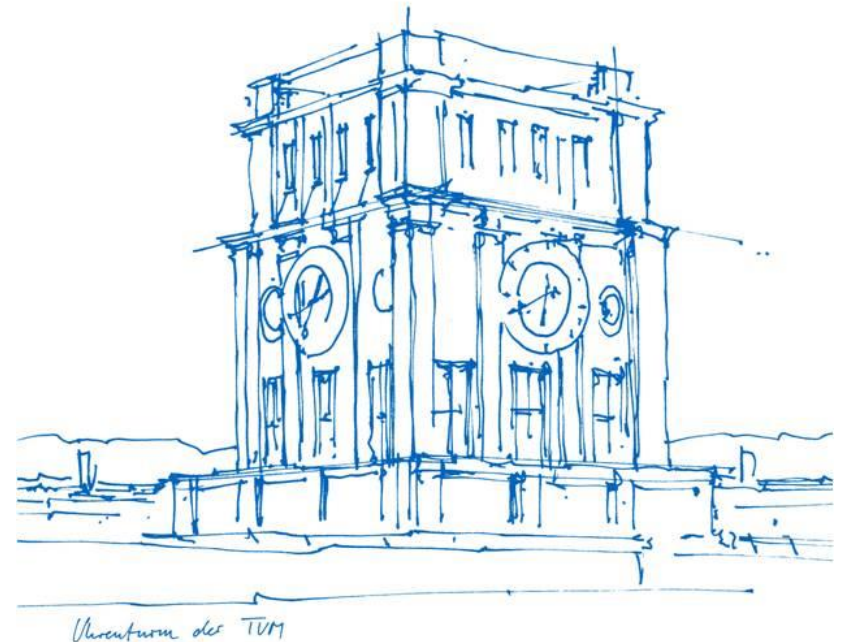
Jonas Fritsch

Technische Universität München

Department of Informatics

Chair for database systems

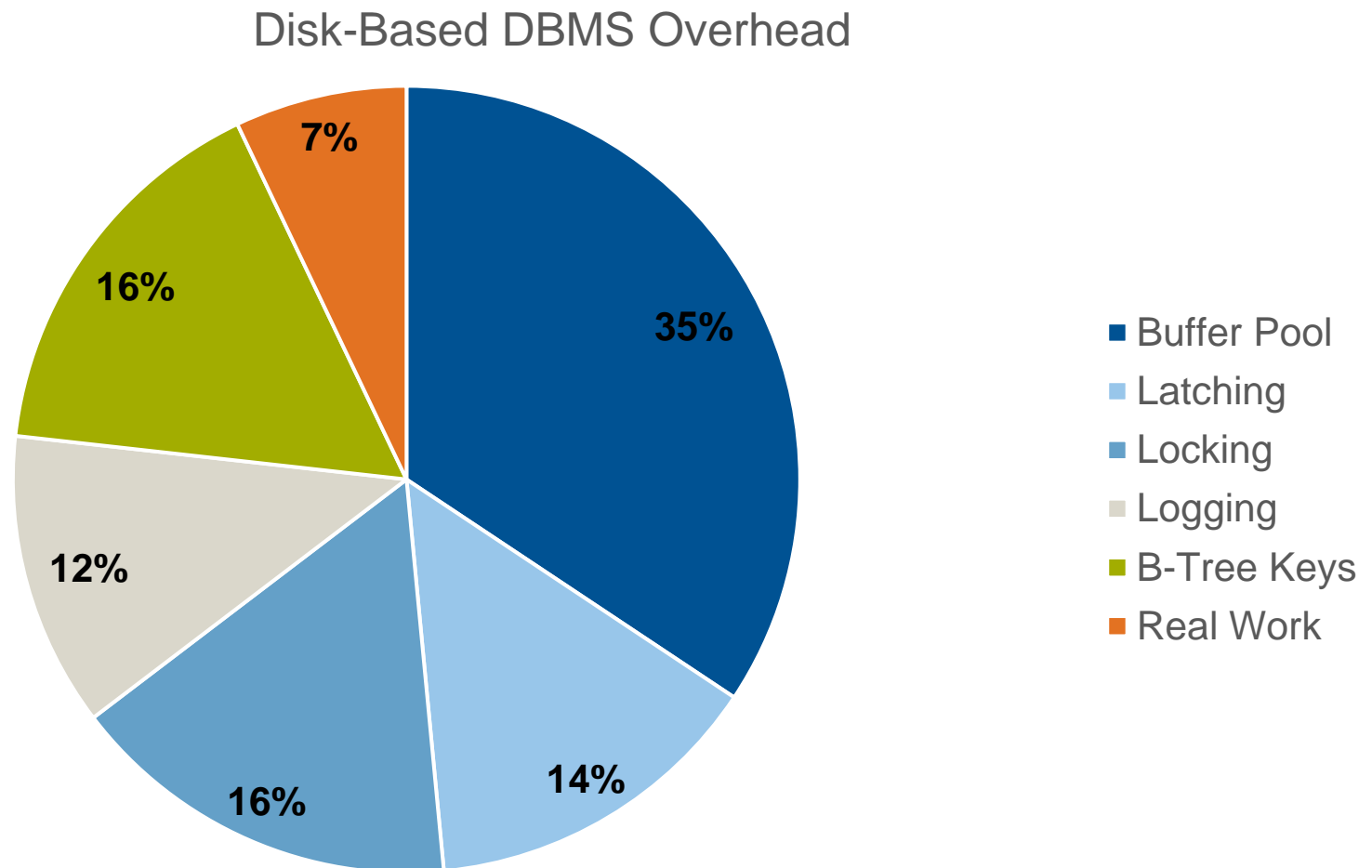
Munich, 11. July 2022



Overview

- Disk-Based vs. Main-Memory DBMS
- Index-Structures in DBMS
- From Tries to ART
- Key Transformations for Bitwise Comparisons
- Benchmarks
- Summary & Conclusion

Disk-Based vs Main-Memory DBMS



Source: [OLTP through the looking glass, and what we found there. SIGMOD, 2008](#)

Index-Structures in DBMS

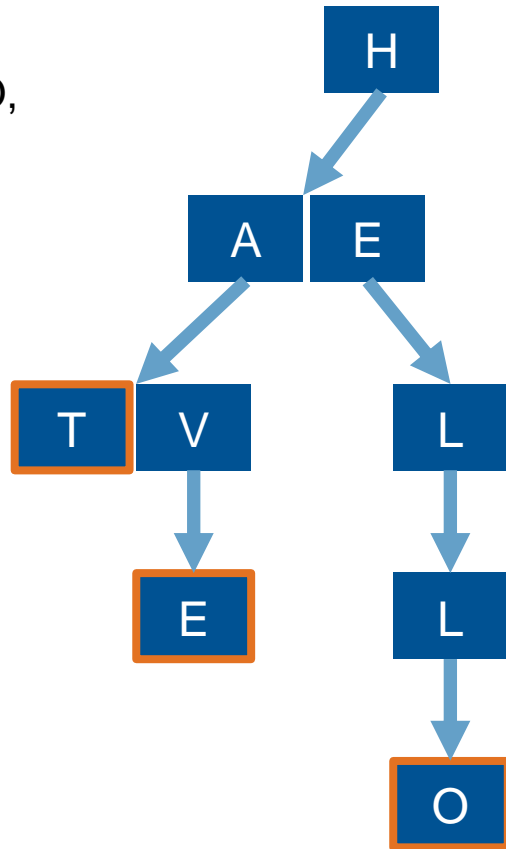
Data structures used to quickly find data in table via a key

2 Index-Types:

- Order Preserving Indexes (Maintains keys in sorted order)
- Hashing Indexes (Associative array mapping hash of key to data record. Only supports equality predicates!)

Tries

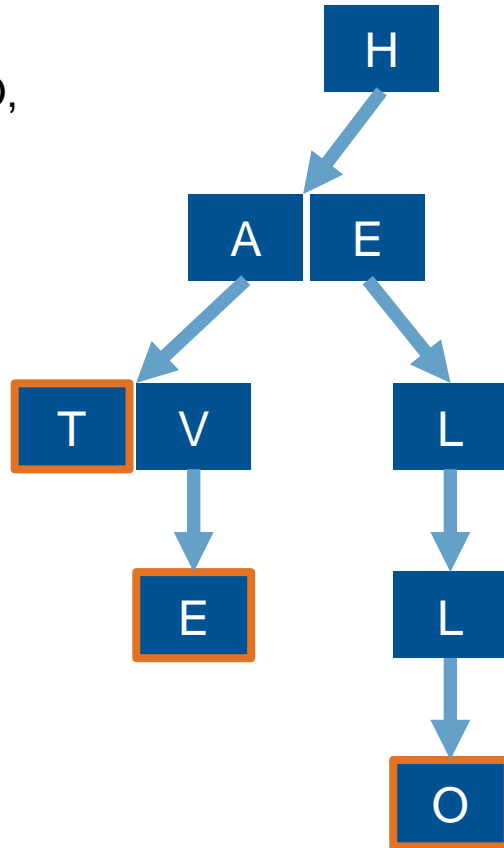
Keys:
HELLO,
HAT,
HAVE



Tries – Properties

Keys:

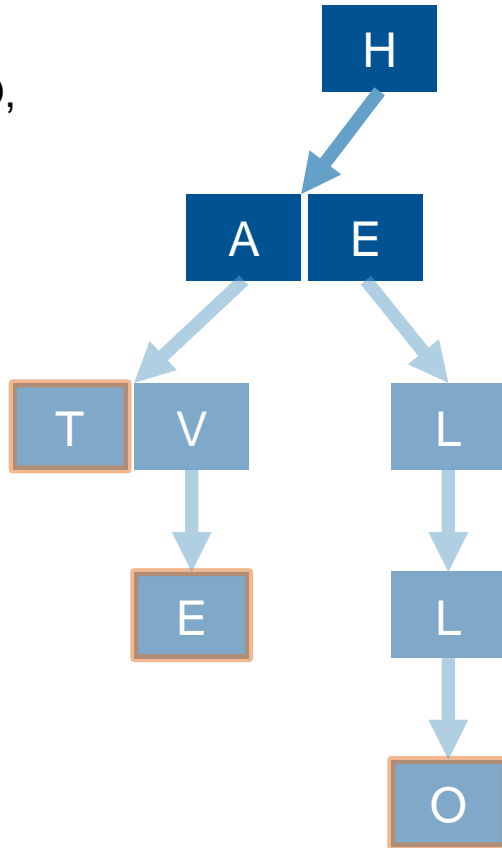
HELLO,
HAT,
HAVE

**Properties:**

- Height/complexity depends on key length k instead of number of elements
- Require no rebalancing
- All insertion orders results in same tree
- Keys stored in lexicographic order
- Keys are stored implicitly along paths

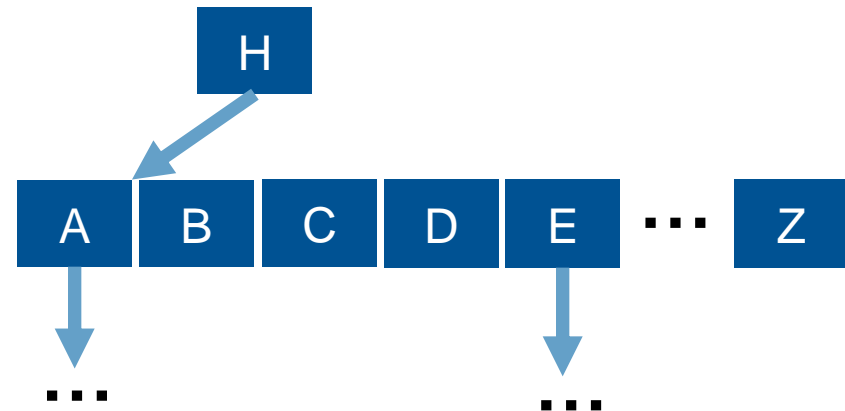
Tries – Implementation

Keys:
HELLO,
HAT,
HAVE



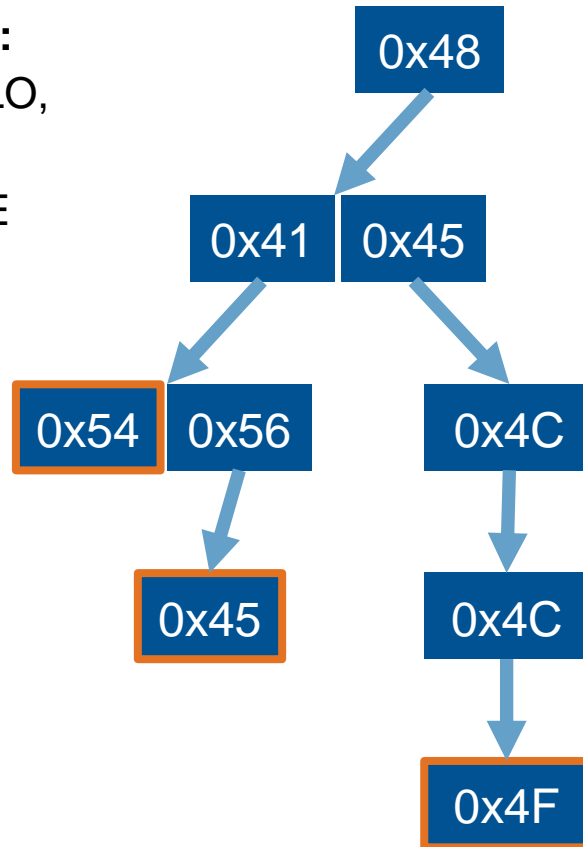
```

class Node {
    bool is_leaf;
    // 'A' to 'Z'
    Node* children[26];
}
  
```



Tries – Key Span, Fanout and Height

Keys:
HELLO,
HAT,
HAVE



Key Span s :

Number of bits each partial key represents (e.g. char is 8 bit span)

Fanout:

Number of children for a node

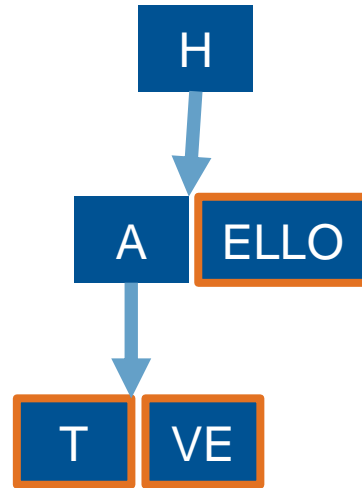
Determined by specific implementation but in general is 2^s

Height:

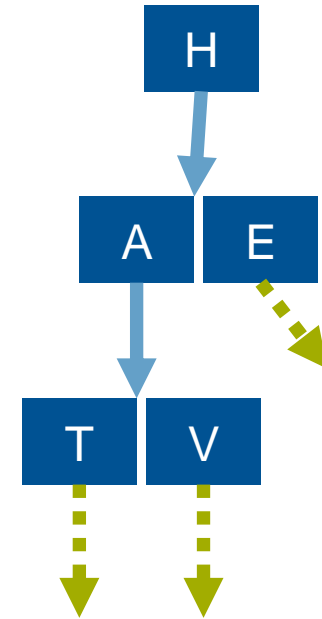
Max Height for k bit keys: $\text{ceil}(k/s)$

Radix Trees

Keys:
HELLO,
HAT,
HAVE

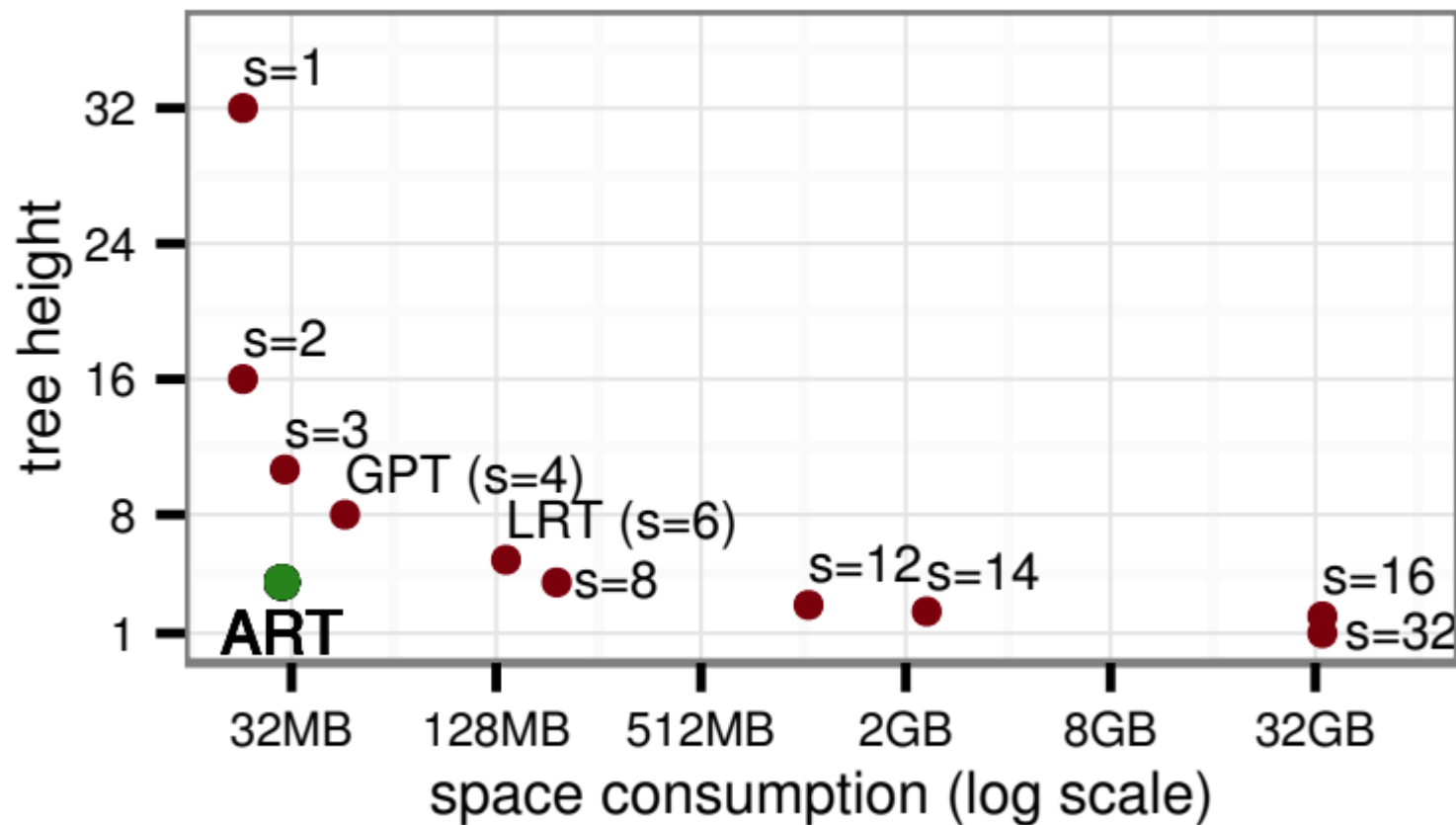


OR



Tuple
Pointer

Height vs. Space Tradeoff

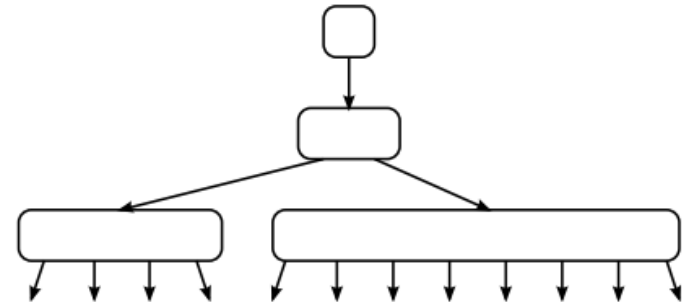
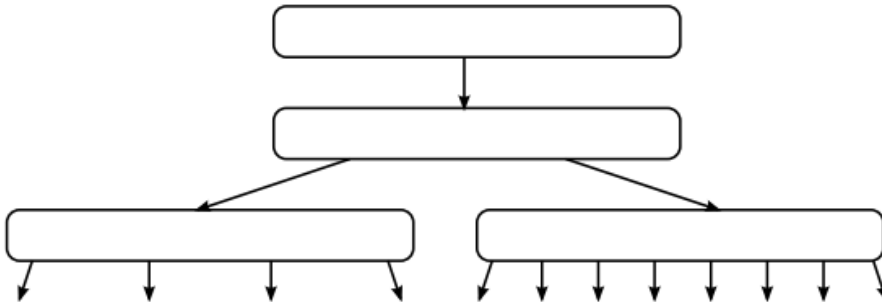


Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

Adaptive Nodes

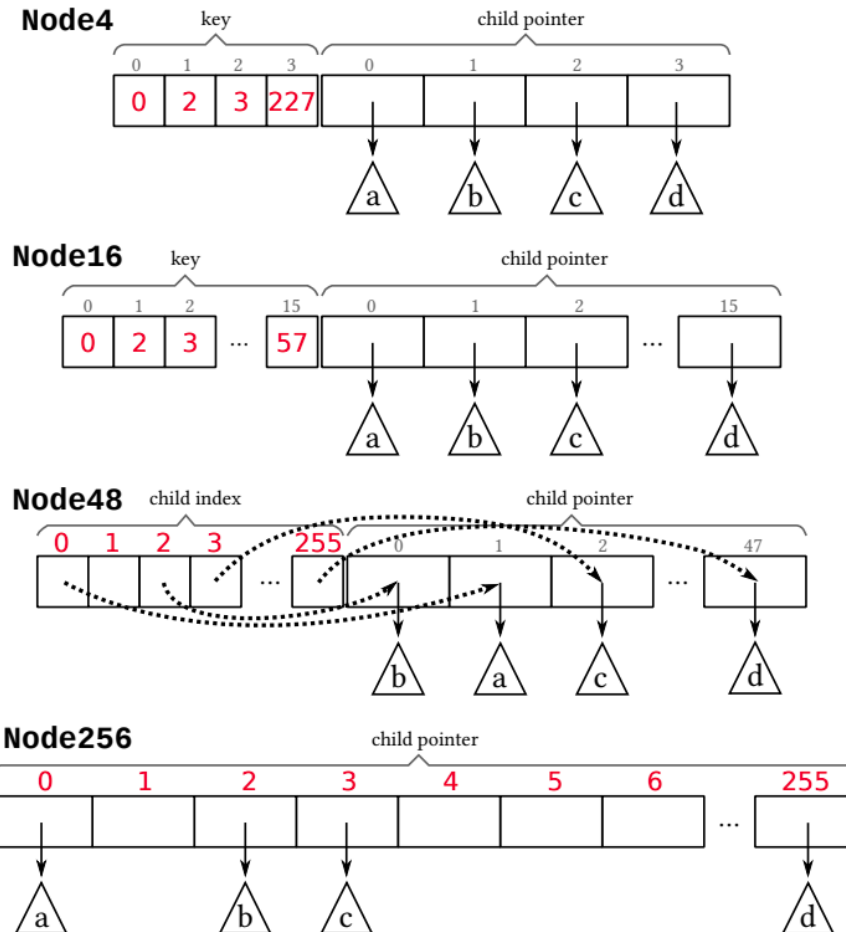
Main Idea:

Use different node types with different fanout based on number of non-null children.



Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

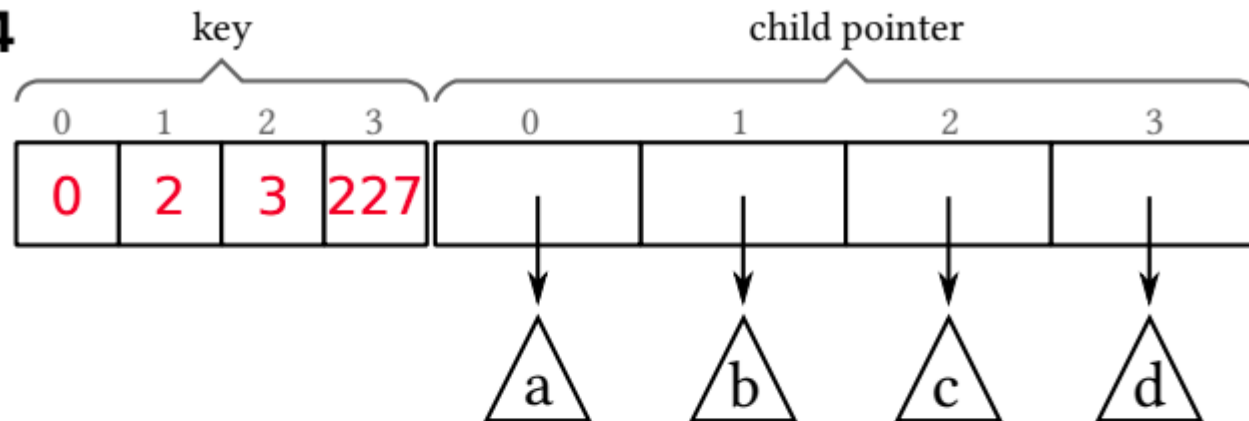
ART – Adaptive Nodes



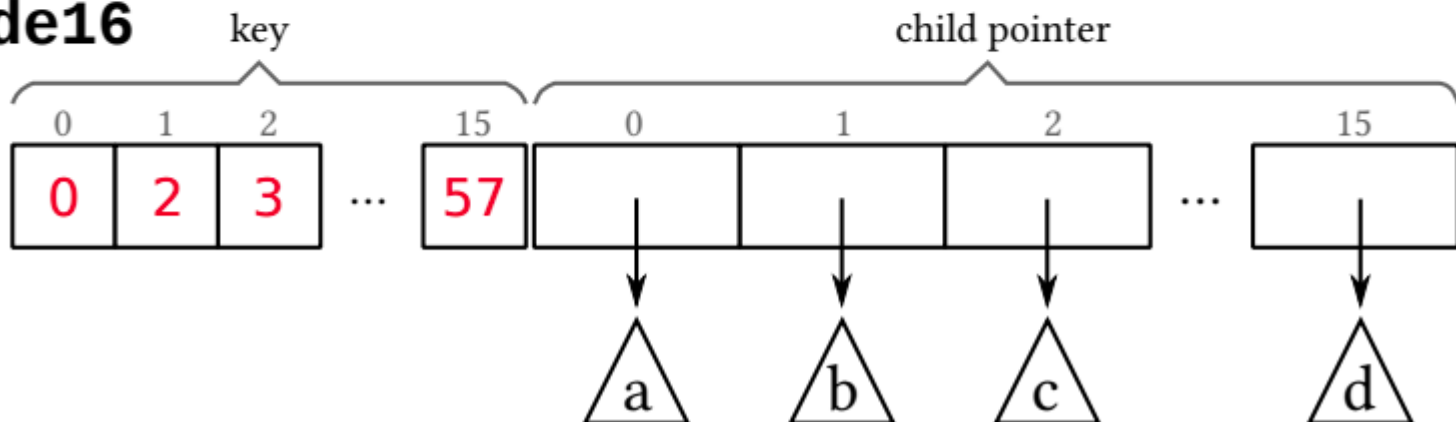
Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Adaptive Nodes

Node4

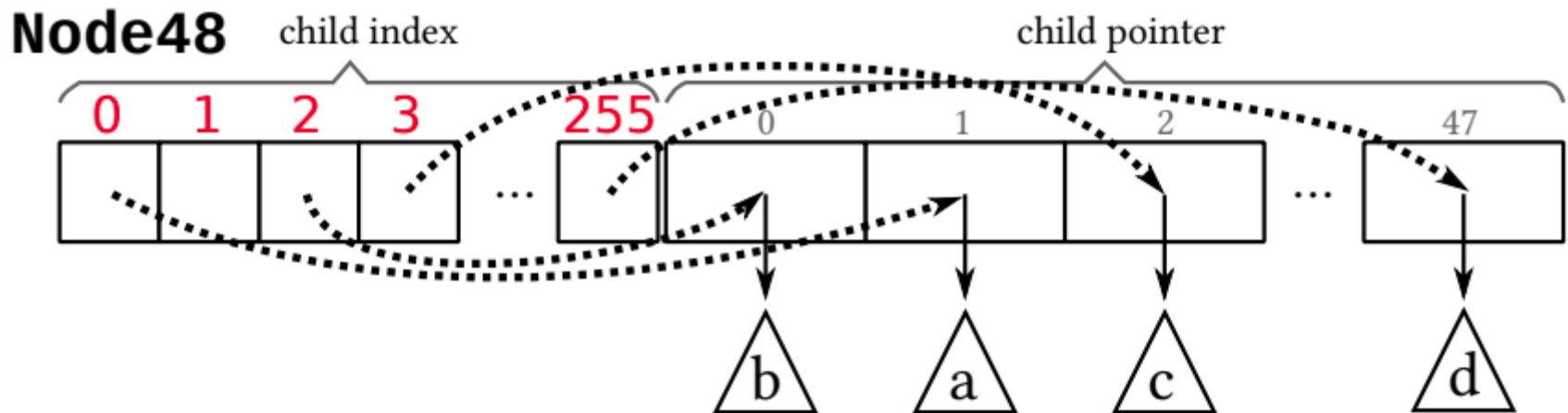


Node16



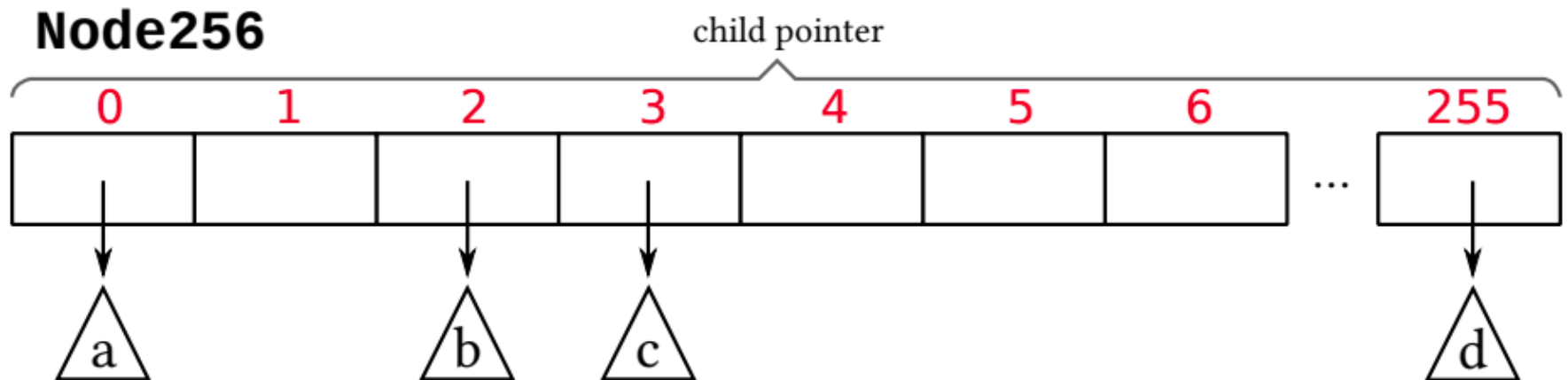
Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Adaptive Nodes



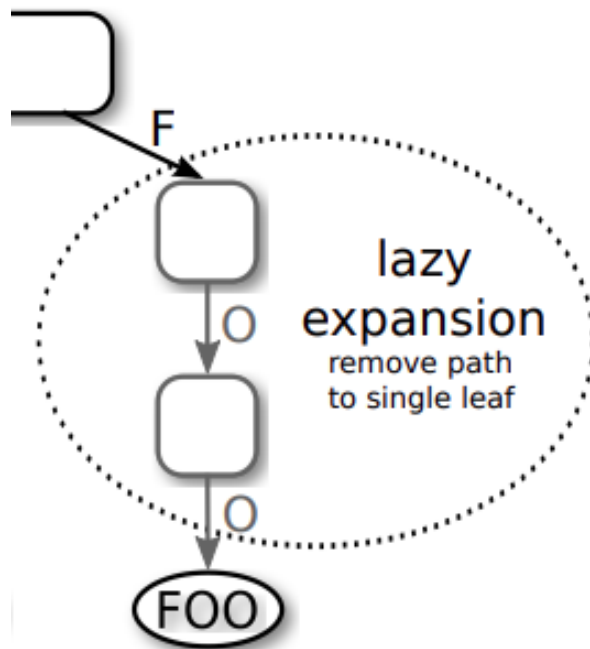
Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Adaptive Nodes



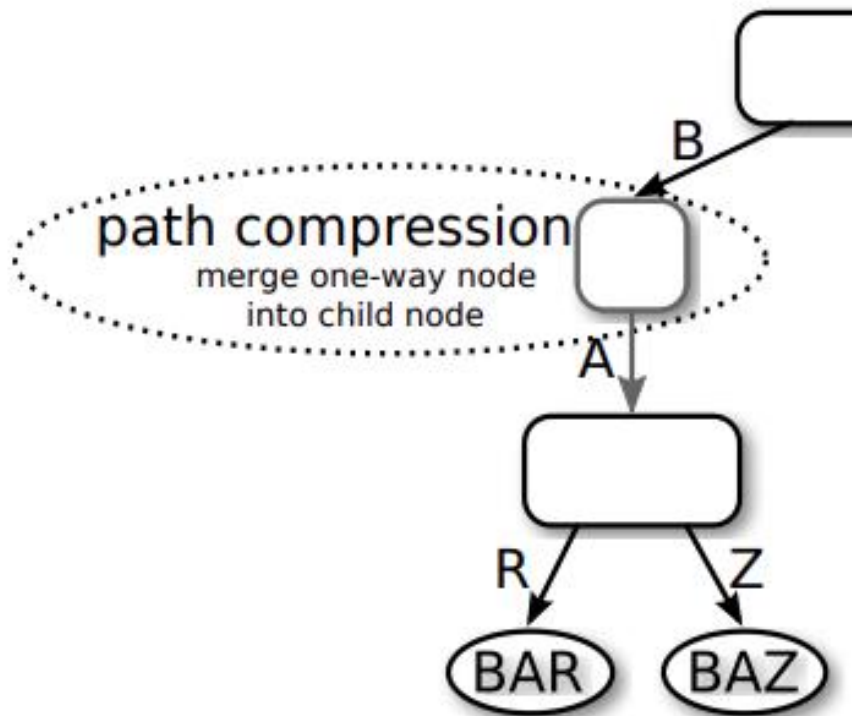
Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Lazy Expansion



Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Path Compression



How to handle omitted partial key?

- **Pessimistic:**
store omitted partial key at parent node
- **Optimistic:**
only store count of omitted nodes and skip over equivalent key part (compare full key when reaching leaf)

Hybrid approach in ART:

Use pessimistic for size of up to 8 bytes for partial key. Then use optimistic.

Key Transformations for Bitwise Comparisons

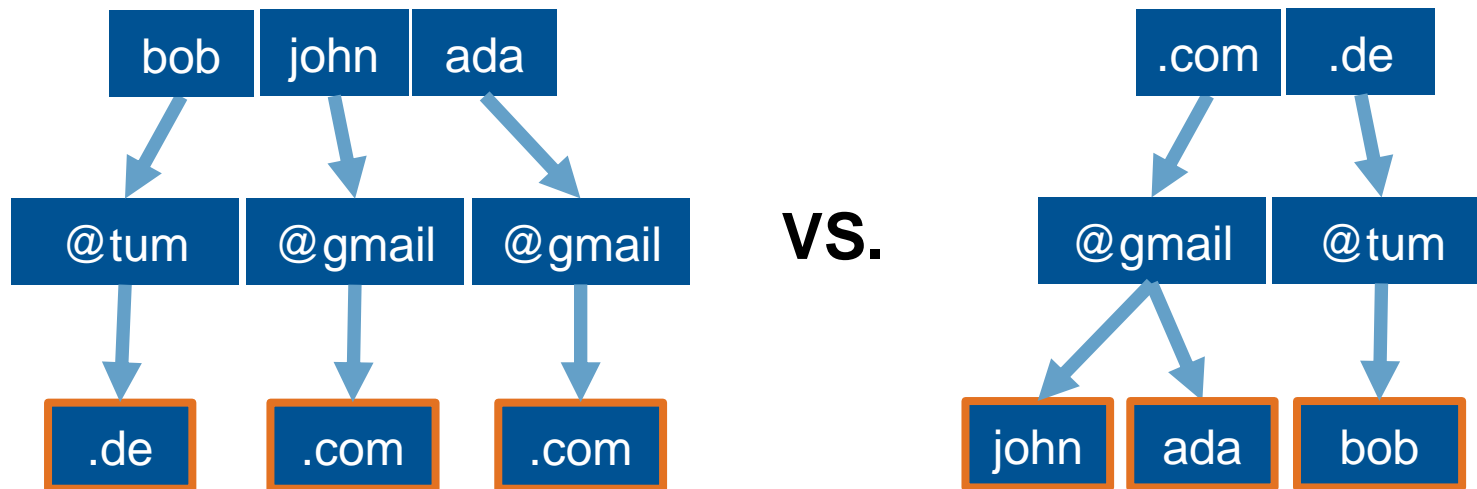
How to translate attribute types to make them binary comparable?

$$a < b \Leftrightarrow \text{bin}(a) < \text{bin}(b) \quad (\text{same for } > \text{ and } =)$$

- **Unsigned Integers:** Flip byte order for little endian machines
- **Signed Integers:** Flip sign bit (so negative values come before positive ones), store as unsigned integer
- **Floats:** Classify into group (neg vs. Pos, normalized vs. denormalized, NaN, ∞ , 0), store as unsigned integer
- **Strings:** Library functions for Unicode Strings
- **Null:** Identify special value
- **Compound Type:** Transform each attribute separately, concatenate results

Store Attributes in Context Sensitive Way

Example: E-Mails from back to front for best compression



ART – Implementation Specifics

ART implementation for storing 32 bit keys (no path compression)

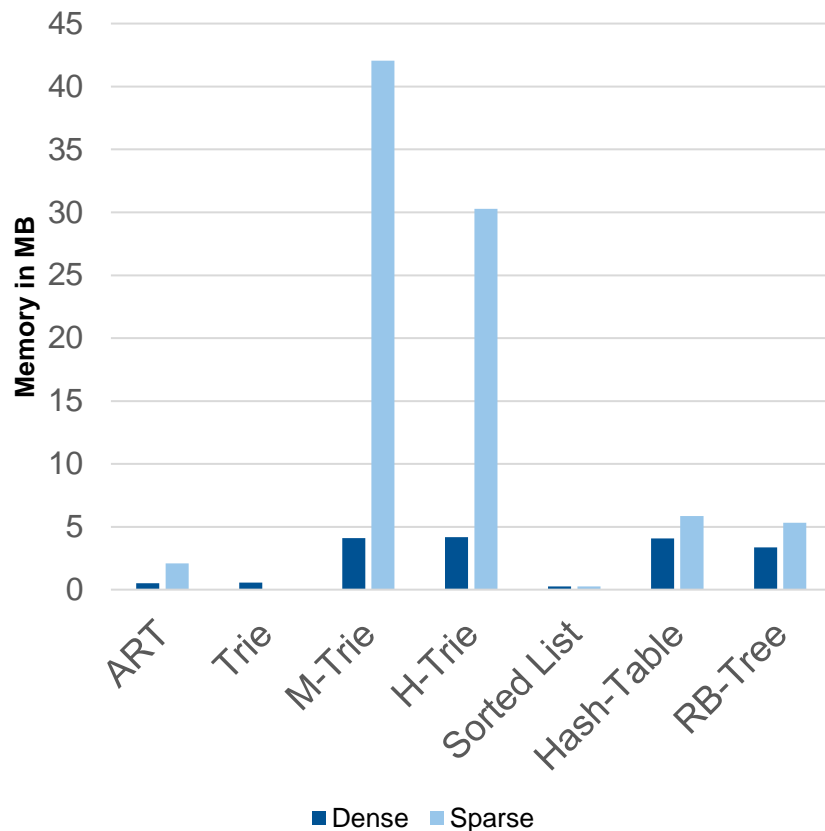
Important Points:

- No recursion! (self explanatory)
- Careful with polymorphie! (C++ utilizes vtables for dynamic dispatch → performance + memory overhead)
- Memory management! ([malloc vs. TCMalloc vs. jemalloc](#) vs. custom memory arenas/pools)
- Node16 search SIMD comparison: used x86-64 SSE2
- Combined value/pointer slots: Pointer Tagging using 3 low bits on 64 bit architecture (pointers are 8-byte aligned)

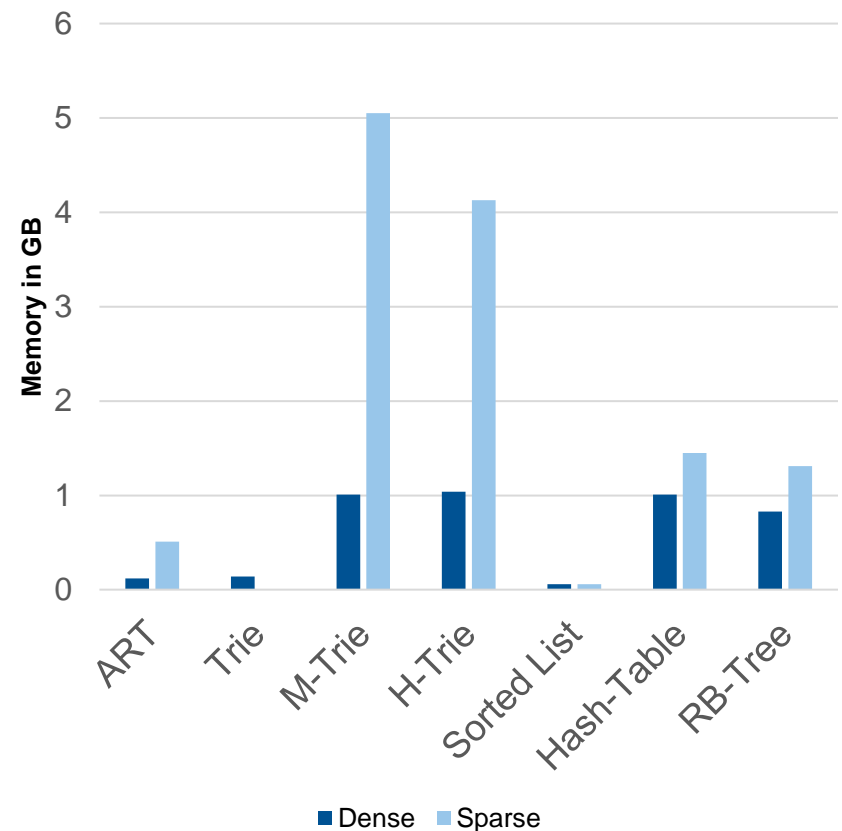
Source: <https://github.com/atalantus/The-Adaptive-Radix-Tree/tree/main/Implementation>

Memory Benchmark

65K random uint32_t keys
(lower is better)



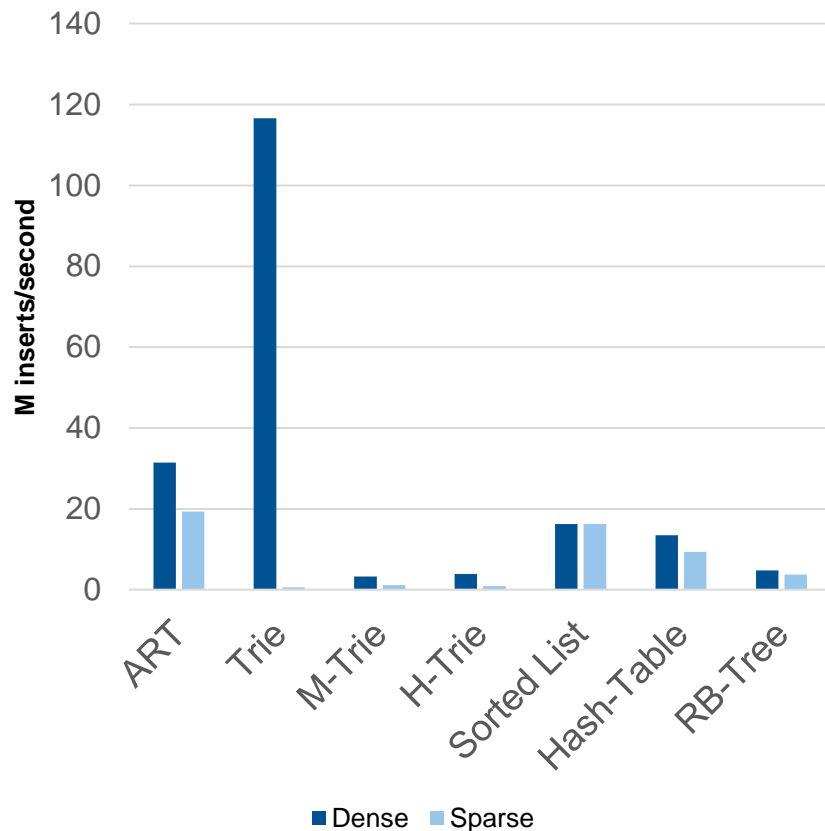
16M random uint32_t keys
(lower is better)



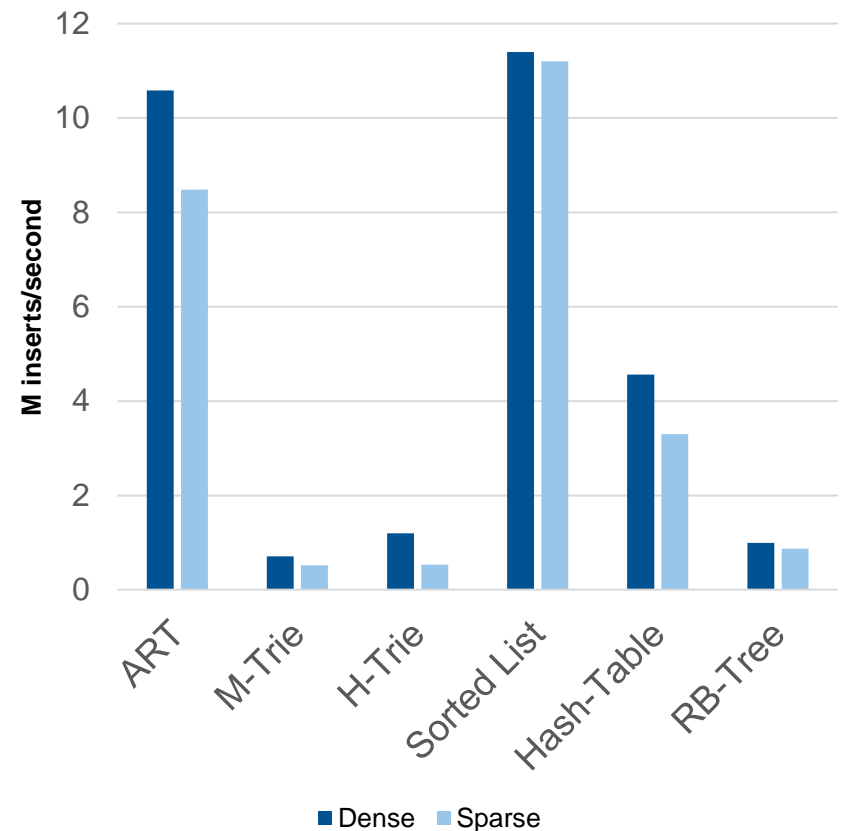
Source: <https://github.com/atalantus/The-Adaptive-Radix-Tree/tree/main/Implementation>

Performance Benchmark (insert)

65K random uint32_t keys
(higher is better)



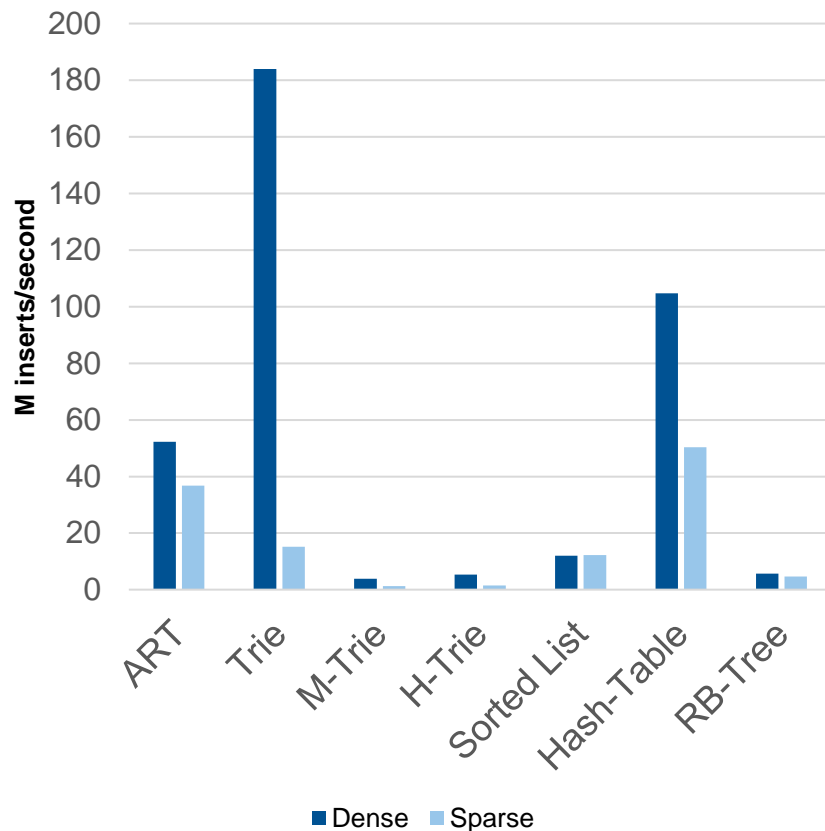
16M random uint32_t keys
(higher is better)



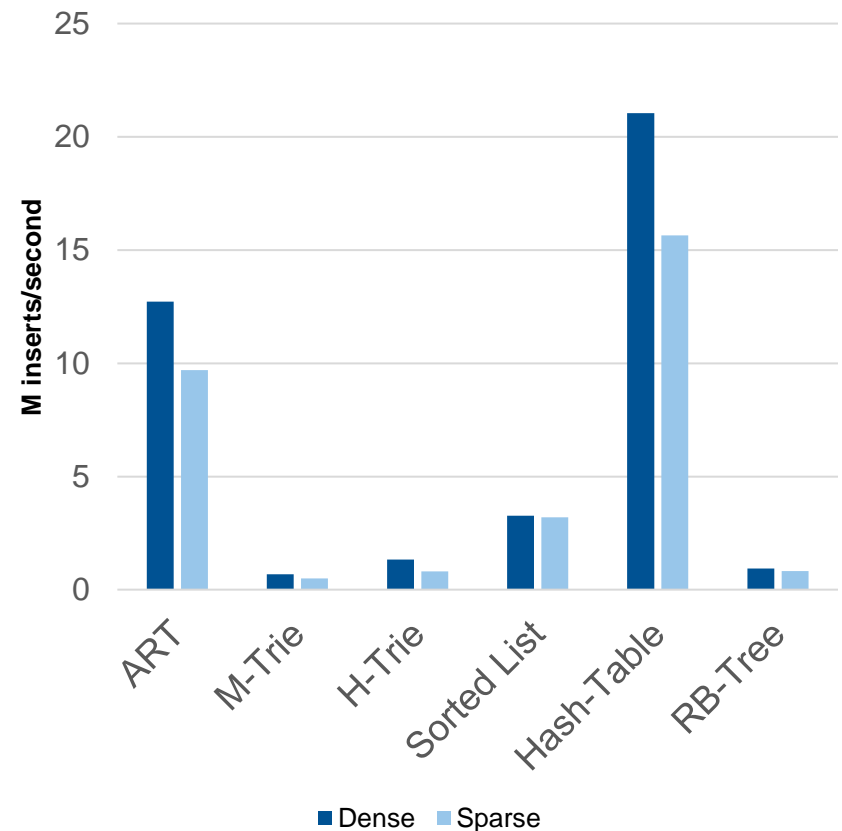
Source: <https://github.com/atalantus/The-Adaptive-Radix-Tree/tree/main/Implementation>

Performance Benchmark (search)

65K random uint32_t keys
(higher is better)

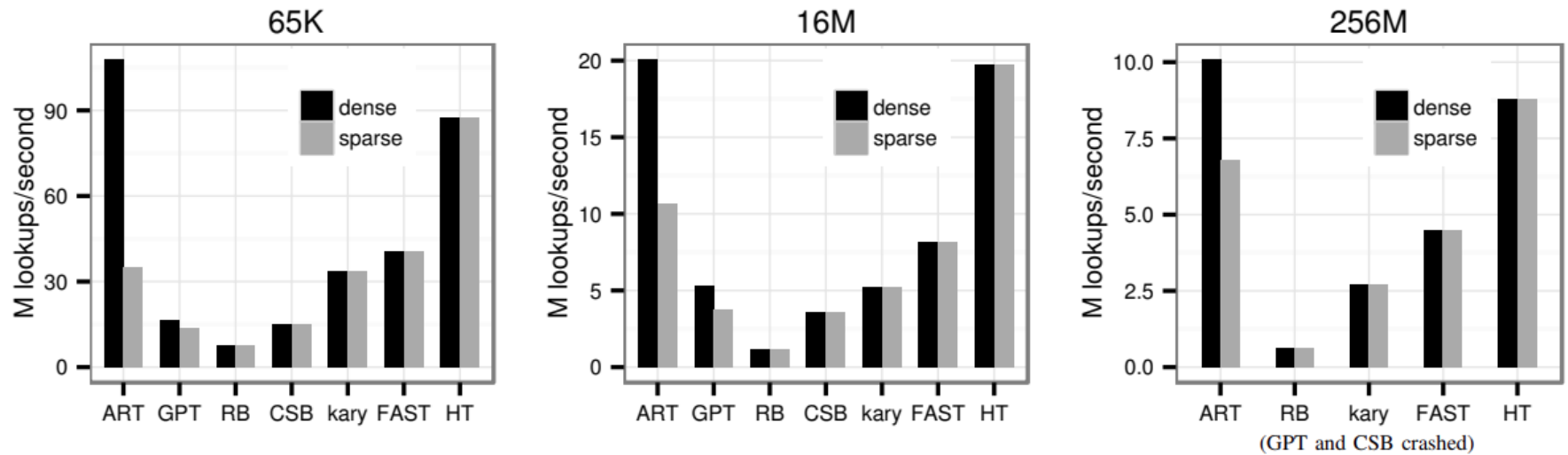


16M random uint32_t keys
(higher is better)



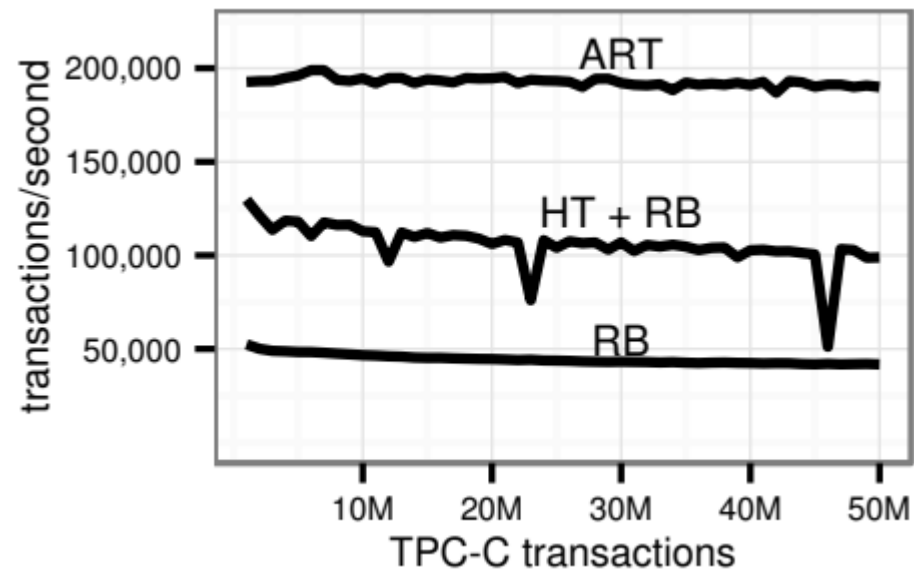
Source: <https://github.com/atalantus/The-Adaptive-Radix-Tree/tree/main/Implementation>

Performance Benchmark (search)



Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

TPC-C Benchmark



Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

Other Benchmarks: Memory Usage

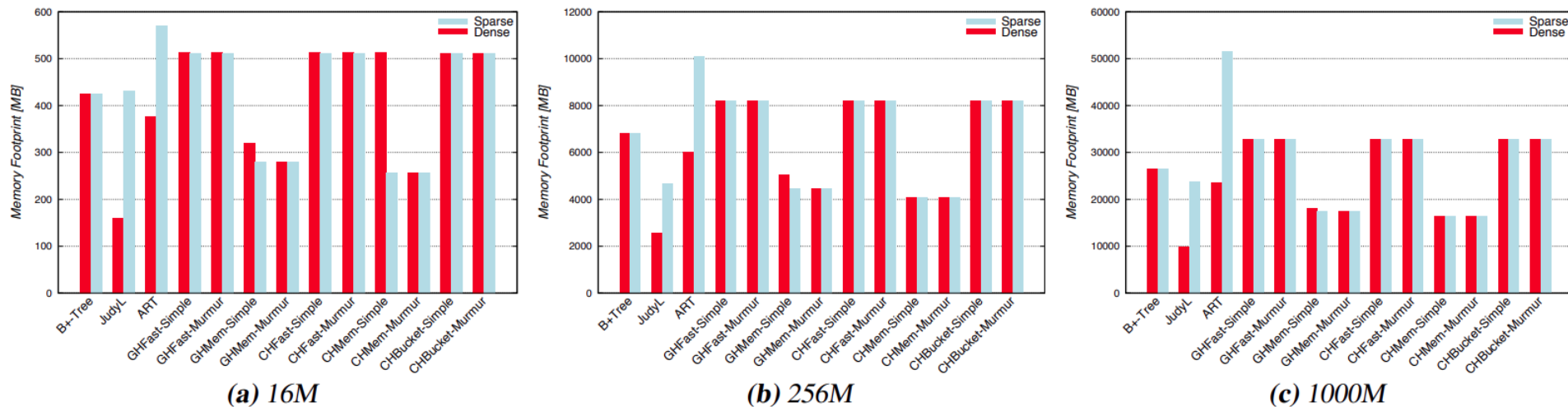


Figure 9: Memory footprint in MB (covering). Lower is better.

Source: [A comparison of adaptive radix trees and hash tables. IEEE 2015](#)

Other Benchmarks: Insert

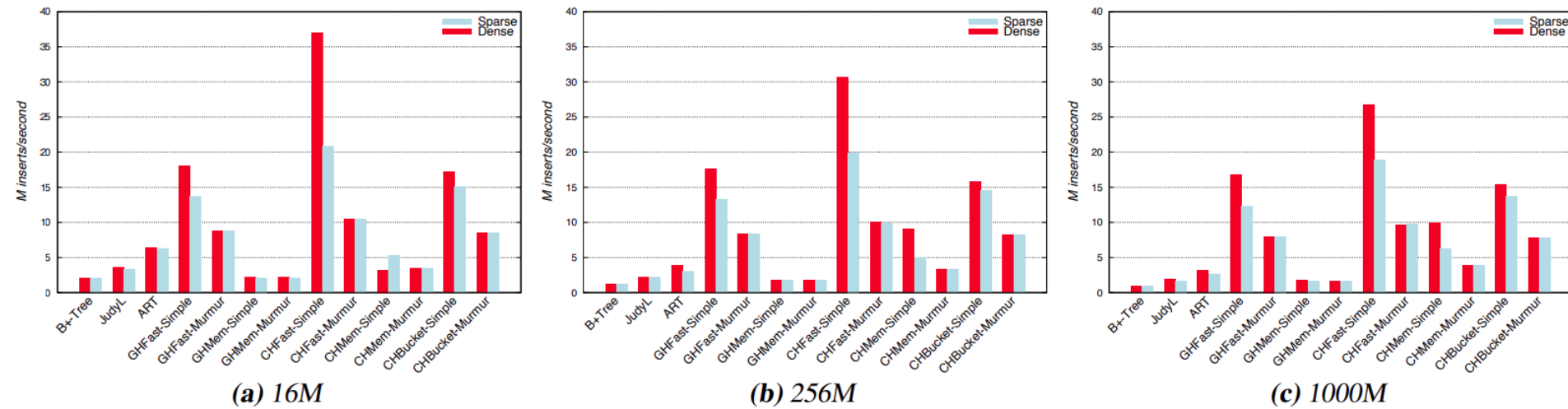


Figure 7: Insertion throughput (covering). Higher is better.

Source: [A comparison of adaptive radix trees and hash tables. IEEE 2015](#)

Other Benchmarks: Search

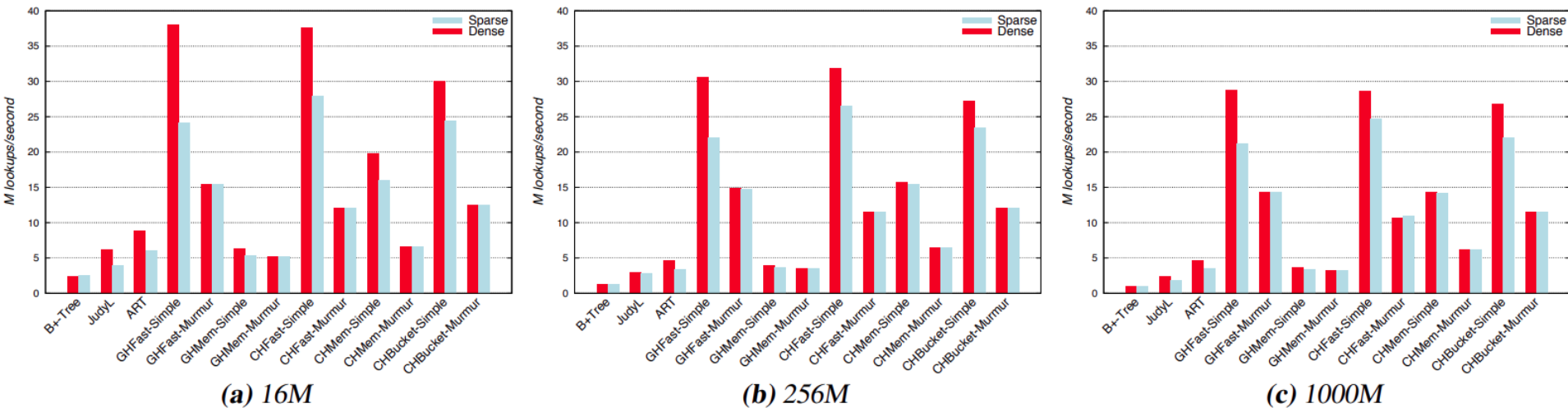
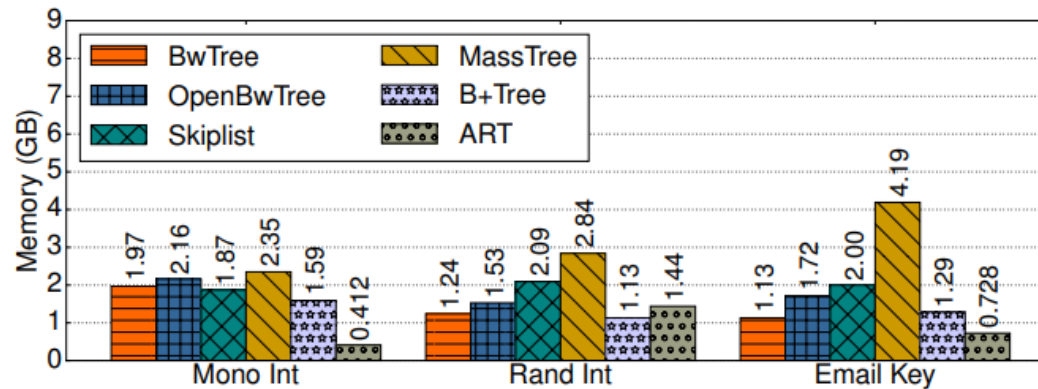


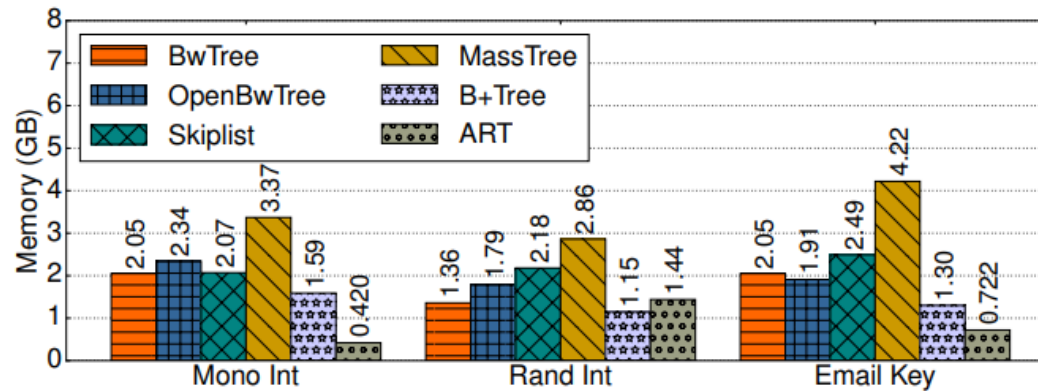
Figure 8: Lookup throughput (covering). Higher is better.

Source: [A comparison of adaptive radix trees and hash tables. IEEE 2015](#)

Other Benchmarks: Memory Usage



(a) Single-Threaded – Read/Update



(b) Multi-Threaded – Read/Update

Source: [Building a Bw-Tree Takes More Than Just Buzz Words. SIGMOD, 2018](#)

Other Benchmarks: Performance

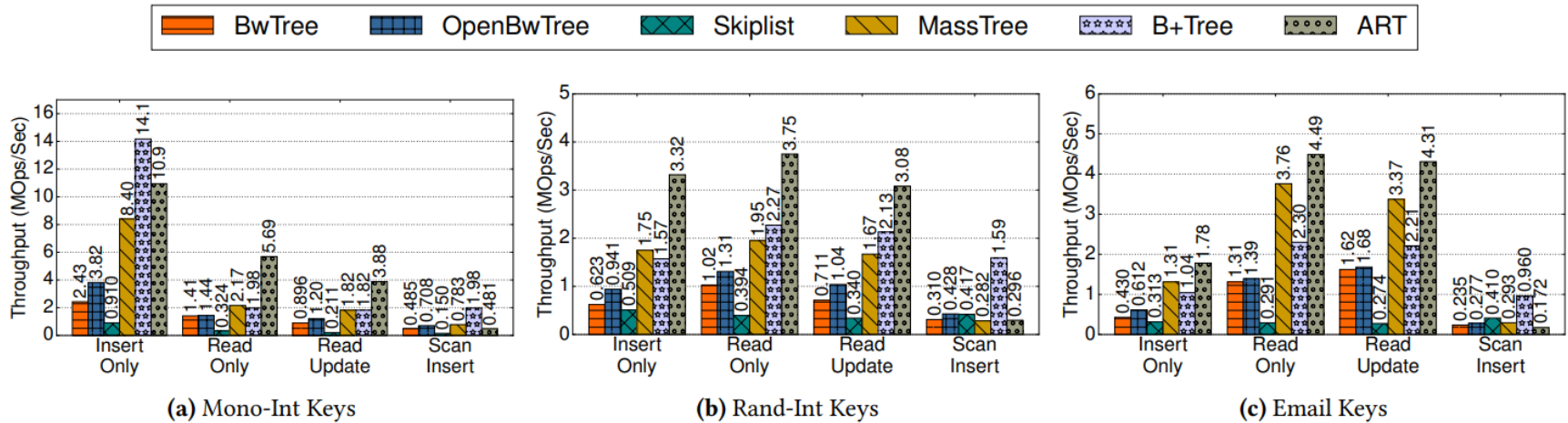


Figure 13: In-Memory Index Comparison (Single-Threaded) – The worker thread is pinned to NUMA node 0.

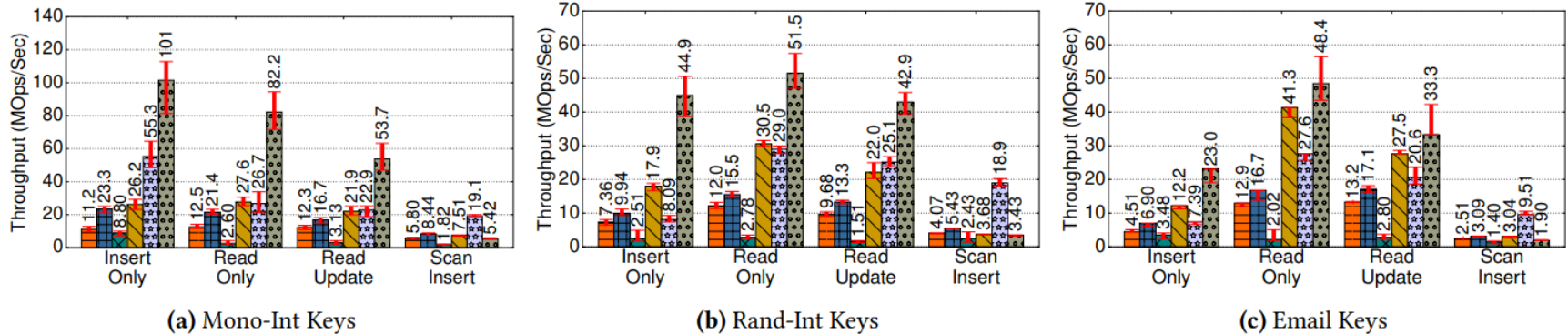


Figure 14: In-Memory Index Comparison (Multi-Threaded) – 20 worker threads. All worker threads are pinned to NUMA node 0.

Source: [Building a Bw-Tree Takes More Than Just Buzz Words. SIGMOD, 2018](#)

Summary & Conclusion

ART as (synchronized) order preserving index structure is

- space efficient
- outperforms other op-index structures in insert, search and update performance
- usable for typical data types
- used as default index in (at least) HyPer and DuckDB

	Range Support	Memory Efficiency	Performance	Researched
B+-Tree	+	–	-- (best for range)	+
ART	+	~	~	~
Judy-Array	+	+	–	–
Hash-Table	–	+	++	~

Summary & Conclusion

Further research being done:

[The ART of practical synchronization](#) (DaMoN, 2016)

- synchronize ART for multithreading using combination of OLP and ROWEX

[HOT: A Height Optimized Trie Index for Main-Memory Database Systems](#) (SIGMOD, 2018)

- ART performs worse on string (more random binary distribution) due to low fanout in lower tree levels
- Idea: Achieve consistently high fanout by adapting the span at each node

[START – Self-Tuning Adaptive Radix Tree](#) (IEEE, 2020)

- Learned indexes outperforming ART in read-mostly benchmark due to hierarchical structure of nodes
- Idea: Optimize lower level nodes using cost model and optimizer

(Judy Arrays patent expired this year)

- even better cache efficiency, complex node transitions (~20k LoC ☺)
- Judy-Pointers

Benchmark Specs

Intel Core i5-8400 CPU @ 2.80GHz

- L1 Instruction Cache: 6 x 32 KB
- L1 Data Cache: 6 x 32 KB
- L2 Cache: 6 x 256 KB
- L3 Cache: 9 MB

16 GB DDR4 RAM

Windows 10 Pro (10.0.19044 Build 19044)

ART – Search

```
    search (node, key, depth)
1  if node==NULL
2      return NULL
3  if isLeaf(node)
4      if leafMatches(node, key, depth)
5          return node
6      return NULL
7  if checkPrefix(node, key, depth) != node.prefixLen
8      return NULL
9  depth=depth+node.prefixLen
10 next=findChild(node, key[depth])
11 return search(next, key, depth+1)
```

Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Search

```
    findChild (node, byte)
1  if node.type==Node4 // simple loop
2      for (i=0; i<node.count; i=i+1)
3          if node.key[i]==byte
4              return node.child[i]
5      return NULL
6  if node.type==Node16 // SSE comparison
7      key=_mm_set1_epi8(byte)
8      cmp=_mm_cmpeq_epi8(key, node.key)
9      mask=(1<<node.count)-1
10     bitfield=_mm_movemask_epi8(cmp)&mask
11     if bitfield
12         return node.child[ctz(bitfield)]
13     else
14         return NULL
15 if node.type==Node48 // two array lookups
16     if node.childIndex[byte]!=EMPTY
17         return node.child[node.childIndex[byte]]
18     else
19         return NULL
20 if node.type==Node256 // one array lookup
21     return node.child[byte]
```

Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Insert

```
    insert (node, key, leaf, depth)
1  if node==NULL // handle empty tree
2      replace(node, leaf)
3      return
4  if isLeaf(node) // expand node
5      newNode=makeNode4()
6      key2=loadKey(node)
7      for (i=depth; key[i]==key2[i]; i=i+1)
8          newNode.prefix[i-depth]=key[i]
9      newNode.prefixLen=i-depth
10     depth=depth+newNode.prefixLen
11     addChild(newNode, key[depth], leaf)
12     addChild(newNode, key2[depth], node)
13     replace(node, newNode)
14     return
15 p=checkPrefix(node, key, depth)
```

Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Insert

```
16  if p!=node.prefixLen // prefix mismatch
17      newNode=makeNode4()
18      addChild(newNode, key[depth+p], leaf)
19      addChild(newNode, node.prefix[p], node)
20      newNode.prefixLen=p
21      memcpy(newNode.prefix, node.prefix, p)
22      node.prefixLen=node.prefixLen-(p+1)
23      memmove(node.prefix,node.prefix+p+1,node.prefixLen)
24      replace(node, newNode)
25      return
26  depth=depth+node.prefixLen
27  next=findChild(node, key[depth])
28  if next // recurse
29      insert(next, key, leaf, depth+1)
30  else // add to inner node
31      if isFull(node)
32          grow(node)
33      addChild(node, key[depth], leaf)
```

Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

ART – Space Consumption

TABLE I
SUMMARY OF THE NODE TYPES (16 BYTE HEADER, 64 BIT POINTERS).

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	$16 + 256 \cdot 8 = 2064$

TABLE II
WORST-CASE SPACE CONSUMPTION PER KEY (IN BYTES) FOR DIFFERENT
RADIX TREE VARIANTS WITH 64 BIT POINTERS.

	$k = 32$	$k \rightarrow \infty$
ART	43	52
GPT	256	∞
LRT	2048	∞
KISS	>4096	NA.

Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)

Performance-Benchmark (search)

TABLE III
PERFORMANCE COUNTERS PER LOOKUP.

	65K			16M		
	ART (d./s.)	FAST	HT	ART (d./s.)	FAST	HT
Cycles	40/105	94	44	188/352	461	191
Instructions	85/127	75	26	88/99	110	26
Misp. Branches	0.0/0.85	0.0	0.26	0.0/0.84	0.0	0.25
L3 Hits	0.65/1.9	4.7	2.2	2.6/3.0	2.5	2.1
L3 Misses	0.0/0.0	0.0	0.0	1.2/2.6	2.4	2.4

Source: [The adaptive radix tree: ARTful indexing for main-memory databases. IEEE, 2013](#)