# The Adaptive Radix Tree

ARTful Indexing for Main-Memory Databases
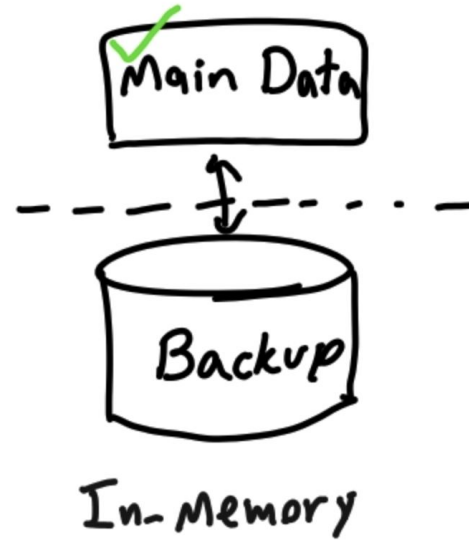
# T-tree

$O(\log_2 N)$



T Node

Parent Ptr

| data₁ | data₂ | data₃ | ••• | dataₙ |

control

Left Child Ptr    Right Child Ptr

T Tree

Do not optimally utilize
on-CPU cache

# Cache Sensitive B+ tree

$O(\log_b N)$

b: branching factor



| 22 | | | |

| 3 | 7 | 13 | 19 |     | 25 | 30 | 33 | |

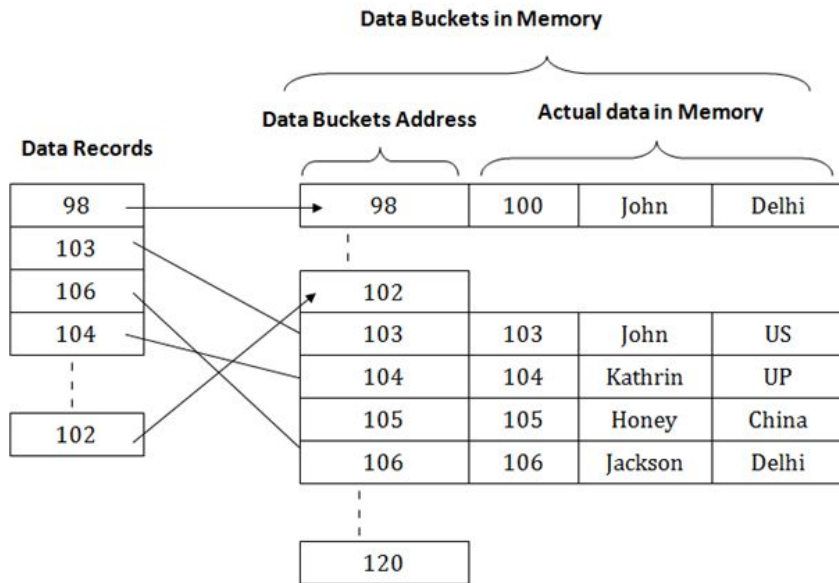| 2 | 3 | | 5 | 7 | | 12 | 13 |     | 16 | 19 | | 20 | 22 | | 24 | 25 | | 27 | 30 | | 31 | 33 | | 36 | 39 |

Require more expensive
update operations

# Is there any data structure that has better performance than O(log n)?



Only support point queries

Cannot handle growth well

Require expansive reorganization upon overflow with O(n) complexity

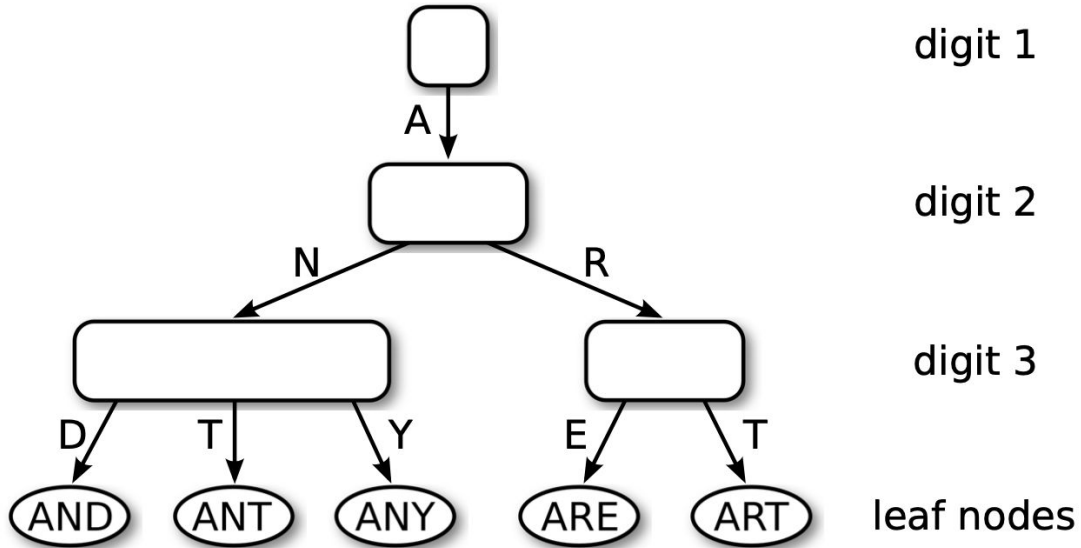**Hash Table !!!**     **O(1)**

# Unfortunate trade-off

Fast hash tables that only allow point queries

Fully-featured, but relatively slow, search trees.

**Any other possible solution with overall better performance?**

# Trie / Prefix Tree / Digital Search Tree
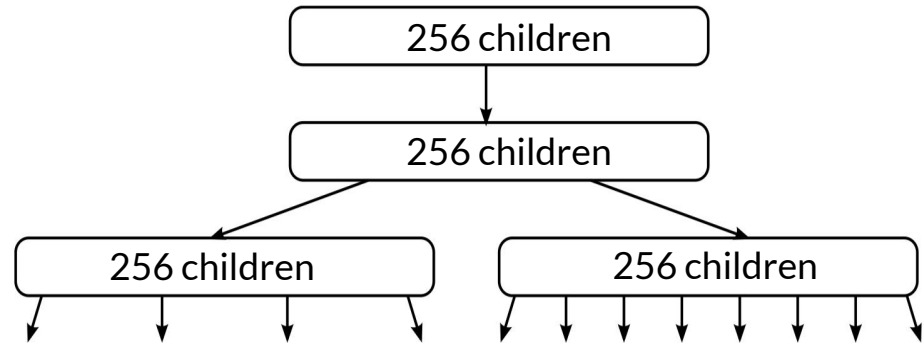


digit 1
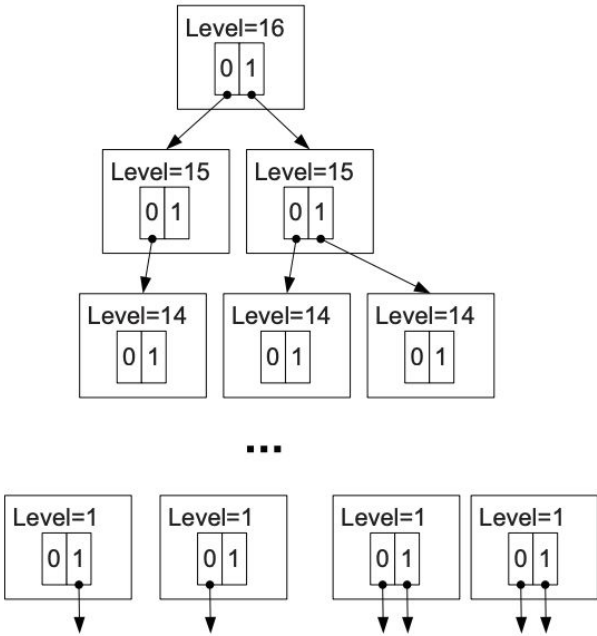
digit 2

O(k)

digit 3

leaf nodes

# Radix Tree

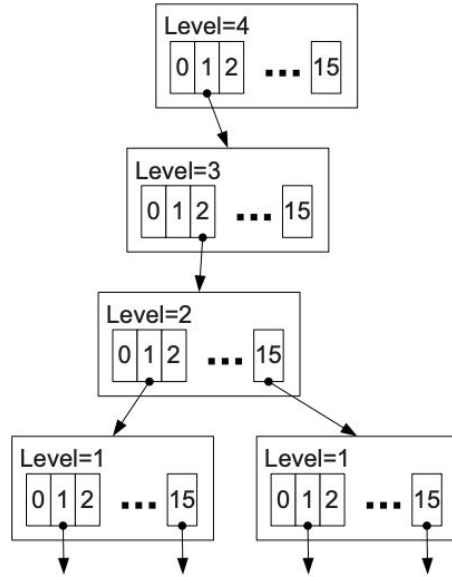00101101: 8-bits span

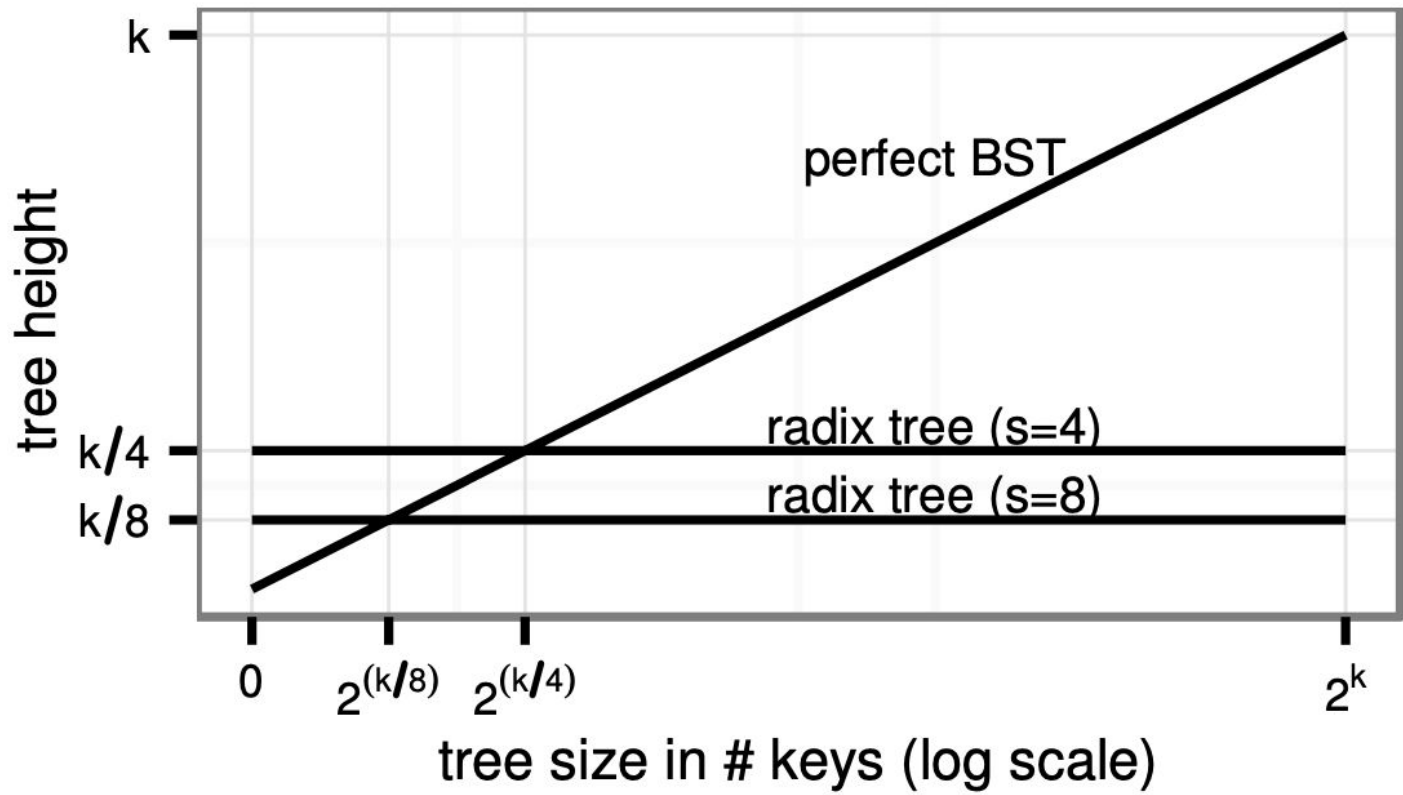00000000
00000001
00000010
...
11111110
11111111
} 2^8 children
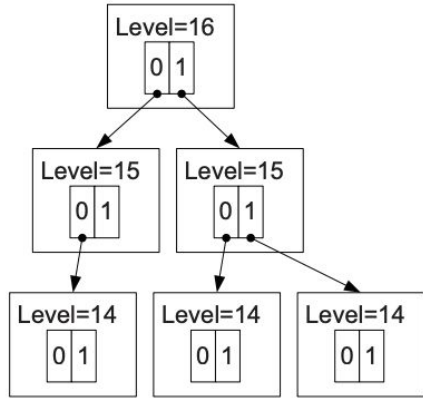
# Generalized Prefix Tree



(a) $k' = 1 \ (h = 16)$

(b) $k' = 4 \ (h = 4)$

k

tree height

k/4

k/8

perfect BST

radix tree (s=4)

radix tree (s=8)

0   $2^{(k/8)}$   $2^{(k/4)}$                                       $2^k$

tree size in # keys (log scale)
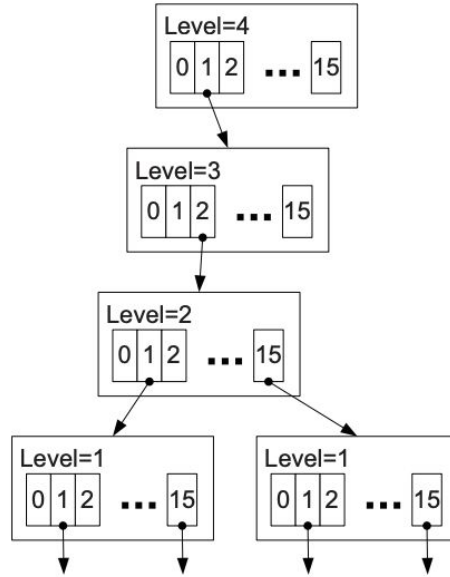
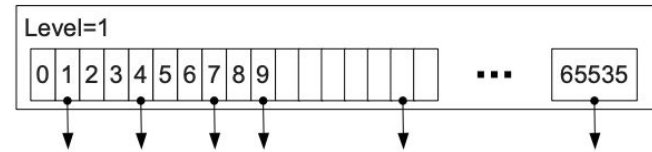BST: Binary Search Tree        s: span parameter        k: # bits of keys

# Generalized Prefix Tree



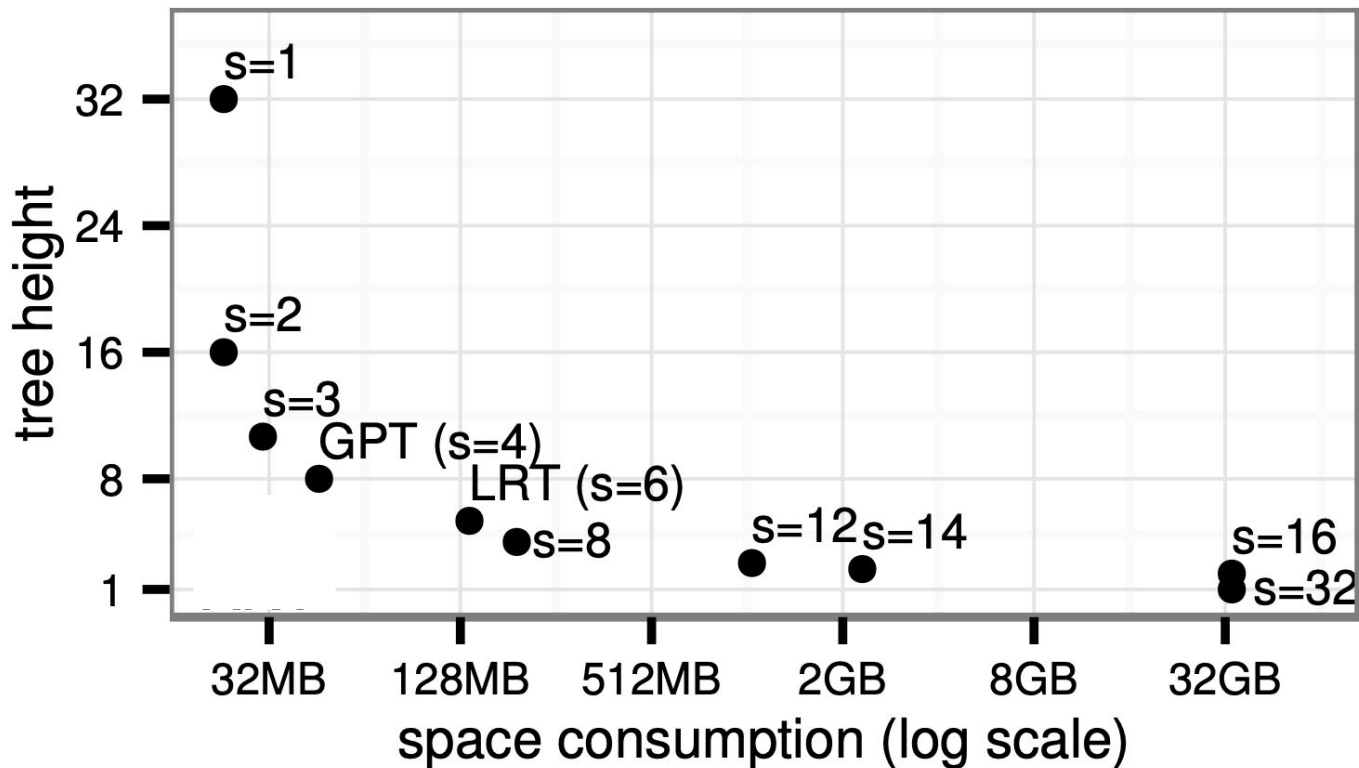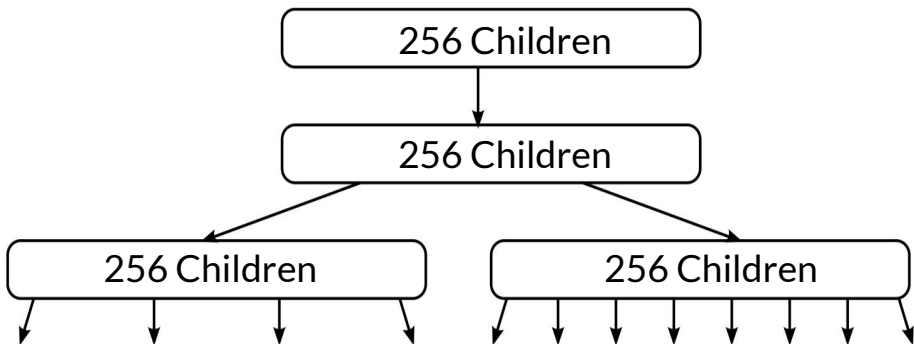(a) $k' = 1$ ($h = 16$)

(b) $k' = 4$ ($h = 4$)

(c) $k' = 16$ ($h = 1$)

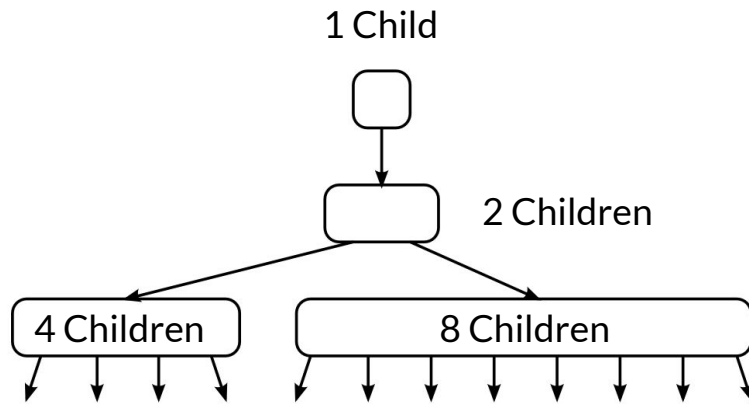# Tree Height vs. Space Consumption



GPT: Generalized Prefix Tree    LRT: Linux kernel radix tree    s: span parameter

# Adaptive Radix Tree (ART)



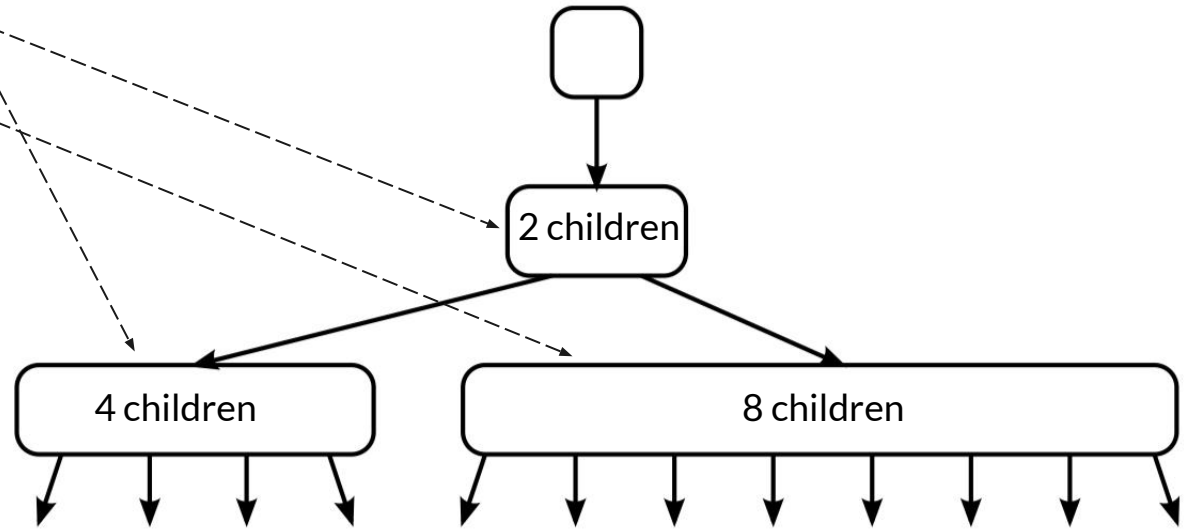Adaptive Node Design

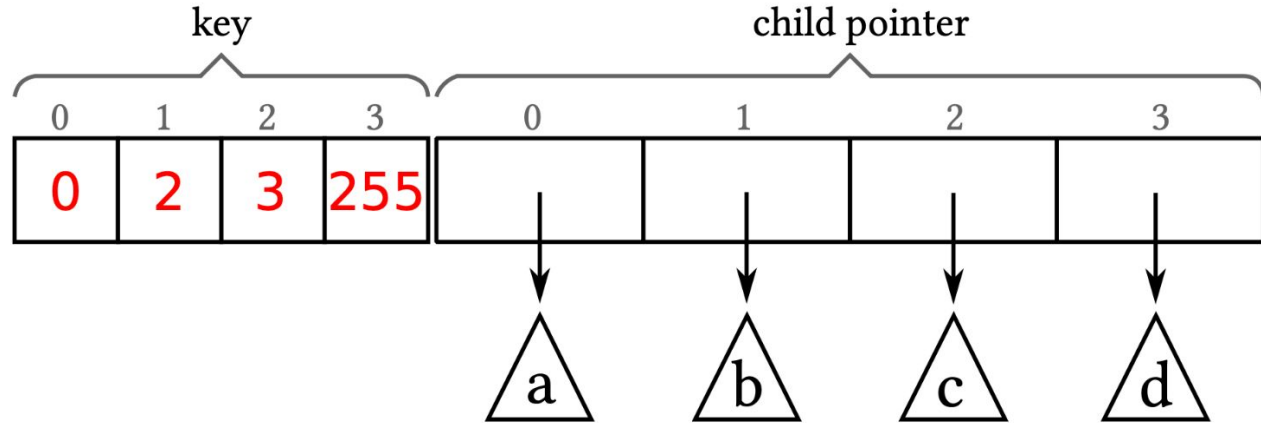# Adaptive Radix Tree (ART)

4 types of nodes:

Node4

Node16

Node48

Node256

# Space Consumption

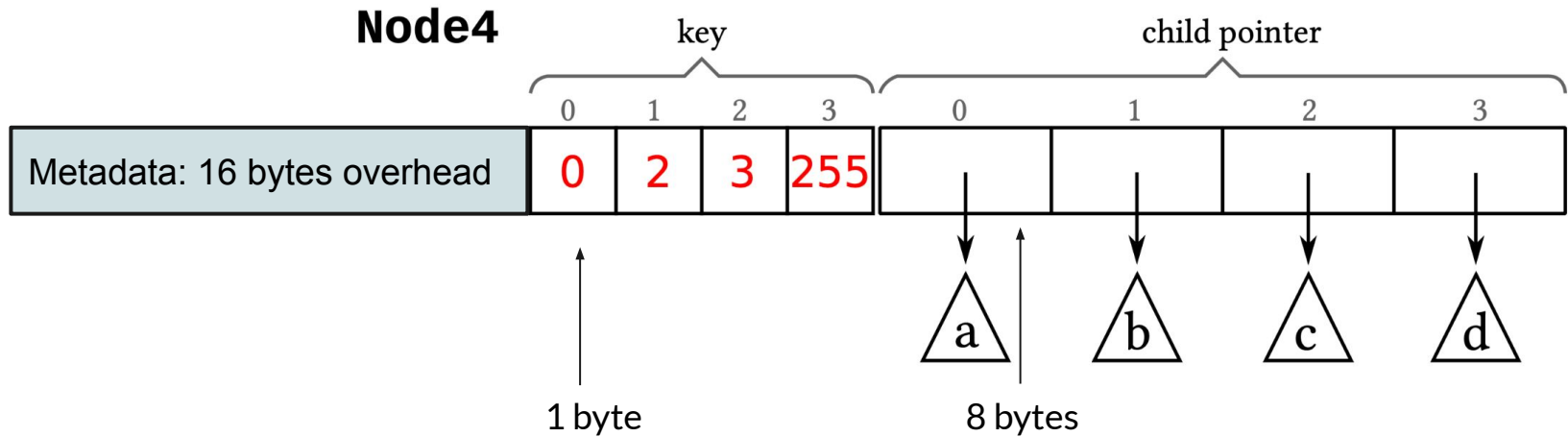**Node4**

key

child pointer

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 0 | 2 | 3 | 255 |

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| a | b | c | d |

Metadata: 16 bytes overhead

1 byte

8 bytes

How to search the key array?

**16 + 4 + 4 x 8 = 52 bytes**

**Node4**

key

child pointer

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 0 | 2 | 3 | 255 |

| 0 | 1 | 2 | 3 |
|---|---|---|---|

a   b   c   d

When search: Loop through

How to search the key array?

# ART Inner Node Type 2

Search in 16 keys

| | | |
|---|---|---|
| 01001100 | 00110101 | 0 |
| 01001100 | 00111110 | 0 |
| 01001100 | compare  01001100 | 1 |
| ... | ... | ... |
| 01001100 | 01110001 | 0 |
| 01001100 | 11001101 | 0 |

Replicate our key 16 times                Compare in parallel

# ART Inner Node Type 3

# ART Inner Node Type 4

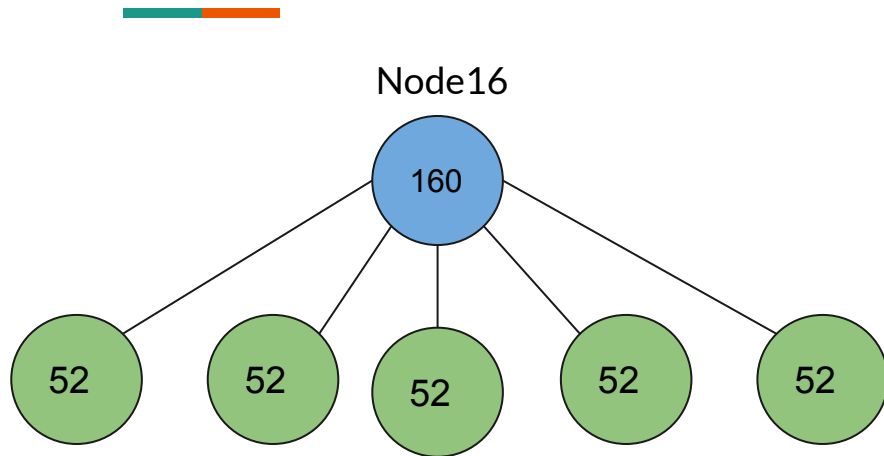| Type | Children | Space (bytes) |
|------|----------|---------------|
| Node4 | 2-4 | $16 + 4 + 4 \cdot 8 = 52$ |
| Node16 | 5-16 | $16 + 16 + 16 \cdot 8 = 160$ |
| Node48 | 17-48 | $16 + 256 + 48 \cdot 8 = 656$ |
| Node256 | 49-256 | $16 + 256 \cdot 8 = 2064$ |

# Space Consumption per key



Node2

52

52          52          52 bytes budget/leaf

Space consumption per key is bounded by 52

| Type | Minimum Children | bytes |
|------|------------------|-------|
| Node4 | 2 | 52 |
| Node16 | 5 | 160 |
| Node48 | 17 | 656 |
| Node256 | 49 | 2064 |

# Space Consumption per key

Node16



| Type | Minimum Children | bytes |
|------|------------------|-------|
| Node4 | 2 | 52 |
| Node16 | 5 | 160 |
| Node48 | 17 | 656 |
| Node256 | 49 | 2064 |

52 x 5 = 260 > 160

52 x 17 = 884 > 656
52 x 49 = 2548 > 2064

# Space Consumption per key



2064

...

2064

52          52

Non-adaptive radix tree:
Unbounded space consumption per key
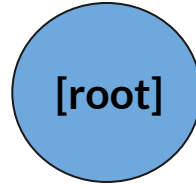
# ART Demo

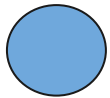Search

Insert

Delete

Lazy Expansion

Path Compression

Key: smile

Value: 10

Inner Node

Leaf Node

[root]

How is ART different from a trie in insertion?

Hint: Which unique technique does ART have for saving space?

# ART Insertion and Path Compression

Key: smile

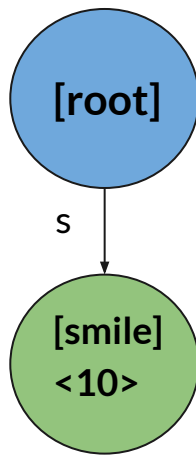Value: 10


Inner Node


Leaf Node

[root]

s

[smile]
<10>

Instead of saving s->m->i->l->e sequence in individual nodes, ART path compression technique saves it in one node and reduces space consumption.
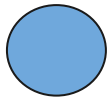
**Path compression:**
One inner node is removed when it has only one or no child node.
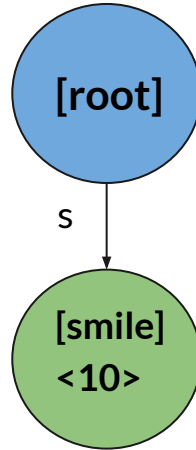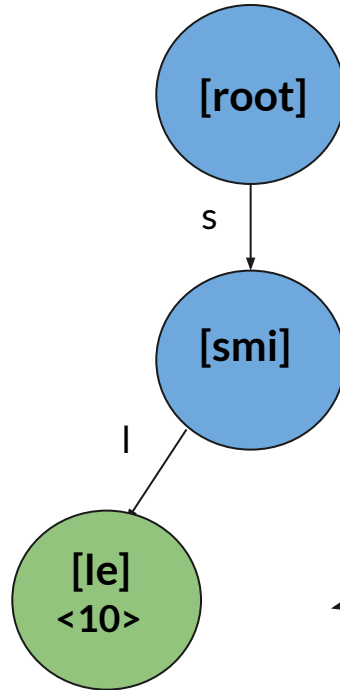
# ART Insertion

Key: smith

Value: 20

[root]

s

[smile]
<10>

Inner Node

Leaf Node

**What will happen next?**

# ART Lazy Expansion



Key: smith

Value: 20

Inner Node

Leaf Node

[root]

s

[smi]

l

[le]
<10>

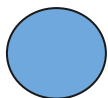How about finding the common part of [smile] and [smith]?

What happened to the previous leaf node [smile]?
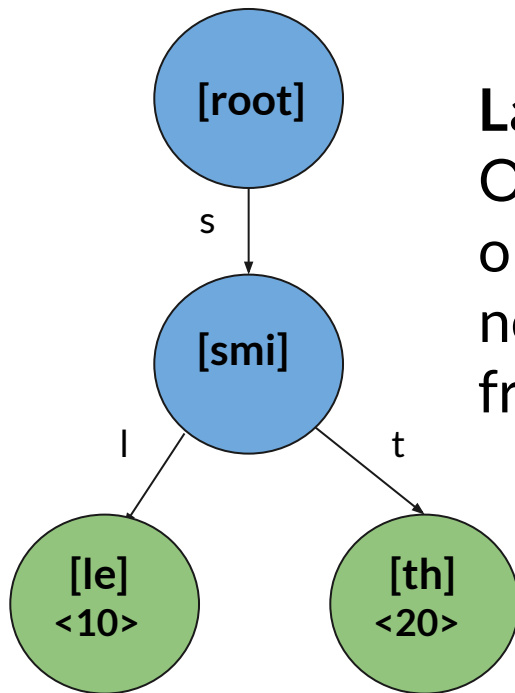
# ART Lazy Expansion

Key: smith

Value: 20



Inner Node

Leaf Node

**Lazy Expansion:**
One new inner node is created only when at least two leaf nodes need it to distinguish from.
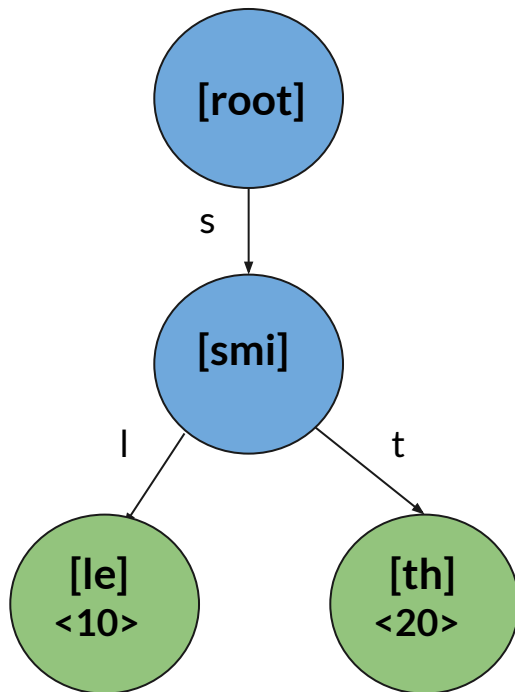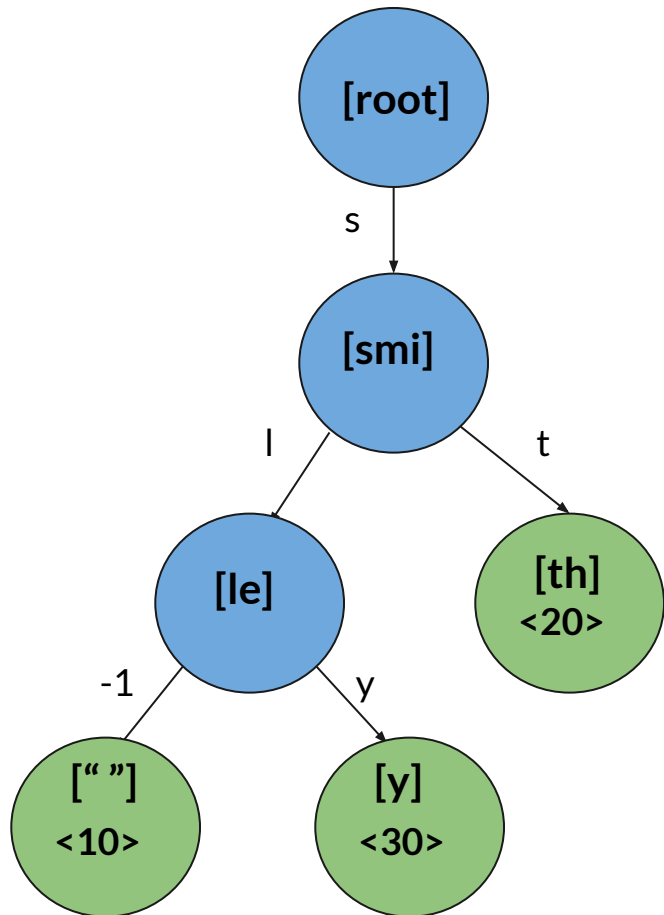
# ART Insertion



Key: smiley

Value: 30

What's next if I want to insert a new key smiley?

Inner Node

Leaf Node

[root]

s

[smi]

l

t

[le]
<10>

[th]
<20>

?

# ART Deletion



Key: smile

Inner Node

Leaf Node

[root]

s

[smi]

l          t

[le]        [th] <20>

-1          y

[" "] <10>     [y] <30>

**Edge case: An empty key.**

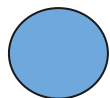**If I delete it, what should happen?**

**What's the first step before the deletion?**

# ART Search



Key: smile

Inner Node

Leaf Node

[root]

s

[smi]

l          t

[le]          [th]
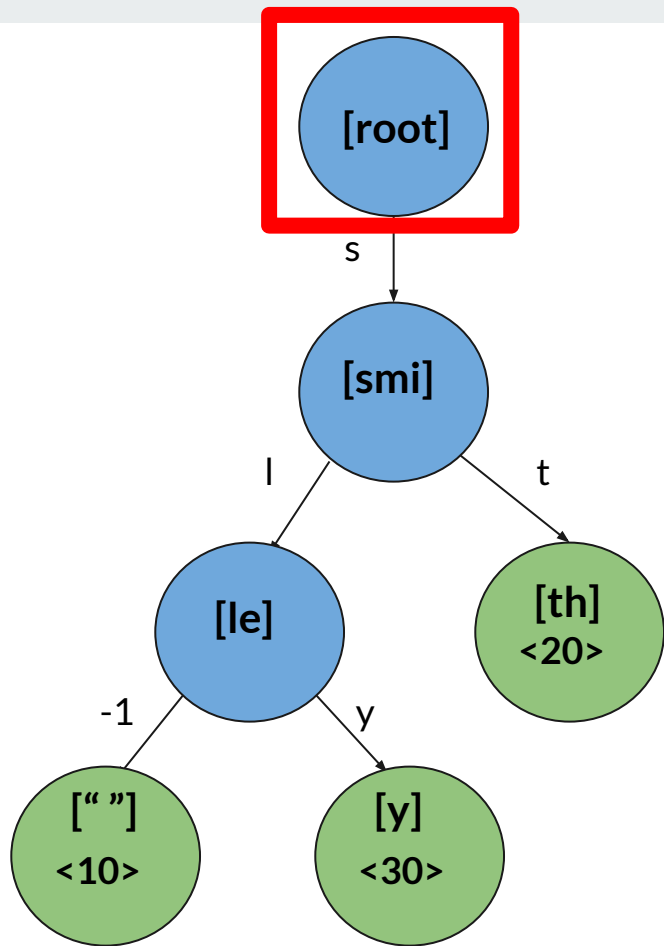              <20>

-1          y

[" "]          [y]
<10>          <30>

Is this node null?

Is this node a leaf node?

Is this node a inner node?

**Searching from the root...**

# ART Search

ART Search

# ART Path Compression

# ART Bulk Loading

Key 1: smile

Key 2: smith

Key 3: smiley



**Bulk load all nodes.**

Starts from the first byte of each key and go through all keys to find the most appropriate inner node type.

Build ART while perform radix sort -> Higher time complexity.

# Evaluation - Random Search
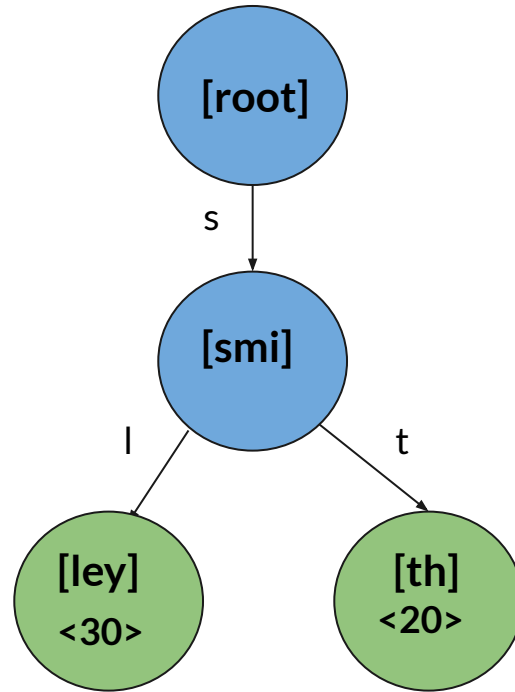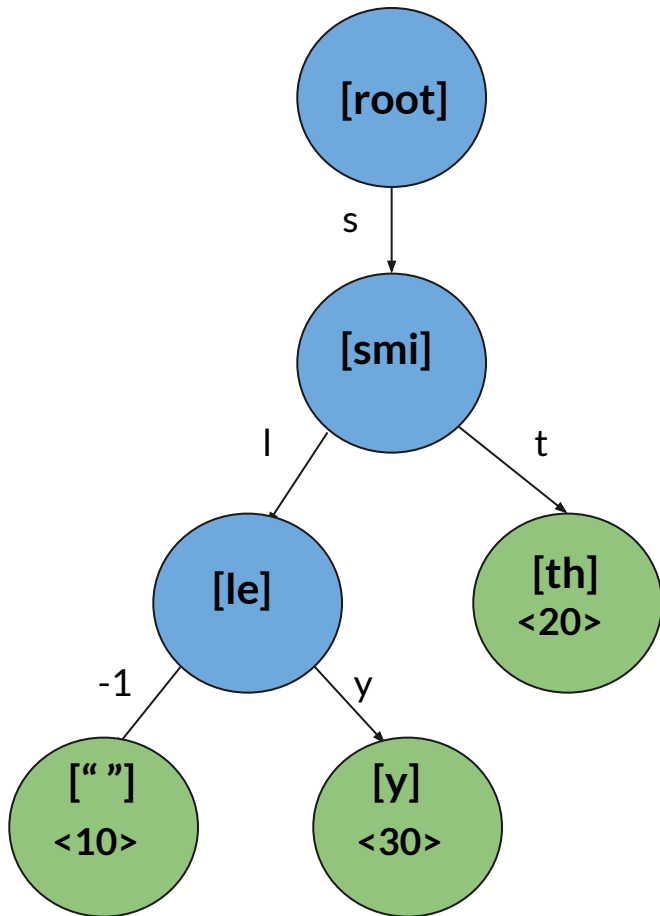
Benchmark against

GPT: Generalized Prefix Tree     RB: red-black tree
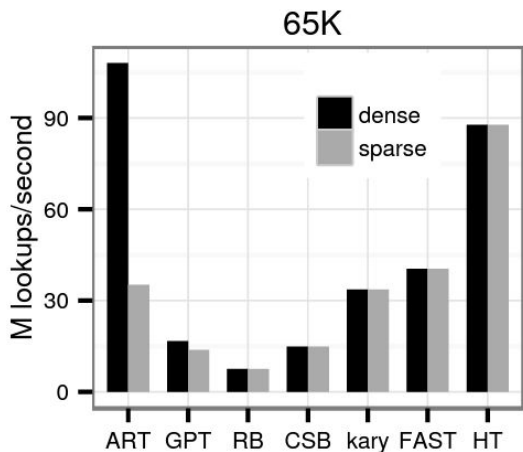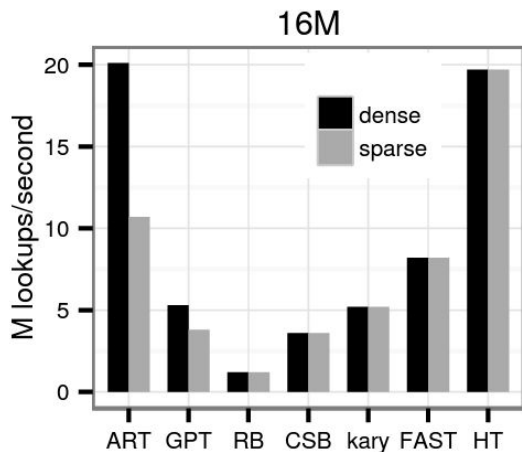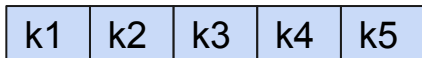
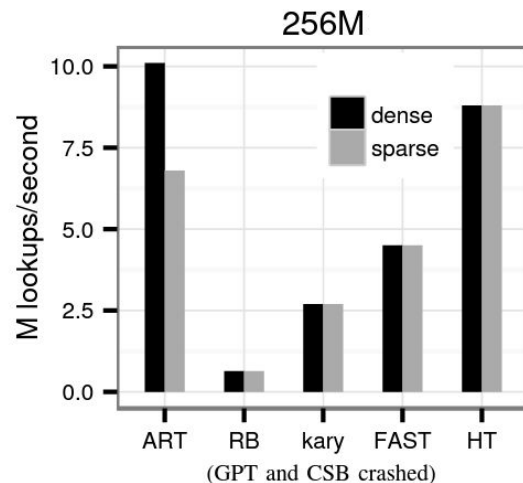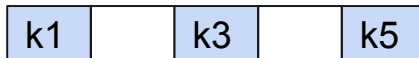CSB: Cache-Sensitive B+-tree     HT: hash table

kary: k-ary search tree     FAST: Fast Architecture Sensitive Tree



Dense keys: Range(1, n)

| k1 | k2 | k3 | k4 | k5 |

Sparse keys: Each bit be randomly 0 or 1

| k1 | | k3 | | k5 |

HT: Hash Table

FAST: Fast Architecture Sensitive Tree

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 |

0 miss prediction branches for ART dense keys!
0.5x cache miss rate for dense keys

To add temporal locality,
Access skewed data to utilize the processor cache

HT: Hash Table

FAST: Fast Architecture Sensitive Tree

Competing Memory Access

HT not affected!



Hit% increases

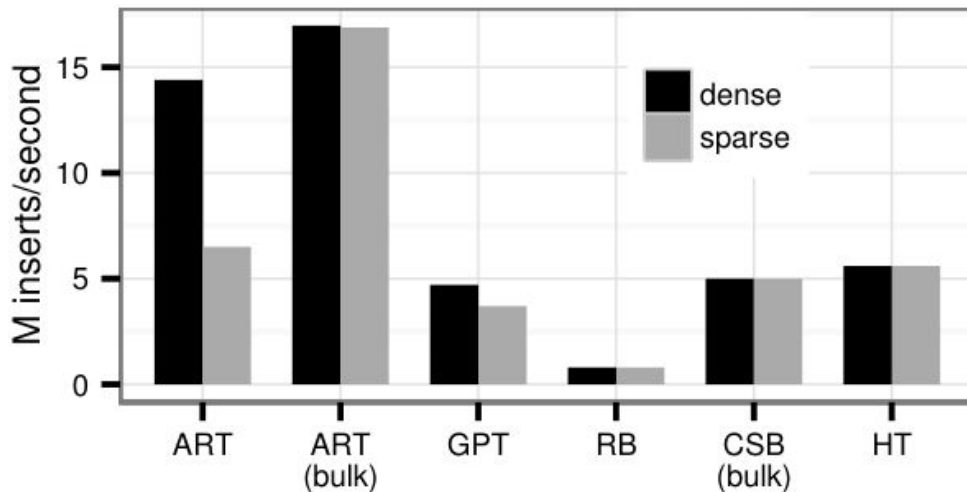ART is more sensitive

GPT: Generalized Prefix Tree     RB: red-black tree

CSB: Cache-Sensitive B+-tree     HT: hash table

Insertion Performance For 16M keys
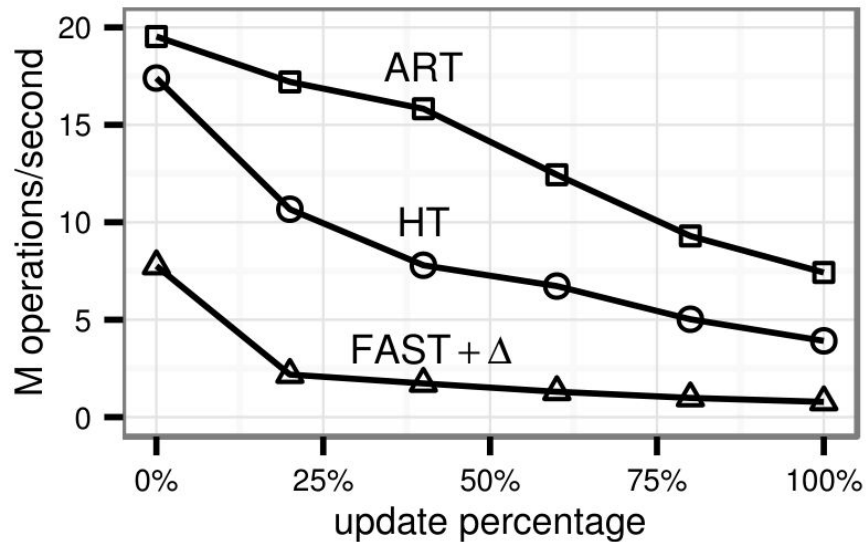


Bulk insertion: Sparse keys -> Dense keys
2.5x faster!

# Evaluation - Search & Update

Random Insertion, Deletion, Lookups

HT: Hash Table

FAST: Fast Architecture Sensitive Tree



Δ is RB tree to help FAST store differences and merge later

## TPC-C OLTP Workload

RB: red-black tree          HT: hash table

Merchandising company

Orders, Customers, Stocks

Manage, Sell, Distribute products



ART ~4x

HT + RB ~2x

RB ~1x

rehashing

| # | Relation | Cardinality | Attribute Types | Space |
|---|----------|-------------|-----------------|-------|
| 1 | item | 100,000 | int | 8.1 |
| 2 | customer | 150,000 | int,int,int | 8.3 |
| 3 | customer | 150,000 | int,int,varchar(16),varchar(16),TID | 32.6 |
| 4 | stock | 500,000 | int,int | 8.1 |
| 5 | order | 22,177,650 | int,int,int | 8.1 |
| 6 | order | 22,177,650 | int,int,int,int,TID | 24.9 |
| 7 | orderline | 221,712,415 | int,int,int,int | 16.8 |

Path compression is good for common prefix!
Both are good for long strings.



Why?
Long unique string!

default
+lazy expansion
+path compression

tree height

index #

Common prefix

# Discussion

Thank You!

Q&A