# 15-721
## DATABASE SYSTEMS

Lecture #08 – Latch-free OLTP Indexes
(Part II)

@Andy_Pavlo // Carnegie Mellon University // Spring 2017

# TODAY'S AGENDA

Bw-Tree Index

ART Index

Profiling in Peloton

# OBSERVATION

We cannot have reverse pointers in a latch-free concurrent Skip List because CaS can only update a single address at a time.

# BW-TREE

Latch-free B+Tree index
→ Threads never need to set latches or block.

**Key Idea #1: Deltas**
→ No updates in place
→ Reduces cache invalidation.

**Key Idea #2: Mapping Table**
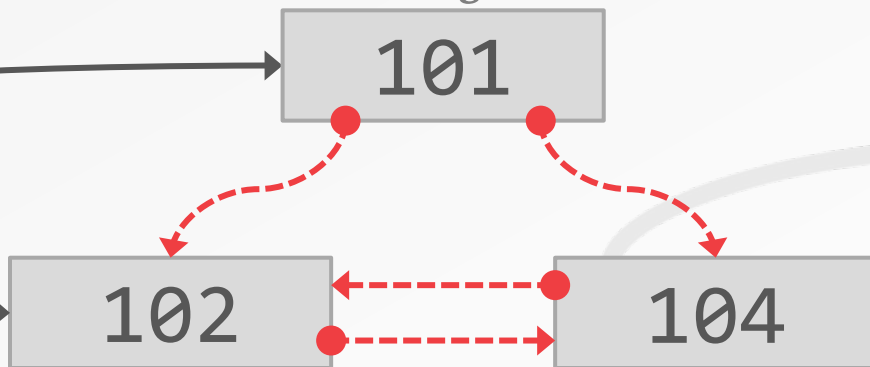→ Allows for CAS of physical locations of pages.

THE BW-TREE: A B-TREE FOR NEW HARDWARE
*ICDE 2013*

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: MAPPING TABLE

**Mapping Table**

| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | |
| 104 | ● |

*Index Page*

101

102    104

Logical
Pointer  - - - →

Physical
Pointer  ——→

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: MAPPING TABLE

**Mapping Table**

| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | |
| 104 | ● |

*Index Page*

102    104

102

104

Logical Pointer

Physical Pointer

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

Each update to a page produces a new delta.

Page 102

*Logical Pointer* ---→

*Physical Pointer* ──→

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | ● |
| 103 | |
| 104 | |

Each update to a page produces a new delta.

▲`Insert 50`

`Page 102`

*Logical Pointer* --->

*Physical Pointer* ⟶

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | ● |
| 103 | |
| 104 | |

*Logical Pointer* ----▶

*Physical Pointer* ——▶

▲ **Insert 50**

Page 102

Each update to a page produces a new delta.

Delta physically points to base page.

Source: Justin Levandoski

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | ● |
| 103 | |
| 104 | |

▲ **Insert 50**

Page 102

*Logical Pointer*  ----→
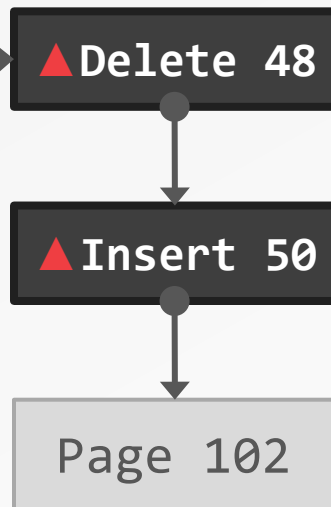
*Physical Pointer*  ——→

Each update to a page produces a new delta.

Delta physically points to base page.

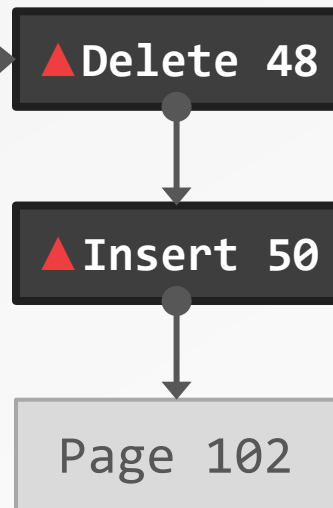Install delta address in physical address slot of mapping table using CAS.

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲ **Insert 50**

Page 102

*Logical Pointer* ----►

*Physical Pointer* ──►

Each update to a page produces a new delta.

Delta physically points to base page.

Install delta address in physical address slot of mapping table using CAS.

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲ Delete 48

▲ Insert 50

Page 102

*Logical Pointer* `----->`

*Physical Pointer* `----->`

Each update to a page produces a new delta.

Delta physically points to base page.

Install delta address in physical address slot of mapping table using CAS.

# BW-TREE: DELTA UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲ Delete 48

▲ Insert 50

Page 102

*Logical Pointer* ---->

*Physical Pointer* ---->

Each update to a page produces a new delta.

Delta physically points to base page.

Install delta address in physical address slot of mapping table using CAS.

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: SEARCH

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Delete 48

▲Insert 50

Page 102

Traverse tree like a regular B+tree.

*Logical Pointer* - - - ➤

*Physical Pointer* ——➤

# BW-TREE: SEARCH

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Delete 48

▲Insert 50

Page 102

*Logical Pointer* ---->

*Physical Pointer* ──>

Traverse tree like a regular B+tree.

If mapping table points to delta chain, stop at first occurrence of search key.

# BW-TREE: SEARCH

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | ● |
| 103 | |
| 104 | |

▲ Delete 48

▲ Insert 50

Page 102

*Logical Pointer* ----▶

*Physical Pointer* ──▶

Traverse tree like a regular B+tree.

If mapping table points to delta chain, stop at first occurrence of search key.

Otherwise, perform binary search on base page.

# BW-TREE: CONTENTION UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 |      |
| 102 |      |
| 103 |      |
| 104 |      |

▲ Insert 50

Page 102

Threads may try to install updates to same state of the page.

*Logical Pointer* ---→

*Physical Pointer* ——→

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: CONTENTION UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Delete 48

▲Insert 16

▲Insert 50

Page 102

Threads may try to install updates to same state of the page.

*Logical Pointer* - - - ▶

*Physical Pointer* ──▶

# BW-TREE: CONTENTION UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Delete 48    ▲Insert 16

▲Insert 50

Page 102

*Logical Pointer* ---→

*Physical Pointer* ──→

Threads may try to install updates to same state of the page.

Winner succeeds, any losers must retry or abort

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: CONTENTION UPDATES

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Delete 48

▲Insert 16

▲Insert 50

Page 102

Threads may try to install updates to same state of the page.

Winner succeeds, any losers must retry or abort

*Logical Pointer* ----→

*Physical Pointer* ——→

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: CONTENTION UPDATES

*Mapping Table*



| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Delete 48

▲Insert 16

▲Insert 50

Page 102

Threads may try to install updates to same state of the page.

Winner succeeds, any losers must retry or abort

*Logical Pointer* ----▶

*Physical Pointer* ───▶

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: DELTA TYPES

**Record Update Deltas**
→ Insert/Delete/Update of record on a page

**Structure Modification Deltas**
→ Split/Merge information

# BW-TREE: CONSOLIDATION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲ **Insert 55**

▲ **Delete 48**

▲ **Insert 50**

Page 102

Consolidate updates by creating new page with deltas applied.

*Logical Pointer* ---->

*Physical Pointer* ——>

# BW-TREE: CONSOLIDATION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲**Insert 55**

▲**Delete 48**

▲**Insert 50**

Page 102

Consolidate updates by creating new page with deltas applied.

*Logical Pointer* ----►

*Physical Pointer* ──►

New 102

# BW-TREE: CONSOLIDATION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Insert 55

▲Delete 48

▲Insert 50

Page 102

Consolidate updates by creating new page with deltas applied.

*Logical Pointer* ╌╌╌▶

*Physical Pointer* ─────▶

New 102

▲Insert 50

# BW-TREE: CONSOLIDATION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 |      |
| 102 |      |
| 103 |      |
| 104 |      |

▲ **Insert 55**

▲ **Delete 48**

▲ **Insert 50**

Page 102

*Logical Pointer* ---->

*Physical Pointer* ——>

New 102

Consolidate updates by creating new page with deltas applied.

CAS-ing the mapping table address ensures no deltas are missed.

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: CONSOLIDATION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲ **Insert 55**

▲ **Delete 48**

▲ **Insert 50**

Page 102

*Logical Pointer* - - - ▶

*Physical Pointer* ────▶

New 102

Consolidate updates by creating new page with deltas applied.
CAS-ing the mapping table address ensures no deltas are missed.

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: CONSOLIDATION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

*Logical Pointer* ----▶

*Physical Pointer* ──▶

New 102

▲ Insert 55
▲ Delete 48
▲ Insert 50
Page 102

Consolidate updates by creating new page with deltas applied.

CAS-ing the mapping table address ensures no deltas are missed.

Old page + deltas are marked as garbage.

# BW-TREE: GARBAGE COLLECTION

Operations are tagged with an **epoch**
→ Each epoch tracks the threads that are part of it and the objects that can be reclaimed.
→ Thread joins an epoch prior to each operation and post objects that can be reclaimed for the current epoch (not necessarily the one it joined)

Garbage for an epoch reclaimed only when all threads have exited the epoch.

# BW-TREE: GARBAGE COLLECTION

**Mapping Table**

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲ Insert 55

▲ Delete 48

▲ Insert 50

CPU1

Page 102

**Epoch Table**

CPU1

*Logical Pointer* ---->

*Physical Pointer* ---->

New 102

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: GARBAGE COLLECTION

# BW-TREE: GARBAGE COLLECTION

# BW-TREE: GARBAGE COLLECTION

# BW-TREE: GARBAGE COLLECTION



*Mapping Table*

*Epoch Table*

| PID | Addr |
|-----|------|
| 101 | |
| 102 | |
| 103 | |
| 104 | |

▲Insert 55

▲Delete 48

▲Insert 50

Page 102

*Logical Pointer*

*Physical Pointer*

New 102

▲Insert 55
▲Delete 48
▲Insert 50
Page 102

# BW-TREE: GARBAGE COLLECTION

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 |      |
| 102 | ● |
| 103 |      |
| 104 |      |

*Logical Pointer* - - - ▶

*Physical Pointer* ——▶

New 102

*Epoch Table*

▲Insert 55

▲Delete 48

▲Insert 50

Page 102

# BW-TREE: STRUCTURE MODIFICATIONS

**Split Delta Record**
→ Mark that a subset of the base page's key range is now located at another page.
→ Use a logical pointer to the new page.

**Separator Delta Record**
→ Provide a shortcut in the modified page's parent on what ranges to find the new page.

# BW-TREE: STRUCTURE MODIFICATIONS

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*



| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | ● |
| 104 | ● |
| 105 | |

*Logical Pointer* - - →

*Physical Pointer* →

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | ● |
| 104 | ● |
| 105 | ● |

101

102

103

104

105

*Logical Pointer* - - →

*Physical Pointer* →

① ②

③ ④ ⑤ ⑥

⑦ ⑧

⑤ ⑥

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

# BW-TREE: STRUCTURE MODIFICATIONS

# BW-TREE: STRUCTURE MODIFICATIONS

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

# BW-TREE: STRUCTURE MODIFICATIONS

# BW-TREE: STRUCTURE MODIFICATIONS

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | ● |
| 104 | ● |
| 105 | ● |

▲ Separator
[5,7)

101

[-∞,3)  [3,7)  [7,∞)

▲ Split

1 2
102

3 4 XX
103

7 8
104

5 6
105

*Logical Pointer* --→

*Physical Pointer* —→

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | ● |
| 104 | ● |
| 105 | ● |

▲ Separator    [5,7)

101

[-∞,3)  [3,7)  [7,∞)

▲ Split

① ②    ③ ④ XX    ⑦ ⑧

102    103    104

⑤ ⑥

105

*Logical Pointer* ----→

*Physical Pointer* ——→

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: STRUCTURE MODIFICATIONS

*Mapping Table*

| PID | Addr |
|-----|------|
| 101 | ● |
| 102 | ● |
| 103 | ● |
| 104 | ● |
| 105 | ● |

▲ Separator

[5,7)

101

[-∞,3)  [3,7)  [7,∞)

▲ Split

1  2

102

3  4  XXX

103

7  8

104

5  6

105

*Logical Pointer*

*Physical Pointer*

CARNEGIE MELLON
DATABASE GROUP

# BW-TREE: PERFORMANCE

Processor: 1 socket, 4 cores w/ 2×HT



Source: Justin Levandoski

# ADAPATIVE RADIX TREE (ART)

Uses digital representation of keys to examine prefixes one-by-one instead of comparing entire key.
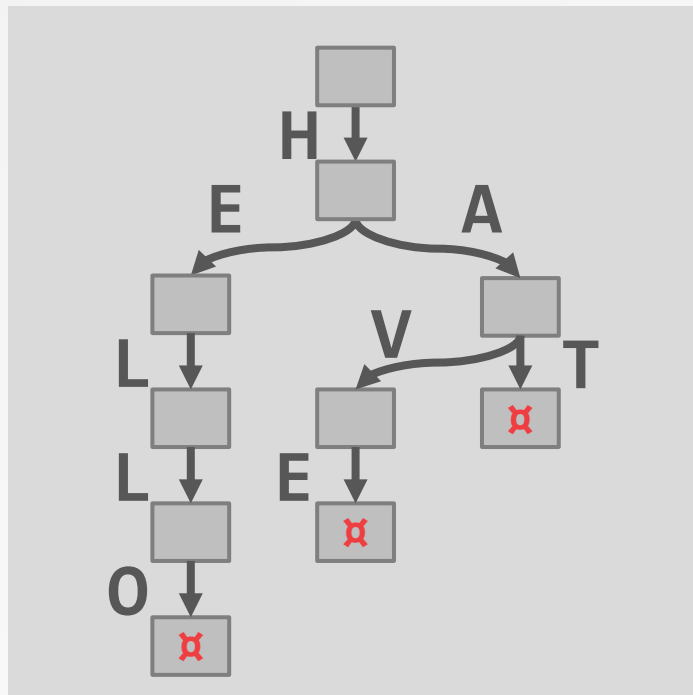
Radix trees properties:
→ The height of the tree depends on the length of keys.
→ Does not require rebalancing
→ The path to a leaf node represents the key of the leaf
→ Keys are stored implicitly and can be reconstructed from paths.

THE ADAPTIVE RADIX TREE: ARTFUL
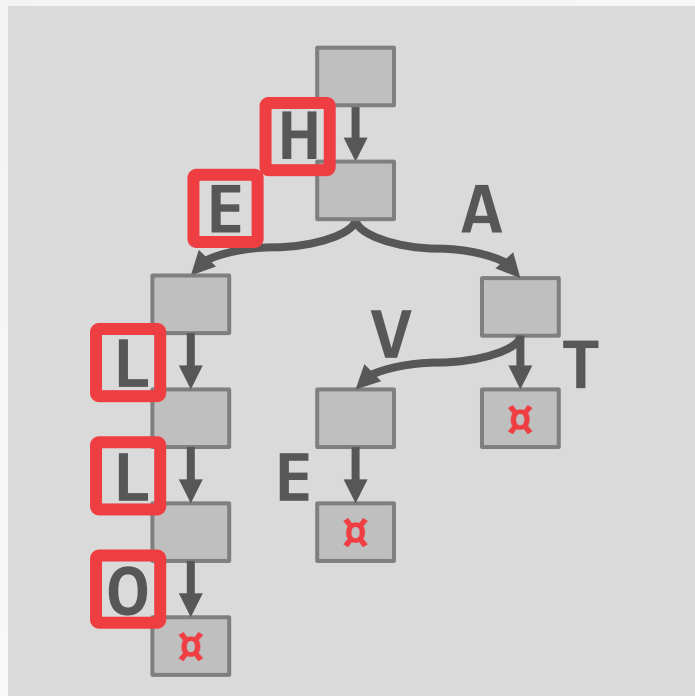INDEXING FOR MAIN-MEMORY DATABASES
*ICDE 2013*

CARNEGIE MELLON
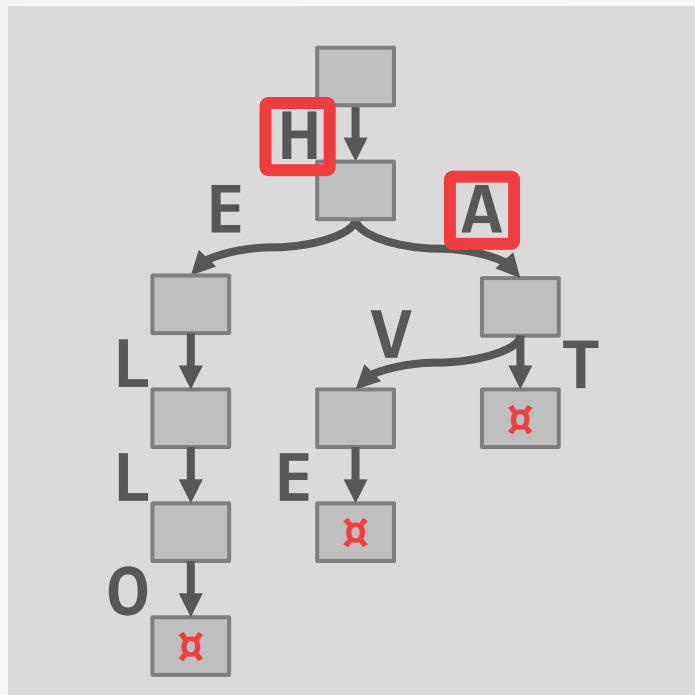DATABASE GROUP

# TRIE VS. RADIX TREE

*Trie*



**Keys:** HELLO, HAT, HAVE

# TRIE VS. RADIX TREE

*Trie*



Keys: HELLO, HAT, HAVE

# TRIE VS. RADIX TREE

*Trie*



**Keys:** HELLO, HAT, HAVE
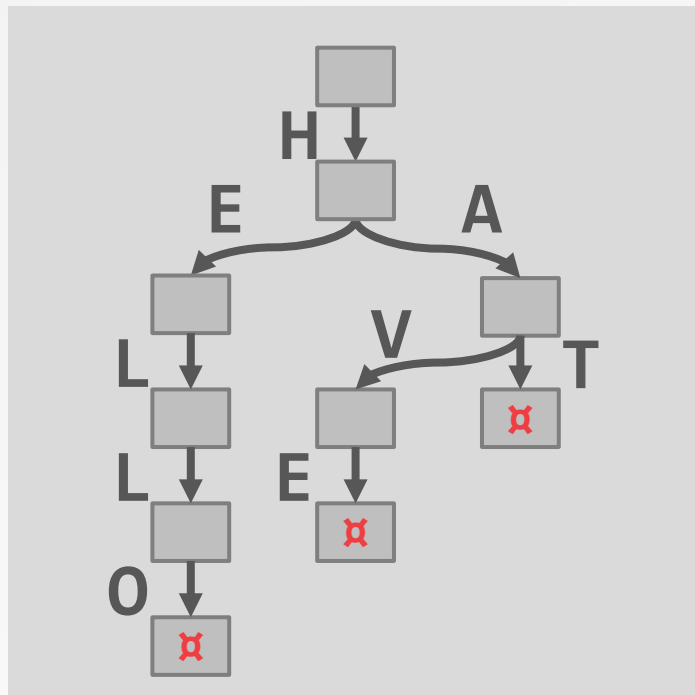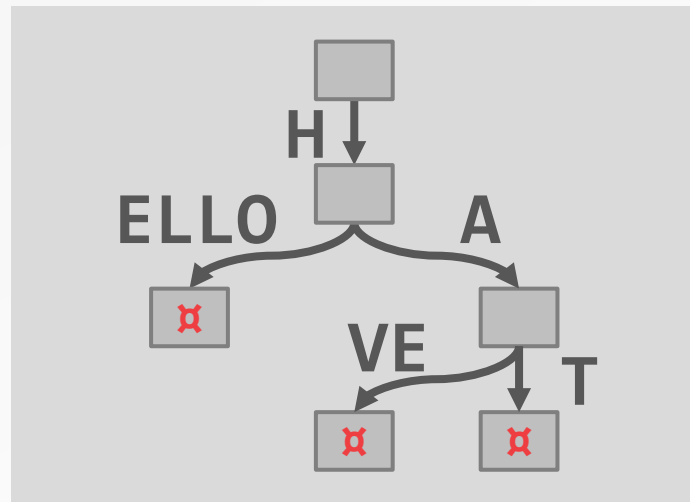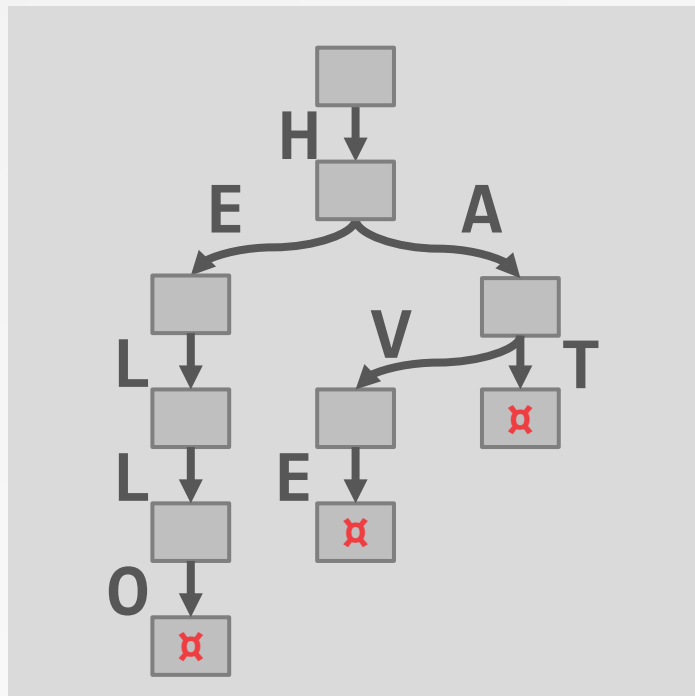
# TRIE VS. RADIX TREE

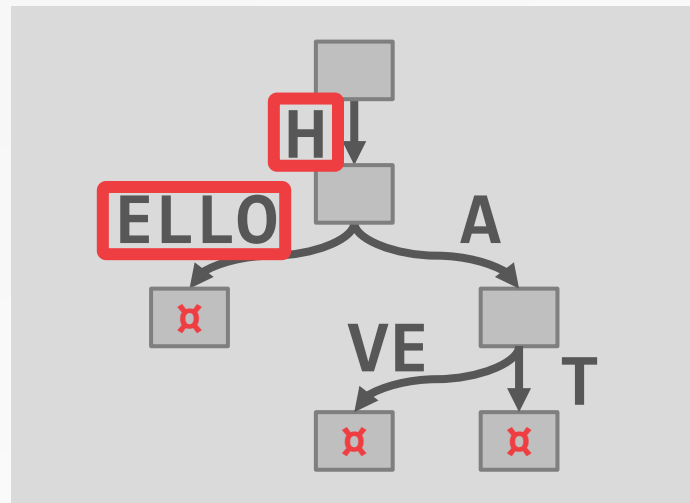

*Trie*

*Radix Tree*

**Keys: HELLO, HAT, HAVE**

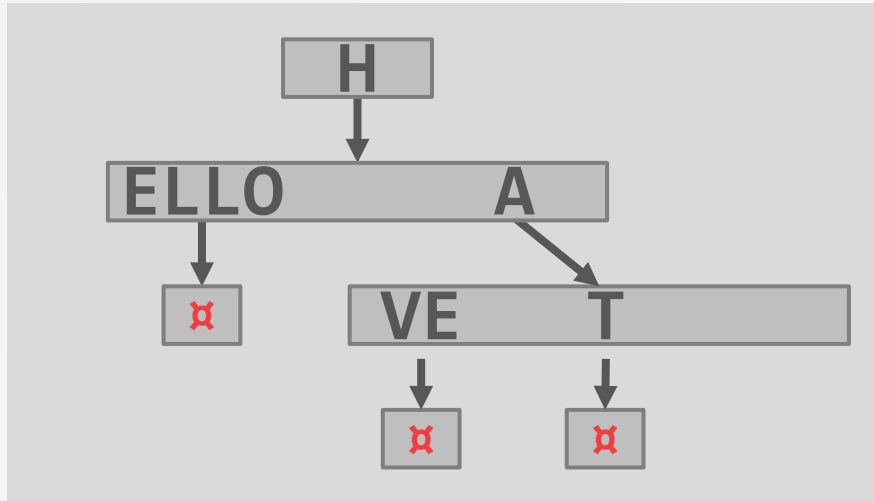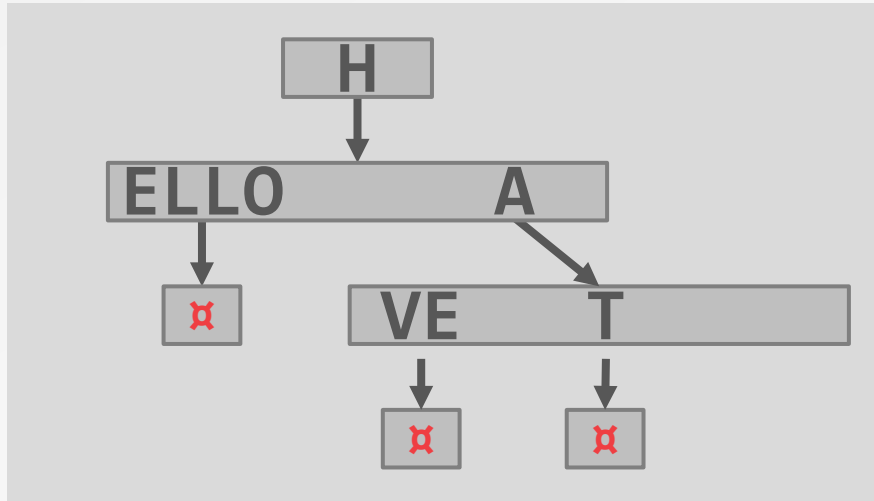# TRIE VS. RADIX TREE

*Trie*

*Radix Tree*



Keys: HELLO, HAT, HAVE

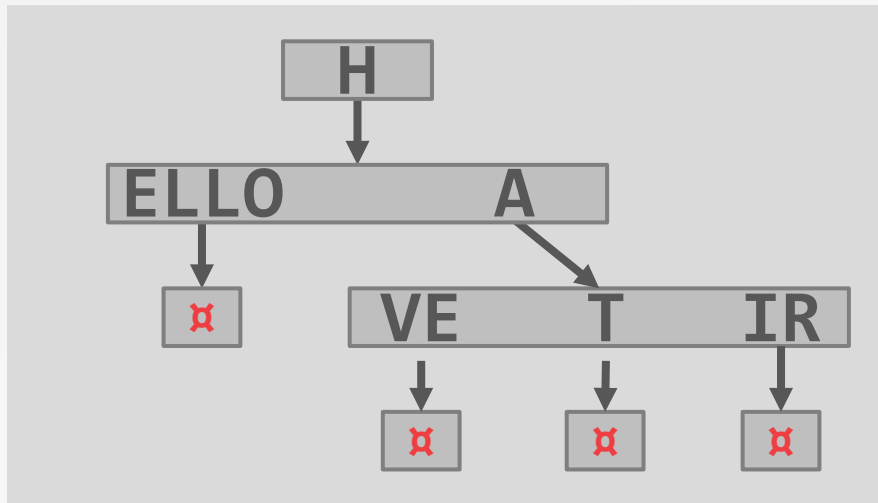# ART INDEX: MODIFICATIONS

# ART INDEX: MODIFICATIONS



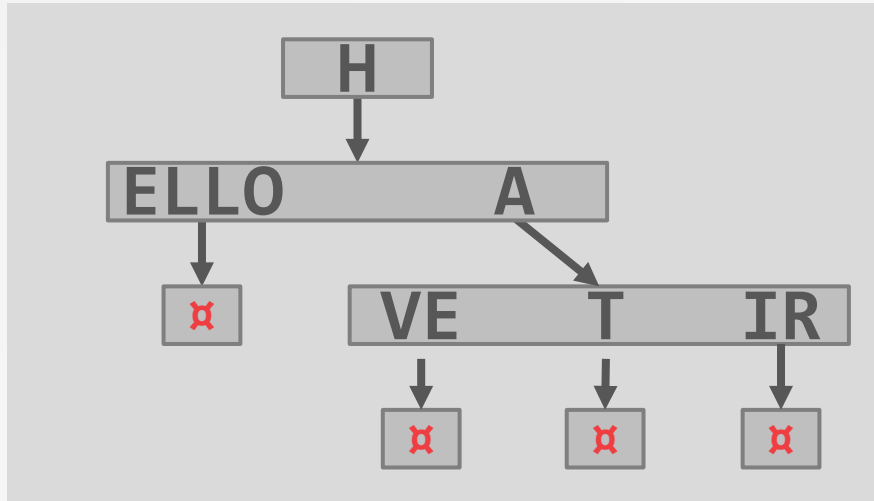**Operation:** Insert **HAIR**

# ART INDEX: MODIFICATIONS



**Operation:** Insert **HAIR**

# ART INDEX: MODIFICATIONS



**Operation:** Insert **HAIR**

**Operation:** Delete **HAT**, **HAVE**

# ART INDEX: MODIFICATIONS



**Operation:** Insert **HAIR**

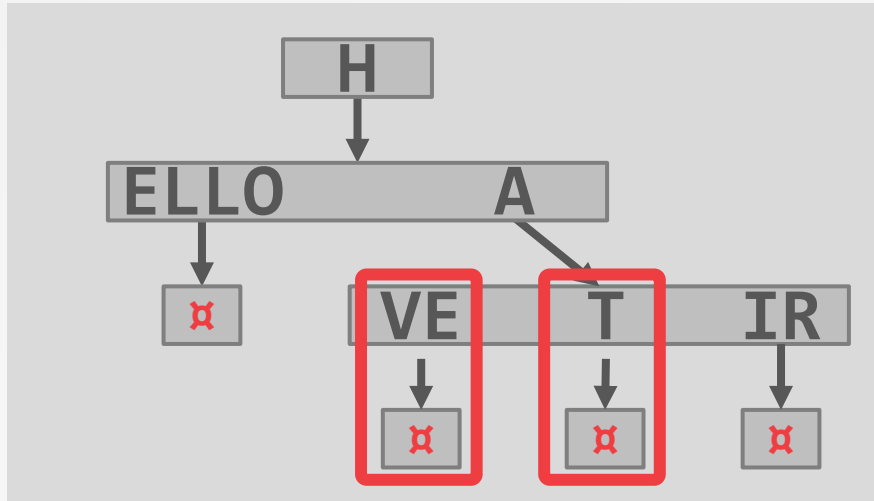**Operation:** Delete **HAT**, **HAVE**

# ART INDEX: MODIFICATIONS
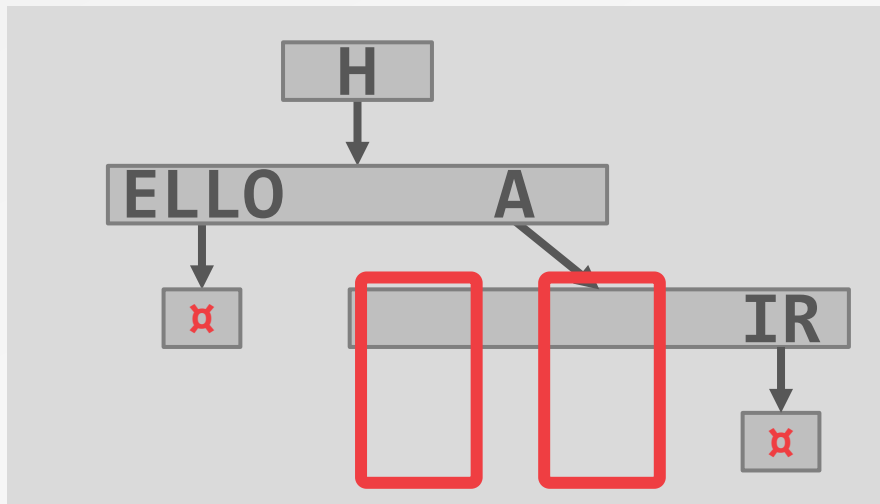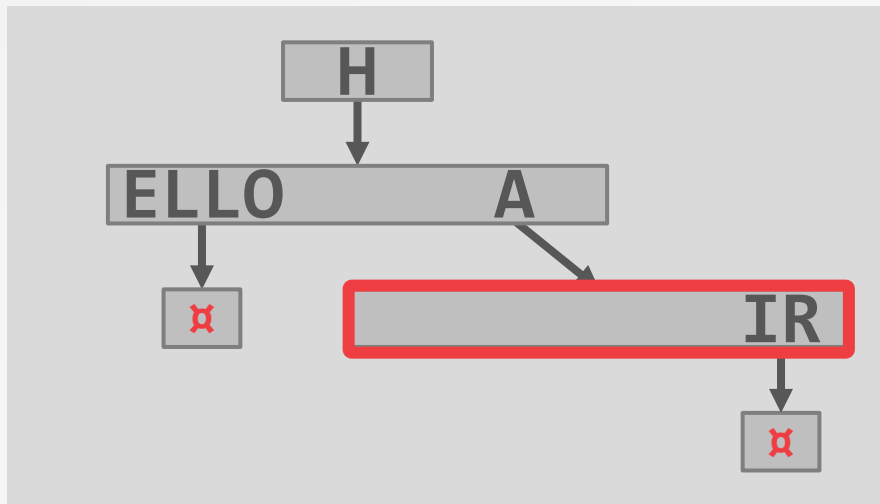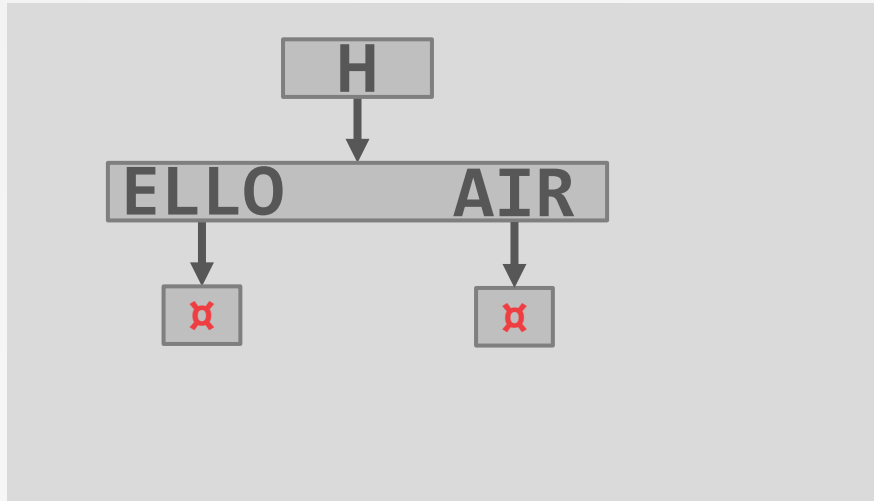


**Operation:** Insert **HAIR**

**Operation:** Delete **HAT**, **HAVE**

# ART INDEX: MODIFICATIONS



**Operation:** Insert **HAIR**

**Operation:** Delete **HAT**, **HAVE**

# ART INDEX: MODIFICATIONS



**Operation:** Insert **HAIR**

**Operation:** Delete **HAT**, **HAVE**

# ART INDEX: BINARY COMPARABLE KEYS

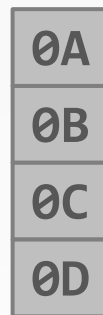Not all attribute types can be decomposed into binary comparable digits for a radix tree.
→ **Unsigned Integers:** Byte order must be flipped for little endian machines.
→ **Signed Integers:** Flip two's-complement so that negative numbers are smaller than positive.
→ **Floats**: Classify into group (neg vs. pos, normalized vs. denormalized), then store as unsigned integer.
→ **Compound**: Transform each attribute separately.

# ART INDEX: BINARY COMPARABLE KEYS

**Int Key: 168496141**

**Hex Key: 0A 0B 0C 0D**

| 0A |
| -- |
| 0B |
| 0C |
| 0D |

*Big Endian*

| 0D |
| -- |
| 0C |
| 0B |
| 0A |

*Little Endian*

# ART INDEX: BINARY COMPARABLE KEYS



**Int Key: 168496141**

**Hex Key: 0A 0B 0C 0D**

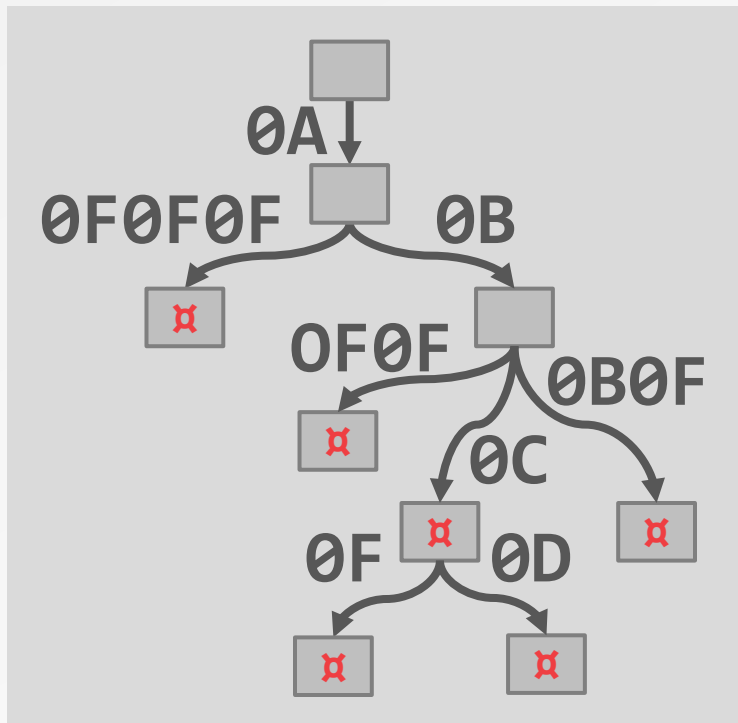| Big Endian | Little Endian |
| --- | --- |
| 0A | 0D |
| 0B | 0C |
| 0C | 0B |
| 0D | 0A |

*Big Endian*          *Little Endian*

# ART INDEX: BINARY COMPARABLE KEYS

**Int Key: 168496141**

**Hex Key: 0A 0B 0C 0D**

| Big Endian |
|:---:|
| 0A |
| 0B |
| 0C |
| 0D |

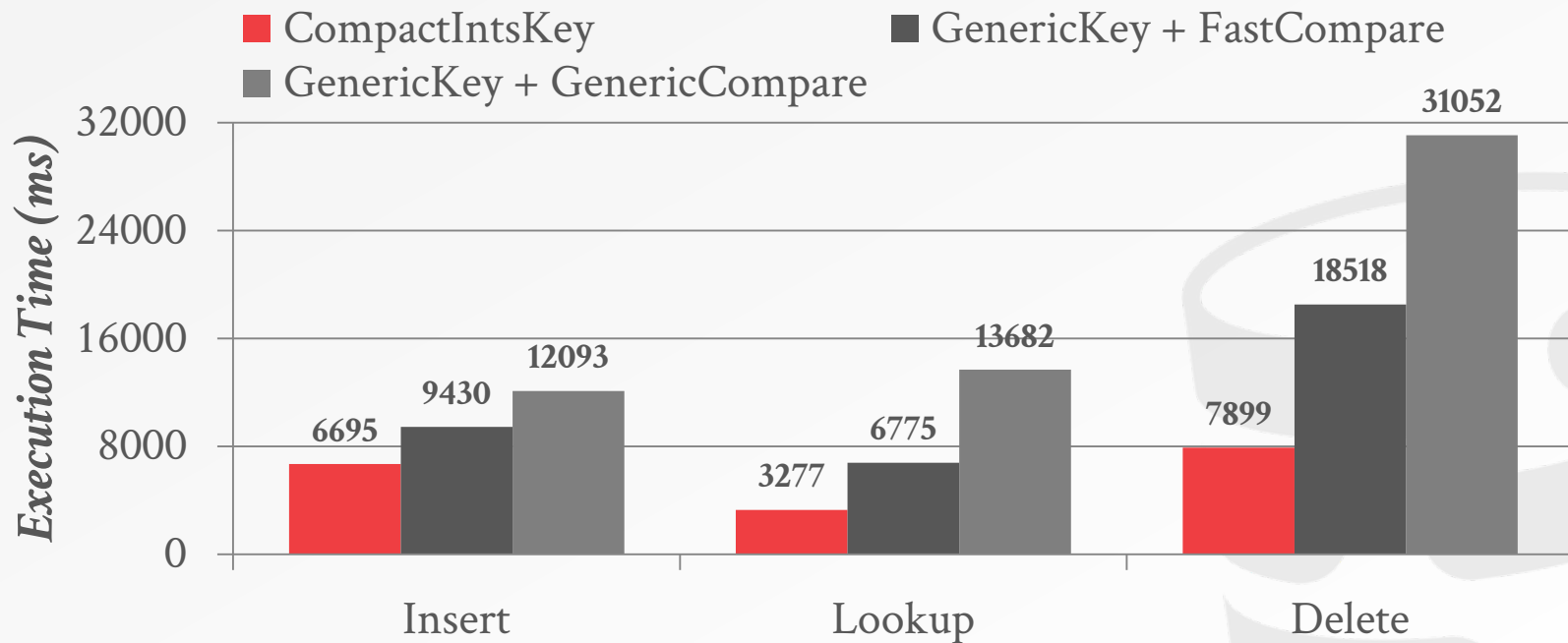| Little Endian |
|:---:|
| 0D |
| 0C |
| 0B |
| 0A |

*Big Endian*    *Little Endian*

# BINARY COMPRABLE KEYS

Peloton w/ Bw-Tree Index
Data Set: 10m keys (three 64-bit ints)

# CONCURRENT ART INDEX

HyPer's ART is **<u>not</u>** latch-free.

Optimistic crabbing scheme where writers are not blocked on readers.
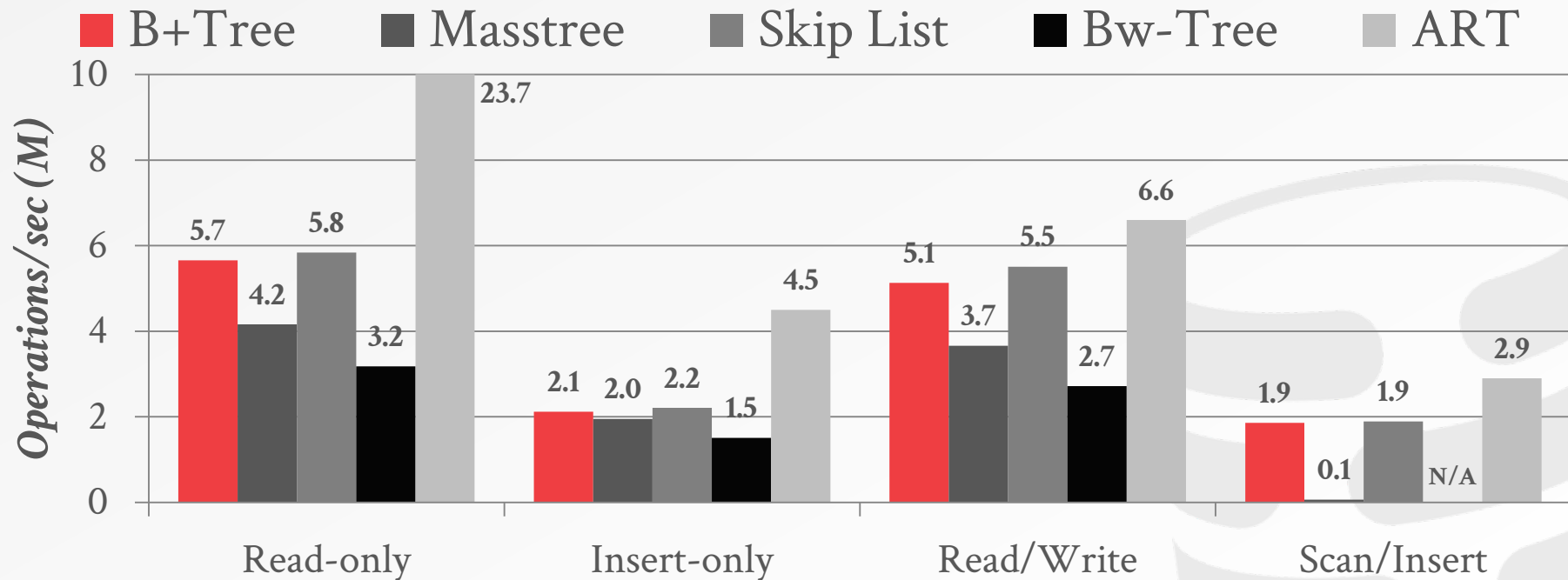→ Writers increment counter when they acquire latch.
→ Readers can proceed if a node's latch is available.
→ It then checks whether the latch's counter has changed from when it checked the latch.

THE ART OF PRACTICAL SYNCHRONIZATION
*DaMoN 2016*

CARNEGIE MELLON
DATABASE GROUP

# SINGLE-THREADED PERFORMANCE

Data Set: 30m Random 64-bit Integers



Source: Huanchen Zhang

# PARTING THOUGHTS

Bw-Tree is probably the most dank latch-free index in recent years.

ART has amazing performance. Need to understand it better.

CARNEGIE MELLON
DATABASE GROUP

*ANDY'S*
# TIPS FOR PROFILING

# MOTIVATION

Consider a program with functions **foo** and **bar**.

How can we speed it up with only a debugger ?
→ Randomly pause it during execution
→ Collect the function call stack

# RANDOM PAUSE METHOD

Consider this scenario
→ Collected 10 call stack samples
→ Say 6 out of the 10 samples were in **foo**

What percentage of time was spent in **foo**?
→ Roughly 60% of the time was spent in **foo**
→ Accuracy increases with # of samples

# AMDAHL'S LAW

Say we optimized **foo** to run 2 times faster

What's the expected overall speedup ?
→ 60% of time spent in **foo** drops in half
→ 40% of time spent in **bar** unaffected

→ *p* = percentage of time spent in optimized task
→ *s* = speed up for the optimized task

→ Overall speedup = ——— = 1.4 times faster

# AMDAHL'S LAW

1  0.6 2 +0.4 1 1  0.6 2 +0.4  0.6 2 0.6 0.6 2 2 0.6
2 +0.4 1  0.6 2 +0.4  = 1.4 times faster

1  0.6 2 +0.4 1 1  0.6 2 +0.4  0.6 2 0.6 0.6 2 2 0.6
2 +0.4 1  0.6 2 +0.4  = 1.4 times faster

1     +(1− )                    +(1−   ) 1
+(1− )

*Say* we optimized **foo** to run 2 times faster

What's the expected overall speedup ?
→ 60% of time spent in **foo** drops in half
→ 40% of time spent in **bar** unaffected

CARNEGIE MELLON
DATABASE GROUP

# PROFILING TOOLS FOR REAL

**Choice #1: Valgrind**
→ Heavyweight instrumentation framework with a lot of tools
→ Sophisticated visualization tools

**Choice #2: Perf**
→ Lightweight tool that can record different kinds of events
→ Console-oriented visualization tools

# CHOICE #1: VALGRIND

Instrumentation framework for building dynamic analysis tools
→ **memcheck**: a memory error detector
→ **callgrind**: a call-graph generating profiler

Using **callgrind** to profile the index test and Peloton in general:

```
$ valgrind --tool=callgrind --trace-children=yes
./tests/skiplist_index_test

$ valgrind --tool=callgrind --trace-children=yes
./bin/peloton &> /dev/null&
```
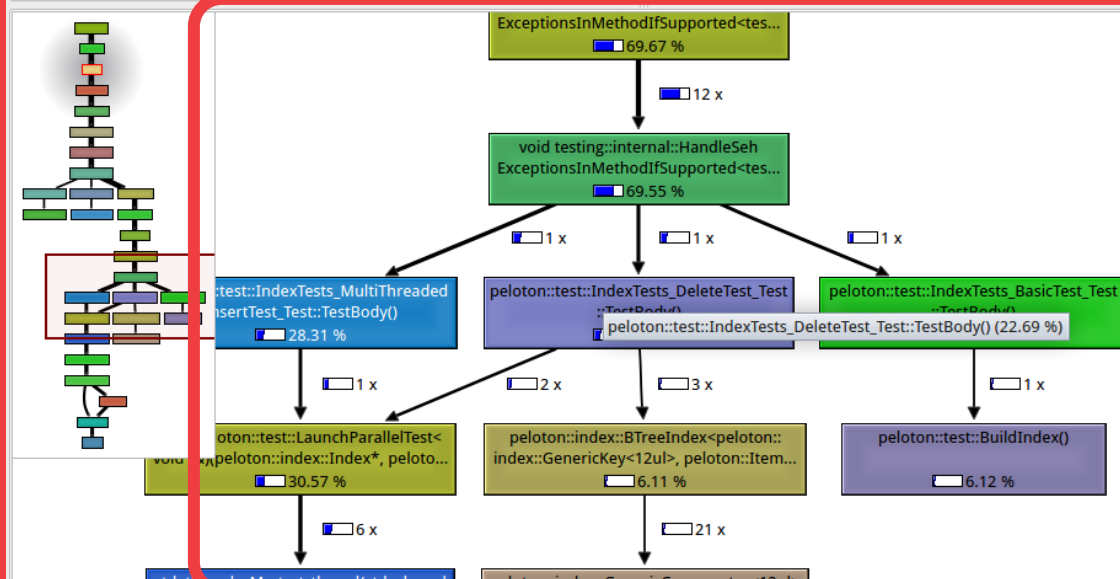
# KCACHEGRIND

Profile data visualization tool

```
$ kcachegrind callgrind.out.12345
```

CARNEGIE MELLON
DATABASE GROUP

# CHOICE #2: PERF

Tool for using the performance counters subsystem in Linux.
→ **-e** = sample the event `cycles` at the user level only
→ **-c** = collect a sample every 2000 occurrences of event

```
$ perf record -e cycles:u -c 2000
./tests/skiplist_index_test
```
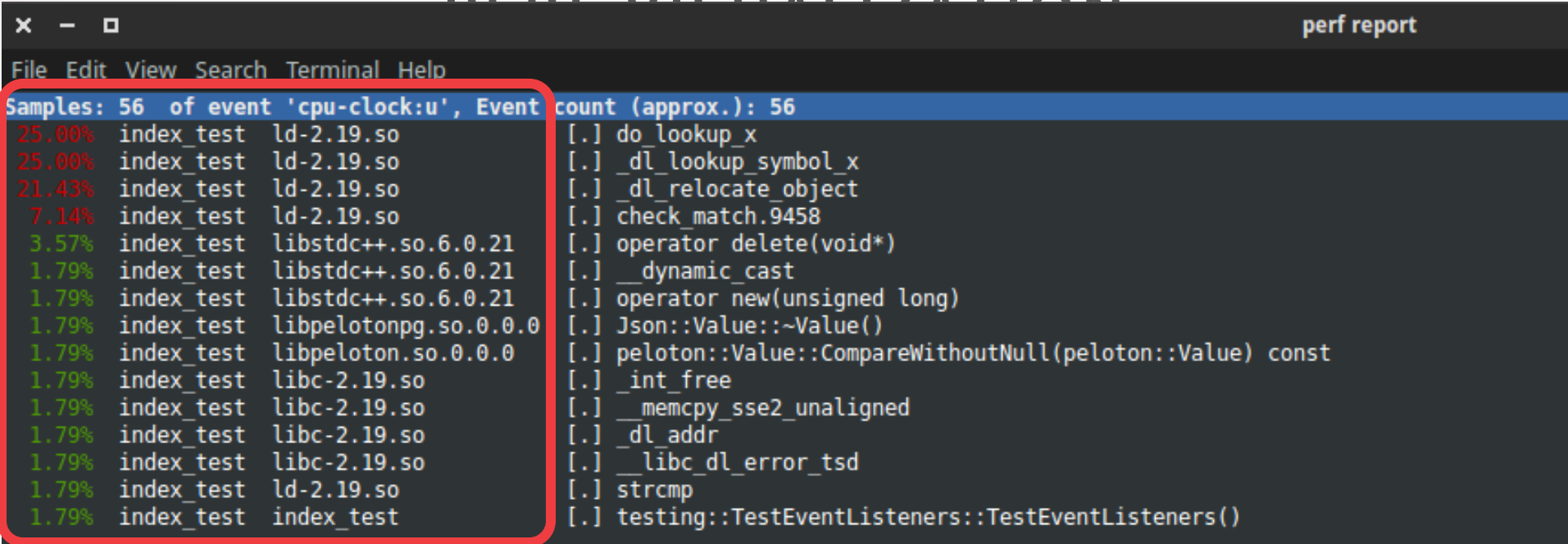
Uses counters for tracking events
→ On counter overflow, the kernel records a sample
→ Sample contains info about program execution

# PERF VISUALIZATION

We can also use **perf** to visualize the generated profile for our application.

```
$ perf report
```

# PERF VISUALIZATION



Cumulative Time Distribution

# PERF EVENTS

Supports several other events like:
→ L1-dcache-load-misses
→ branch-misses

To see a list of events:

```
$ perf list
```

Another usage example:

```
$ perf record -e cycles,LLC-load-misses -c 2000
./tests/skiplist_index_test
```

# REFERENCES

**Valgrind**
→ The Valgrind Quick Start Guide
→ Callgrind
→ Kcachegrind
→ Tips for the Profiling/Optimization process

**Perf**
→ Perf Tutorial
→ Perf Examples
→ Perf Analysis Tools

# NEXT CLASS

Indexing for OLAP workloads.
→ More from Microsoft Research…