

# **DATABASE SYSTEM IMPLEMENTATION**

## GT 4420/6422 // SPRING 2019 // @JOY\_ARULRAJ

**LECTURE #12: OLTP INDEXES (PART II)** 

CREATING THE NEXT°

### TODAY'S AGENDA

B+Tree Overview Index Implementation Issues ART Index

### LOGISTICS

**Reminder:** Problem set due on Feb 21<sup>st</sup>.

**Reminder:** Mid-term Exam on Feb 26<sup>th</sup>.

Reminder: Project Proposals due on Feb 28<sup>th</sup>.

## B+TREE

A **B+Tree** is a self-balancing tree data structure that keeps data sorted and allows searches, sequential access, insertions, and deletions in **O(log n)**.

- $\rightarrow$  Generalization of a binary search tree in that a node can have more than two children.
- $\rightarrow$  Optimized for systems that read and write large blocks of data.



## **B+TREE PROPERTIES**

A B+tree is an *M*-way search tree with the following properties:

- $\rightarrow$  It is perfectly balanced (i.e., every leaf node is at the same depth).
- → Every inner node other than the root, is at least half-full
  M/2-1 ≤ #keys ≤ M-1
- $\rightarrow$  Every inner node with **k** keys has **k+1** non-null children













## NODES

Every node in the B+Tree contains an array of key/value pairs.

- $\rightarrow$  The keys will always be the column or columns that you built your index on
- $\rightarrow$  The values will differ based on whether the node is classified as **<u>inner nodes</u>** or <u>**leaf nodes.**</u>

The arrays are (usually) kept in sorted key order.

## LEAF NODE VALUES

#### Approach #1: Record Ids

 $\rightarrow$  A pointer to the location of the tuple that the index entry corresponds to.

#### Approach #2: Tuple Data

- $\rightarrow$  The actual contents of the tuple is stored in the leaf node.
- $\rightarrow$  Secondary indexes have to store the record id as their values.

## LEAF NODE VALUES

#### Approach #1: Record Ids

 $\rightarrow$  A pointer to the location of the tuple that the index entry corresponds to. ORACLE<sup>®</sup>

#### Approach #2: Tuple Data

- $\rightarrow$  The actual contents of the tuple is stored in the leaf node.
- $\rightarrow$  Secondary indexes have to store the record id as their values.



**PostgreSQL** 















## **B+TREE INSERT**

Find correct leaf L.

Put data entry into L in sorted order.

If L has enough space, done!

Else, must split L into L and a new node L2

- $\rightarrow$  Redistribute entries evenly, copy up middle key.
- $\rightarrow$  Insert index entry pointing to L2 into parent of L.

To split inner node, redistribute entries evenly, but push up middle key.

### **B+TREE VISUALIZATION**

<u>https://www.cs.usfca.edu/~galles/visualizati</u> <u>on/BPlusTree.html</u>

Source: David Gales (Univ. of San Francisco)

## **B+TREE DELETE**

Start at root, find leaf L where entry belongs. Remove the entry.

- If L is at least half-full, done!
- If L has only M/2-1 entries,
- $\rightarrow$  Try to re-distribute, borrowing from sibling (adjacent node with same parent as L).
- $\rightarrow$  If re-distribution fails, merge L and sibling.

If merge occurred, must delete entry (pointing to L or sibling) from parent of L.

### **B+TREES IN PRACTICE**

Typical Fill-Factor: 67%.  $\rightarrow$  Average Fanout = 2\*100\*0.67 = 134

Typical Capacities:  $\rightarrow$  Height 4: 1334 = 312,900,721 entries  $\rightarrow$  Height 3: 1333 = 2,406,104 entries

Pages per level:  $\rightarrow$  Level 1 = 1 page = 8 KB  $\rightarrow$  Level 2 = 134 pages = 1 MB  $\rightarrow$  Level 3 = 17,956 pages = 140 MB

## **CLUSTERED INDEXES**

The table is stored in the sort order specified by the primary key.

 $\rightarrow$  Can be either heap- or index-organized storage.

Some DBMSs always use a clustered index.

 $\rightarrow$  If a table doesn't include a pkey, the DBMS will automatically make a hidden row id pkey.

Other DBMSs cannot use them at all.

## SELECTION CONDITIONS

The DBMS can use a B+Tree index if the query provides any of the attributes of the search key.

Example: Index on <a,b,c>

- → Supported: (a=5 AND b=3)
- $\rightarrow$  Supported: (b=3).

Not all DBMSs support this.

For hash index, we must have all attributes in search key.

## **B+TREE DESIGN CHOICES**

Node Size Merge Threshold Intra-Node Search Variable Length Keys Non-Unique Indexes



## NODE SIZE

The slower the disk, the larger the optimal node size for a B+Tree.

- $\rightarrow$  HDD ~1MB
- $\rightarrow$  SSD: ~10KB
- $\rightarrow$  In-Memory: ~512B

Optimal sizes can vary depending on the workload  $\rightarrow$  Leaf Node Scans vs. Root-to-Leaf Traversals

### MERGE THRESHOLD

Some DBMSs don't always merge nodes when it is half full.

Delaying a merge operation may reduce the amount of reorganization.

May be better to just let underflows to exist and then periodically rebuild entire tree.

#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### **Approach #3: Interpolation**

#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### **Approach #3: Interpolation**

→ Approximate location of desired key based on known distribution of keys. Find Key=8

#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation


#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation

→ Approximate location of desired key based on known distribution of keys.



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation

→ Approximate location of desired key based on known distribution of keys.



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation

→ Approximate location of desired key based on known distribution of keys.



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### Approach #3: Interpolation

→ Approximate location of desired key based on known distribution of keys.



#### Approach #1: Linear

 $\rightarrow$  Scan node keys from beginning to end.

#### Approach #2: Binary

→ Jump to middle key, pivot left/right depending on comparison.

#### **Approach #3: Interpolation**

→ Approximate location of desired key based on known distribution of keys.



## INDEX IMPLEMENTATION ISSUES

Bulk Insert Pointer Swizzling Prefix Compression Memory Pools Garbage Collection Non-Unique Keys Variable-length Keys Prefix Compression

The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up.

The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up. *Keys: 3, 7, 9, 13, 6, 1* 

The fastest/best way to build a B+Tree is to firstsort the keys and then build the index from thebottom up.Keys: 3, 7, 9, 13, 6, 1

Sorted Keys: 1, 3, 6, 7, 9, 13

The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up. *Keys: 3, 7, 9, 13, 6, 1* 

Sorted Keys: 1, 3, 6, 7, 9, 13



The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up. *Keys: 3, 7, 9, 13, 6, 1* 

Sorted Keys: 1, 3, 6, 7, 9, 13



Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.

Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.





Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.





Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.







Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.







Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.



Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.





Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.



Nodes use page ids to reference other nodes in the index. The DBMS has to get the memory location from the page table during traversal.





## **MEMORY POOLS**

We don't want to be calling **malloc** and **free** anytime we need to add or delete a node.

This could lead to a system call.

- → If you call **malloc** to request 10 bytes of memory, the allocator may invoke the **sbrk** (or **mmap**) system call to request 4K bytes from OS.
- → Then, when you call **malloc** next time to request another 10 bytes, it may not have to issue a system call; instead, it may return a pointer within allocated memory.

## **MEMORY POOLS**

If all the nodes are the same size, then the index can maintain a pool of available nodes.  $\rightarrow$  Insert: Grab a free node, otherwise create a new one.  $\rightarrow$  Delete: Add the node back to the free pool.

Need some policy to decide when to retract the pool size (garbage collection & de-fragmentation).

- $\rightarrow$  Reference Counting
- $\rightarrow$  Epoch-based Reclamation
- $\rightarrow$  Hazard Pointers
- $\rightarrow$  Many others...



- $\rightarrow$  Reference Counting
- $\rightarrow$  Epoch-based Reclamation
- $\rightarrow$  Hazard Pointers
- $\rightarrow$  Many others...



- $\rightarrow$  Reference Counting
- $\rightarrow$  Epoch-based Reclamation
- $\rightarrow$  Hazard Pointers
- $\rightarrow$  Many others...



- $\rightarrow$  Reference Counting
- $\rightarrow$  Epoch-based Reclamation
- $\rightarrow$  Hazard Pointers
- $\rightarrow$  Many others...



- $\rightarrow$  Reference Counting
- $\rightarrow$  Epoch-based Reclamation
- $\rightarrow$  Hazard Pointers
- $\rightarrow$  Many others...



- $\rightarrow$  Reference Counting
- $\rightarrow$  Epoch-based Reclamation
- $\rightarrow$  Hazard Pointers
- $\rightarrow$  Many others...



# **REFERENCE COUNTING**

Maintain a counter for each node to keep track of the number of threads that are accessing it.

- $\rightarrow$  Increment the counter before accessing.
- $\rightarrow$  Decrement it when finished.
- $\rightarrow$  A node is only safe to delete when the count is zero.

This has bad performance for multi-core CPUs → Incrementing/decrementing counters causes a lot of cache coherence traffic.

## OBSERVATION

We don't actually care about the actual value of the reference counter. We only need to know when it reaches zero.

We don't have to perform garbage collection immediately when the counter reaches zero.

# EPOCH GARBAGE COLLECTION

Maintain a global epoch counter that is periodically updated (e.g., every 10 ms).

 $\rightarrow$  Keep track of what threads enter the index during an epoch and when they leave.

Mark the current epoch of a node when it is marked for deletion.

 $\rightarrow$  The node can be reclaimed once all threads have left that epoch (and all preceding epochs).

Also known as *Read-Copy-Update* (RCU) in Linux.

# NON-UNIQUE INDEXES

#### Approach #1: Duplicate Keys

 $\rightarrow$  Use the same node layout but store duplicate keys multiple times.

#### Approach #2: Value Lists

 $\rightarrow$  Store each key only once and maintain a linked list of unique values.



### NON-UNIQUE: DUPLICATE KEYS



### NON-UNIQUE: DUPLICATE KEYS



### NON-UNIQUE: DUPLICATE KEYS



#### NON-UNIQUE: VALUE LISTS


# NON-UNIQUE: VALUE LISTS



# VARIABLE LENGTH KEYS

#### Approach #1: Pointers

 $\rightarrow$  Store the keys as pointers to the tuple's attribute.

#### Approach #2: Variable Length Nodes

- $\rightarrow$  The size of each node in the index can vary.
- $\rightarrow$  Requires careful memory management.

#### Approach #3: Padding

 $\rightarrow$  Always pad the key to be max length of the key type.

#### Approach #4: Key Map / Indirection

 $\rightarrow$  Embed an array of pointers that map to the key + value list within the node.











Store a minimum prefix that is needed to correctly route probes into the index.



Store a minimum prefix that is needed to correctly route probes into the index.



Store a minimum prefix that is needed to correctly route probes into the index.



Store a minimum prefix that is needed to correctly route probes into the index.



Store a minimum prefix that is needed to correctly route probes into the index.



# ADAPATIVE RADIX TREE (ART)

Uses digital representation of keys to examine prefixes 1-by-1 instead of comparing entire key.

Radix trees properties:

- $\rightarrow$  The height of the tree depends on the length of keys.
- (unlike B+tree where height depends on the number of keys)
- $\rightarrow$  Does not require rebalancing
- $\rightarrow$  The path to a leaf node represents the key of the leaf
- $\rightarrow$  Keys are stored implicitly and can be reconstructed from paths.

 $\rightarrow$  Structure does not depend on order of key insertion ADAPTIVE RADIX TREE: ARTFUL KING FOR MAIN-MEMORY DATABASES



'CDF 2013

#### Trie (Re`trie'val - 1959)



#### Keys: HELLO, HAT, HAVE

#### Trie (Re`trie'val - 1959)





#### Trie (Re`trie'val - 1959)





#### Trie (Re`trie'val - 1959)





#### Trie (Re`trie'val - 1959)



#### Keys: HELLO, HAT, HAVE





# ART: ADAPTIVELY SIZED NODES

The index supports four different internal node types with different capacities.

Pack in multiple digits into a single node to improve cache locality.





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



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











Not all attribute types can be decomposed into binary comparable digits for a radix tree.

- → **Unsigned Integers:** Byte order must be flipped to big endian representation for little endian machines (x86).
- $\rightarrow$  **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.

# Int Key: 168496141

# Int Key: 168496141 Hex Key: 0A 0B 0C 0D

# Int Key: 168496141 Hex Key: 0A 0B 0C 0D









# Int Key: 168496141 Hex Key: 0A 0B 0C 0D
#### **ART: BINARY COMPARABLE KEYS**





Lookup: 658205 Hex: 0A 0B 1D

## BINARY COMPARABLE KEYS

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



# CONCURRENT ART INDEX

111

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

 $\rightarrow$  The authors argue that it would be a significant amount of work to make it latch-free.

Approach #1: Optimistic Lock Coupling Approach #2: Read-Optimized Write Exclusion

Optimistic crabbing scheme where writers are not blocked on readers.

- $\rightarrow$  Writers increment counter when they acquire latch.
- $\rightarrow$  Readers can proceed if a node's latch is available.
- $\rightarrow$  It then checks whether the latch's counter has changed from when it checked the latch.











































# **READ-OPTIMIZED WRITE EXCLUSION**

Each node includes an exclusive lock that blocks only other writers and not readers.

- $\rightarrow$  Readers proceed without checking versions or locks.
- $\rightarrow$  Every writer must ensure that reads are always consistent.

Requires fundamental changes to how threads make modifications to the data structure.

### **IN-MEMORY INDEXES**

Processor: 1 socket, 10 cores w/ 2×HT Workload: 50m Random Integer Keys (64-bit)

■ Open Bw-Tree ■ B+Tree ■ Skip List ■ Masstree ■ ART



## PARTING THOUGHTS

Andy was wrong about the Bw-Tree and latchfree indexes.

#### NEXT CLASS

Query Compilation

Reminder: Mid-term Exam on Feb 26<sup>th</sup>.

Reminder: Project Proposals due on Feb 26<sup>th</sup>.