Database Applications (15-415)

DBMS Internals- Part III Lecture 11, February 19, 2014

Mohammad Hammoud



Today...

Last Session:

- DBMS Internals- Part II
 - Brief summaries of disks, disk space management, and buffer management
 - Files and Access Methods (file organizations and motivation for indexing)
- Today's Session:
 - DBMS Internals- Part III
 - Tree-based indexes: ISAM, B and B+ (*if time allows*) trees

Announcements:

- PS3 is now posted. It is due on March 02, 2014
- The midterm exam is on Wednesday Feb 26 (all material are included)

Carnegie Mellon University Qatar



Outline





Motivation

Consider a file of student records *sorted* by GPA

Page 1	Page 2	Page 3	Page N	Data File
]]

How can we answer a range selection (E.g., "Find all students with a GPA higher than 3.0")?

- What about doing a *binary search* followed by a *scan*?
 - Yes, but...
- What if the file becomes "very" large?
 - Cost is proportional to the number of pages fetched
 - Hence, may become very slow!

Motivation

What about creating an *index file* (with one entry per page) and do binary search there?



But, what if the index file becomes also "very" large?

Motivation

Repeat recursively!



Each tree page is a disk block and all data records reside (*if chosen to be part of the index*) in ONLY leaf pages

How else data records can be stored?

Where to Store Data Records?

- In general, 3 alternatives for "data records" (each referred to as k*) can be pursued:
 - Alternative (1): K* is an actual data record with key k
 - Alternative (2): K* is a <k, rid> pair, where rid is the record id of a data record with search key k
 - Alternative (3): K* is a <k, rid-list> pair, where rid-list is a list of rids of data records with search key k

Where to Store Data Records?

In general, 3 alternatives for "data records" (each referred to as k*) can be pursued:

A (1): Leaf pages contain the actual data (i.e., the data records)

A (2): Leaf pages contain the <key, rid> pairs and actual data records are stored in a separate file

A (3): Leaf pages contain the <key, rid-list> pairs and actual data records are stored in a separate file

The choice among these alternatives is orthogonal to the *indexing technique*.

Outline



جامعة کارنیدی میلود فی قطر Carnegie Mellon University Qatar

ISAM Trees: Page Overflows

Now, what if there are a lot of insertions?



This structure is referred to as *Indexed Sequential Access Method* (ISAM)

ISAM File Creation

- How to create an ISAM file?
 - All leaf pages are allocated *sequentially* and *sorted* on the search key value
 - If Alternative (2) or (3) is used, the data records are created and *sorted* before allocating leaf pages
 - The non-leaf pages are subsequently allocated

An Example of ISAM Trees



2 Entries Per Page.

ISAM: Searching for Entries

- Search begins at root, and key comparisons direct it to a leaf
- Search for 27*





















ISAM: Some Issues

- Once an ISAM file is created, insertions and deletions affect only the contents of leaf pages (i.e., *ISAM is a <u>static</u> structure*!)
- Since index-level pages are *never* modified, there is no need to lock them during insertions/deletions
 - Critical for concurrency!
- Long overflow chains can develop easily
 - The tree can be initially set so that ~20% of each page is free
- If the data distribution and size are relatively static, ISAM might be a good choice to pursue!

Outline



جامعۃ کارنے جی سیلوں فی قطر Carnegie Mellon University Qatar

Dynamic Trees

- ISAM indices are static
 - Long overflow chains can develop as the file grows, leading to poor performance
- This calls for more flexible, *dynamic* indices that adjust gracefully to insertions and deletions
 - No need to allocate the leaf pages sequentially as in ISAM
- Among the most successful dynamic index schemes are B and B+ trees

B and B+ Trees

- B and B+ trees are designed to work on disks
 - A B/B+ tree node is usually as large as a whole disk page
 - B/B+ trees copy selected pages from disk into main memory as needed
 - Only a constant number of pages exit in memory at any time; hence, the size of the memory does not limit the size of a B/B+ tree that can be handled
 - O(log (N)) for any operation, assuming N-key B/B+ tree!

Outline



جامعۃ کارنے جے میلوں فی قطر Carnegie Mellon University Qatar

B Tree Properties

- Each node in a B tree of order *d* (this is a measure of the capacity of a tree):
 - Has at most 2d keys
 - Has at least *d* keys (except the root, which may have just 1 key)
 - All leaves are on the same level
 - Has exactly *n-1* keys if the number of pointers is *n*

Points to a sub-tree in which all keys are less than k₁



B Tree Properties

- Each node in a B tree of order *d* (this is a measure of the capacity of a tree):
 - Has at most 2d keys
 - A variant of a B tree, known as B* tree, requires each internal node to be at least 2/3 full, rather than half full, as a B tree requires.
 - All leaves are on the same level
 - Has exactly *n-1* keys if the number of pointers is *n*

Points to a sub-tree in which all keys are less than k₁



A B Tree Example

Below is a B Tree example with order *d* = 1



B trees are *balanced* search trees (they generalize binary trees)










B Tree: Queries

What about range queries (E.g., 5<salary<8)</p>



B Tree: Queries

What about range queries (E.g., 5<salary<8)</p>



B Tree: Insertions

- Search begins at root, and key comparisons direct it to a leaf (as in ISAM)
- Insert the given entry in the located leaf
- If an <u>overflow</u> occurs, <u>split</u> the node and <u>push up</u> the <u>middle key</u> (recursively)

How do we define an overflow in a B tree?

Easy case: insert '8'



Easy case: insert '8'



Hard case: insert '2'



Hard case: insert '2'



Push up the middle key and split

Hard case: insert '2'



Hard case: insert '2'



Automatic, incremental re-organization (contrast with ISAM!)

Pseudo-code: B Tree Insertions

```
INSERTION OF KEY 'K'
find the correct leaf node 'L';
if ('L' overflows ){
    split 'L', and push middle key to parent node 'P';
    if ('P' overflows){
        repeat the split recursively; }
    else{
        add the key 'K' in node 'L';
        /* maintaining the key order in 'L' */ }
```

- Delete begins at root, and key comparisons direct it to a leaf (as in ISAM)
- Delete entry (*if found*) in the located leaf
- If an underflow occurs, merge nodes

How do we define an underflow in a B tree?

- Four cases:
 - Case1: delete a key at a leaf no underflow
 - Case2: delete non-leaf key no underflow
 - Case3: delete leaf-key; underflow, and 'rich sibling'
 - Case4: delete leaf-key; underflow, and 'poor sibling'

- Four cases:
 - Case1: delete a key at a leaf no underflow
 - Case2: delete non-leaf key no underflow
 - Case3: delete leaf-key; underflow, and 'rich sibling'
 - Case4: delete leaf-key; underflow, and 'poor sibling'

Easy case: delete '3'



Easy case: delete '3'



- Four cases:
 - Case1: delete a key at a leaf no underflow

Case2: delete non-leaf key – no underflow

- Case3: delete leaf-key; underflow, and 'rich sibling'
- Case4: delete leaf-key; underflow, and 'poor sibling'

Delete '6'



Delete '6'



How to promote?

Pick the largest key from the left sub-tree (or the smallest from the right sub-tree).

Delete '6'



Delete '6'



- Four cases:
 - Case1: delete a key at a leaf no underflow
 - Case2: delete non-leaf key no underflow
 - Case3: delete leaf-key; underflow, and 'rich sibling'
 - Case4: delete leaf-key; underflow, and 'poor sibling'





Delete '7'



'Borrow' = can happen ONLY through the parent









- Four cases:
 - Case1: delete a key at a leaf no underflow
 - Case2: delete non-leaf key no underflow
 - Case3: delete leaf-key; underflow, and 'rich sibling'
 - Case4: delete leaf-key; underflow, and 'poor sibling'

Delete '13'



Delete '13'



Delete '13'



'Poor' = can host a key, without overflowing

Delete '13'



Delete '13'



Merge, by pulling a key from the **parent** (*the opposite of insertions*!)

Merge &

Pull from

Parent

Delete '13'



Delete '13'



But, what if the parent underflows? *Repeat recursively!*
Pseudo-code: B Tree Deletions

```
DELETION OF KEY 'K'
 locate key 'K', in node 'N'
 if( 'N' is a non-leaf node) {
   delete 'K' from 'N';
   find the immediately largest key 'K1';
     /* which is guaranteed to be on a leaf node 'L' */
   copy 'K1' in the old position of 'K';
   invoke this DELETION routine on 'K1' from the leaf node
 else {
/* 'N' is a leaf node */
  (next slide..)
```

Pseudo-code: B Tree Deletions

/* 'N' is a leaf node */ if('N' underflows){ let 'N1' be the sibling of 'N'; if('N1' is "rich"){ /* ie., N1 can lend us a key */ borrow a key from 'N1' THROUGH the parent node; }else{ /* N1 is 1 key away from underflowing */ MERGE: pull the key from the parent 'P', and merge it with the keys of 'N' and 'N1' into a new node; if ('P' underflows) { repeat recursively }

Outline



جا مہۃ کارنیدی ہیلوں فی قطر Carnegie Mellon University Qatar

Clustered vs. Un-clustered Indexes

- Indexes can be either clustered or un-clustered
- Clustered Indexes:
 - When the ordering of data records is the same as (or close to) the ordering of entries in some index
- Un-clustered Indexes:
 - When the ordering of data records differs from the ordering of entries in some index

Clustered vs. Un-clustered Indexes

- Is an index that uses Alternative (1) clustered or un-clustered?
 - Clustered
- Is an index that uses Alternatives (2) or (3) clustered or un-clustered?
 - Clustered if data records are sorted on the search key field
- In practice:
 - A clustered index is an index that uses Alternative (1)
 - Indexes that use Alternatives (2) or (3) are un-clustered

Outline



جا مدة کارنیدی میلود فی قطر Carnegie Mellon University Qatar

For clustering indexes, data records are scattered



<u>IDEA 1:</u> *replicate* keys from non-leaf nodes, to make sure every key appears at only one leaf node!

How can we facilitate sequential operations?



IDEA 2: String all leaf nodes together!

Towards B+ trees



Is this enough?

B+ trees



Vital for clustering indexes!

B+ Tree: Searching for Entries

- Search begins at root, and key comparisons direct it to a leaf (as in ISAM)
- Example 1: Search for entry 5*



B+ Tree: Searching for Entries

- Search begins at root, and key comparisons direct it to a leaf (as in ISAM)
- Example 2: Search for entry 15*



B+ Trees: Inserting Entries

- Find correct leaf L
- Put data entry onto L
 - If L has enough space, done!
 - Else, must <u>split</u> L into L and a new node L₂
 - Redistribute entries evenly, <u>copying up</u> the middle key
- Parent node may overflow
 - <u>Push up</u> middle key (splits "grow" trees; a root split increases the height of the tree)

Insert entry 8*



Leaf is *full*; hence, split!









Insert entry 8*



Splitting the root lead to an increase of height by 1!

What about re-distributing entries instead of splitting nodes?





Insert entry 8*



"Copy up" the new low key value!

But, when to redistribute and when to split?

Splitting vs. Redistributing

Leaf Nodes

- Previous and next-neighbor pointers must be updated upon insertions (*if splitting is to be pursued*)
- Hence, checking whether redistribution is possible does not increase I/O
- Therefore, if a sibling can spare an entry, re-distribute

Non-Leaf Nodes

- Checking whether redistribution is possible usually increases I/O
- Splitting non-leaf nodes typically pays off!

B+ Insertions: Keep in Mind

- Every data entry must appear in a leaf node; hence, "copy up" the middle key upon splitting
- When splitting index entries, simply "push up" the middle key
- Apply splitting and/or redistribution on leaf nodes
- Apply only splitting on non-leaf nodes

B+ Trees: Deleting Entries

- Start at root, find leaf L where entry belongs
- Remove the entry
 - If L is at least half-full, done!
 - If L underflows
 - Try to re-distribute (i.e., borrow from a "rich sibling" and "copy up" its *lowest key*)
 - If re-distribution fails, <u>merge</u> L and a "poor sibling"
 - Update parent
 - And possibly merge, recursively

Delete 19*



Removing 19* does not entail an underflow

Delete 19*



FINAL TREE!

Delete 20*



Deleting 20* entails an underflow; hence, check a sibling for redistribution

Delete 20*



The sibling is 'rich' (i.e., can lend an entry); hence, remove 20* and redistribute!

Delete 20*



"Copy up" 27*, the lowest value in the leaf from which we borrowed 24*

Delete 20*



"Copy up" 27*, the lowest value in the leaf from which we borrowed 24*

Delete 20*



FINAL TREE!

Delete 24*



The affected leaf will contain only 1 entry and the sibling cannot lend any entry (i.e., redistribution is not applicable); hence, <u>merge</u>!



Delete 24*



Almost there...




Delete 24*



Delete 24*



Delete 24*



Delete 24*



FINAL TREE!

Delete 24*



Assume (instead) the above tree *during* deleting 24*

Now we can <u>re-distribute</u> (*instead of merging*) keys!

Delete 24*



DONE! It suffices to re-distribute only 20; 17 was redistributed for illustration.

