

CompSci 516 Database Systems

Lecture 9 and 10
Storage
and
Index

Instructor: Sudeepa Roy

Duke CS, Fall 2019 CompSci 516: Database Systems 1

Announcements

- HW1 Deadlines!
 - Today: Q4
 - Q5: next to next Tuesday 10/01

Duke CS, Fall 2019 CompSci 516: Database Systems 2

Storage

- How are pages stored in a file?
 - Heap file (no particular order of records)
 - Sorted file (records sorted on any given field)
- How are records stored in a page?
 - Fixed length records
 - Variable length records
- How are fields stored in a record?
 - Fixed length fields/records
 - Variable length fields/records

The following slides give you the basic ideas, exact implementation may vary

Duke CS, Fall 2019 CompSci 516: Database Systems 3

Heap File Implemented as a List

- The header page id and Heap file name must be stored someplace
- Each page contains 2 `pointers' plus data
- But to insert a new record, we may need to scan several pages on the free list to find one with sufficient space

Duke CS, Fall 2019 CompSci 516: Database Systems 4

Heap File Using a Page Directory

- The entry for a page can include the number of free bytes on the page.
- The directory is a collection of pages
 - linked list implementation of directory is just one alternative
 - Much smaller than linked list of all heap file pages!

Duke CS, Fall 2019 CompSci 516: Database Systems 5

Storage

- How are pages stored in a file?
- How are records stored in a page?
 - Fixed length records
 - Variable length records
- How are fields stored in a record?
 - Fixed length fields/records
 - Variable length fields/records

Duke CS, Fall 2019 CompSci 516: Database Systems 6

How do we arrange a collection of records on a page?

- Each page contains several **slots**
 - one for each record
- Record is identified by **record id or rid = <page-id, slot-number>**
- Fixed-Length Records
- Variable-Length Records
- For both, there are options for
 - Record formats** (how to organize the fields within a record)
 - Page formats** (how to organize the records within a page)

Duke CS, Fall 2019 CompSci 516: Database Systems 7

Page Formats: Fixed Length Records

- Record id = <page id, slot #>
- Packed: moving records for free space management changes rid; may not be acceptable or may be slow to reorganize
- Unpacked: use a bitmap – scan the bit array to find an empty slot
- Each page also may contain additional info like the id of the next page (not shown)

Duke CS, Fall 2019 CompSci 516: Database Systems 8

Page Formats: Variable Length Records

- Need to find a page with the right amount of space
 - Too small – cannot insert
 - Too large – waste of space
- if a record is deleted, need to move the records so that all free space is contiguous
 - need ability to move records within a page
 - Changes record id
- Can maintain a **directory of slots** (next slide)

Duke CS, Fall 2019 CompSci 516: Database Systems 9

Page Formats: Variable Length Records Directory of Slots

- Each slot contains <record-offset, record-length>
 - deletion = set record-offset to -1
- Record-id rid = <page, slot-in-directory> remains unchanged
 - Can move records on page without changing rid
 - so, attractive for fixed-length records too

Duke CS, Fall 2019 CompSci 516: Database Systems 10

Storage

- How are pages stored in a file?
- How are records stored in a page?
 - Fixed length records
 - Variable length records
- How are fields stored in a record?
 - Fixed length fields/records
 - Variable length fields/records

Duke CS, Fall 2019 CompSci 516: Database Systems 11

Record Formats: Fixed Length

- Each field has a fixed length
 - for all records
 - the number of fields is also fixed
 - fields can be stored consecutively
- Information about field types same for all records in a file
 - stored in **system catalogs**
- Finding i-th field does not require scan of record
 - given the address of the record, address of a field can be obtained easily

Duke CS, Fall 2019 CompSci 516: Database Systems 12

Record Formats: Variable Length

- Cannot use fixed-length slots for records
- Two alternative formats (note: # fields is fixed for relational data)

Fields Delimited by Special Symbols

1. use delimiters

Array of Field Offsets

2. use offsets at the start of each record

- Second offers direct access to i-th field, efficient storage of nulls (special don't know value); small directory overhead

Duke CS, Fall 2019
CompSci 516: Database Systems
13

Main takeaways: storage

- Disk is slow but large and persistent
- Main memory or buffer is fast but small and not persistent
- If a page is edited in memory, needs to be written back to disk
- Unit of cost = page I/O (read and write)
- A record (= tuple) is accessed by rid (record id): gives the address of the page and the slot

Duke CS, Fall 2019
CompSci 516: Database Systems
14

Indexes

Duke CS, Fall 2019
CompSci 516: Database Systems
15

Indexes

- An index on a file speeds up selections on the search key fields for the index
 - Any subset of the fields of a relation can be the search key for an index on the relation.
 - “Search key” is not the same as “key”!
- An index contains a collection of data entries, and supports efficient retrieval of all data entries k^* with a given key value k
 - Why multiple entries for a given k ?

Duke CS, Fall 2019
CompSci 516: Database Systems
16

Remember: Terminology

- Index search key (key): k
 - Used to search a record
- Data entry : k^*
 - Pointed to by k
 - Contains record id(s) or record itself
- Records or data
 - Actual tuples
 - Pointed to by record ids

INDEX does this

Duke CS, Fall 2019
CompSci 516: Database Systems
17

Alternatives for Data Entry k^* in Index k

Advantages/ Disadvantages?

- In a data entry k^* we can store:
 1. (Alternative 1) The actual data record with key value k , or
 2. (Alternative 2) $\langle k, rid \rangle$
 - rid = record of data record with search key value k , or
 3. (Alternative 3) $\langle k, rid\text{-list} \rangle$
 - list of record ids of data records with search key k
- Choice of alternative for data entries is orthogonal to the indexing technique used to locate data entries with a given key value k

Duke CS, Fall 2019
CompSci 516: Database Systems
18

Alternatives for Data Entries: **Alternative 1**

- In a data entry k^* we can store:
 - The actual data record with key value k
 - $\langle k, rid \rangle$
 - rid = record of data record with search key value k
 - $\langle k, rid-list \rangle$
 - list of record ids of data records with search key k

- Index structure is a file organization for data records
 - instead of a Heap file or sorted file
- At most one index can use Alternative 1
 - Otherwise, data records are duplicated, leading to redundant storage and potential inconsistency
- Problem with Alt-1: If data records are very large, #pages with data entries is high
 - Implies size of auxiliary information in the index is also large

Alternatives for Data Entries: **Alternative 2, 3**

- In a data entry k^* we can store:
 - The actual data record with key value k
 - $\langle k, rid \rangle$
 - rid = record of data record with search key value k
 - $\langle k, rid-list \rangle$
 - list of record ids of data records with search key k

- Data entries typically much smaller than data records
 - So, better than Alternative 1 with large data records
 - Especially if search keys are small.
- Alternative 3 more compact than Alternative 2
 - but leads to variable-size data entries even if search keys have fixed length.

Index Classification

- Primary vs. secondary
- Clustered vs. unclustered
- Tree-based vs. Hash-based

Primary vs. Secondary Index

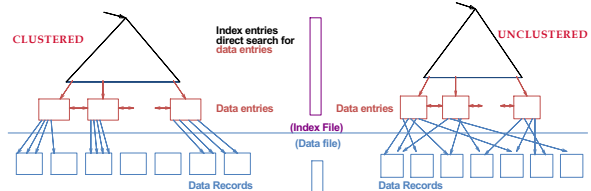
- If search key contains primary key, then called **primary index**, otherwise **secondary**
 - Unique index**: Search key contains a candidate key
- Duplicate data entries:
 - if they have the same value of search key field k
 - Primary/unique index never has a duplicate
 - Other secondary index can have duplicates

Clustered vs. Unclustered Index

- If order of data records in a file is the same as, or 'close to', order of data entries in an index, then clustered, otherwise unclustered
- A file can be clustered on at most one search key
- Cost of retrieving data records (range queries) through index varies greatly based on whether index is clustered or not

Clustered vs. Unclustered Index

- Suppose that Alternative (2) is used for data entries, and that the data records are stored in a Heap file
- To build clustered index, first sort the Heap file
 - with some free space on each page for future inserts
 - Overflow pages may be needed for inserts
 - Thus, data records are 'close to', but not identical to, sorted



Methods for indexing

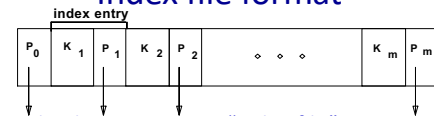
- Tree-based
- Hash-based

Tree-based Index and B⁺-Tree

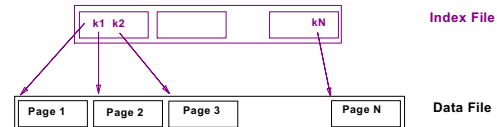
Range Searches

- *“Find all students with gpa > 3.0”*
 - If data is in sorted file, do “binary search” to find first such student, then scan to find others.
 - Cost of binary search can be quite high.

Index file format



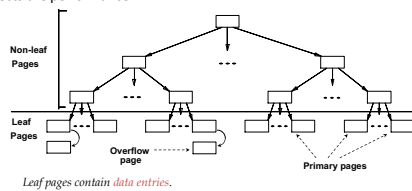
- Simple idea: Create an “index file”
 - <first-key-on-page, pointer-to-page>, sorted on keys



Can do binary search on (smaller) index file
but may still be expensive: apply this idea repeatedly

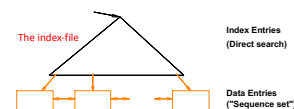
Indexed Sequential Access Method (ISAM)

- Leaf-pages contain data entry – also some overflow pages
- DBMS organizes layout of the index – a static structure
- If a number of inserts to the same leaf, a long overflow chain can be created
 - affects the performance



B+ Tree

- Most Widely Used Index: a dynamic structure
- Insert/delete at $\log_e N$ cost = height of the tree (cost = I/O)
 - $F = \text{fanout}$, $N = \text{no. of leaf pages}$
 - tree is maintained **height-balanced**
- Minimum 50% occupancy
 - Each node contains $d \leq m \leq 2d$ entries
 - Root contains $1 \leq m \leq 2d$ entries
 - The parameter d is called the **order** of the tree
- Supports **equality and range-searches** efficiently



B+ Tree Indexes

Non-leaf Pages

Leaf Pages (Sorted by search key)

- Leaf pages contain **data entries**, and are chained (prev & next)
- Non-leaf pages have **index entries**; only used to direct searches:

index entry

P₀ K₁ P₁ K₂ P₂ ... K_m P_m

Duke CS, Fall 2019 CompSci 516: Database Systems 31

Example B+ Tree

- Search begins at root, and key comparisons direct it to a leaf
- Search for 5*, 15*, all data entries $\geq 24^*$...

Based on the search for 15, we know it is not in the tree!*

Duke CS, Fall 2019 CompSci 516: Database Systems 32

Example B+ Tree

Note how data entries in leaf level are sorted

- Find
 - 28*?
 - 29*?
 - All > 15* and < 30*

Duke CS, Fall 2019 CompSci 516: Database Systems 33

B+ Trees in Practice

- Typical order: $d = 100$. Typical fill-factor: 67%
 - average fanout $F = 133$
- Typical capacities:
 - Height 4: $133^4 = 312,900,700$ records
 - Height 3: $133^3 = 2,352,637$ records
- Can often hold top levels in buffer pool:
 - Level 1 = 1 page = 8 Kbytes
 - Level 2 = 133 pages = 1 Mbyte
 - Level 3 = 17,689 pages = 133 Mbytes

Duke CS, Fall 2019 CompSci 516: Database Systems 34

Inserting a Data Entry into a B+ Tree

- Find correct leaf L
- Put data entry onto L
 - If L has enough space, **done**
 - Else, must **split** L
 - into L and a new node L2
 - Redistribute entries evenly, **copy up** middle key.
 - Insert index entry pointing to L2 into parent of L.
- This can happen recursively
 - To **split index node**, redistribute entries evenly, but **push up** middle key
 - Contrast with leaf splits
- Splits "grow" tree; **root split** increases height.
 - Tree growth: gets **wider** or **one level taller** at top.

See this slide later, First, see examples on the next few slides

Duke CS, Fall 2019 CompSci 516: Database Systems 35

Inserting 8* into Example B+ Tree

STEP-1

- Copy-up: 5 appears in leaf and the level above
- Observe how minimum occupancy is guaranteed

Entry to be inserted in parent node. (Note that 5 is copied up and continues to appear in the leaf.)

Duke CS, Fall 2019 CompSci 516: Database Systems 36

Inserting 8* into Example B+ Tree

Need to split parent STEP-2

- Note difference between copy-up and push-up
- What is the reason for this difference?
- All data entries must appear as leaves
 - (for easy range search)
- no such requirement for indexes
 - (so avoid redundancy)

Entry to be inserted in parent node. (Note that 17 is pushed up and only appears once in the index. Contrast this with a leaf split.)

Duke CS, Fall 2019 CompSci 516: Database Systems 38

Example B+ Tree After Inserting 8*

- Notice that root was split, leading to increase in height.
- In this example, we can avoid split by re-distributing entries (insert 8 to the 2nd leaf node from left and copy it up instead of 13)
 - however, this is usually not done in practice - since need to access 1-2 extra pages always (for two siblings), and average occupancy may remain unaffected as the file grows

Duke CS, Fall 2019 CompSci 516: Database Systems 38

Deleting a Data Entry from a B+ Tree

Each non-root node contains $d \leq m \leq 2d$ entries

- Start at root, find leaf L where entry belongs
- Remove the entry
 - If L is at least half-full, done!
 - If L has only $d-1$ entries,
 - Try to re-distribute, borrowing from sibling (adjacent node with same parent as L)
 - If re-distribution fails, merge L and sibling
- If merge occurred, must delete entry (pointing to L or sibling) from parent of L
- Merge could propagate to root, decreasing height

See this slide later, First, see examples on the next few slides

Duke CS, Fall 2019 CompSci 516: Database Systems 39

Example Tree: Delete 19*

Before deleting 19*

- We had inserted 8*
- Now delete 19*
- Easy

Duke CS, Fall 2019 CompSci 516: Database Systems 40

Example Tree: Delete 19*

After deleting 19*

Duke CS, Fall 2019 CompSci 516: Database Systems 41

Example Tree: Delete 20*

Before deleting 20*

Duke CS, Fall 2019 CompSci 516: Database Systems 42

Example Tree: Delete 20*

After deleting 20*
- step 1

- < 2 entries in leaf-node
- Redistribute

Duke CS, Fall 2019 CompSci 516: Database Systems 43

Example Tree: Delete 20*

After deleting 20*
- step 2

- Notice how middle key is copied up

Duke CS, Fall 2019 CompSci 516: Database Systems 44

Example Tree: ... And Then Delete 24*

Before deleting 24*

Duke CS, Fall 2019 CompSci 516: Database Systems 45

Example Tree: ... And Then Delete 24*

After deleting 24*
- Step 1

- Once again, imbalance at leaf
- Can we borrow from sibling(s)?
- No – d-1 and d entries (d = 2)
- Need to merge

Duke CS, Fall 2019 CompSci 516: Database Systems 46

Example Tree: ... And Then Delete 24*

After deleting 24*
- Step 2

Observe "toss" of old index entry 27

- Imbalance at parent
- Merge again
- But need to "pull down" root index entry

because, three index 5, 13, 30
but five pointers to leaves

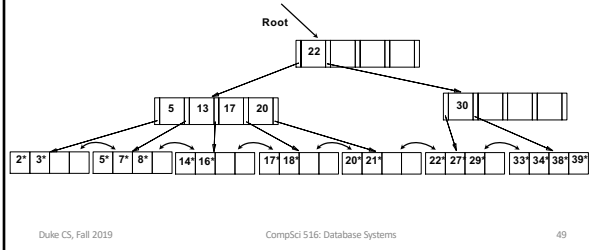
Duke CS, Fall 2019 CompSci 516: Database Systems 47

Final Example Tree

Duke CS, Fall 2019 CompSci 516: Database Systems 48

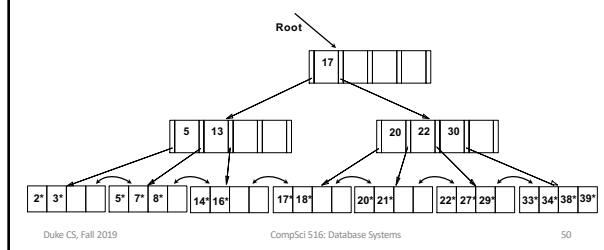
Example of Non-leaf Re-distribution

- An intermediate tree is shown
- In contrast to previous example, can re-distribute entry from left child of root to right child



After Re-distribution

- Intuitively, entries are re-distributed by 'pushing through' the splitting entry in the parent node.
 - It suffices to re-distribute index entry with key 20; we've re-distributed 17 as well for illustration.



Duplicates

- **First Option:**
 - The basic search algorithm assumes that all entries with the same key value resides on the same leaf page
 - If they do not fit, use overflow pages (like ISAM)
- **Second Option:**
 - Several leaf pages can contain entries with a given key value
 - Search for the left most entry with a key value, and follow the leaf-sequence pointers
 - Need modification in the search algorithm
- if $k^* = \langle k, rid \rangle$, several entries have to be searched
 - Or include rid in k – becomes unique index, no duplicate
 - If $k^* = \langle k, rid-list \rangle$, same solution, but if the list is long, again a single entry can span multiple pages

A Note on 'Order'

- **Order (d)**
 - denotes minimum occupancy
- replaced by physical space criterion in practice ('at least half-full')
 - Index pages can typically hold many more entries than leaf pages
 - Variable sized records and search keys mean different nodes will contain different numbers of entries.
 - Even with fixed length fields, multiple records with the same search key value (duplicates) can lead to variable-sized data entries (if we use Alternative (3))

Summary

- Tree-structured indexes are ideal for range-searches, also good for equality searches
- ISAM is a static structure
 - Only leaf pages modified; overflow pages needed
 - Overflow chains can degrade performance unless size of data set and data distribution stay constant
- B+ tree is a dynamic structure
 - Inserts/deletes leave tree height-balanced; $\log_p N$ cost
 - High fanout (F) means depth rarely more than 3 or 4
 - Almost always better than maintaining a sorted file
 - Most widely used index in database management systems because of its versatility.
 - One of the most optimized components of a DBMS
- Next: Hash-based index

Hash-based Index

Hash-Based Indexes

- Records are grouped into buckets
 - Bucket = **primary page** plus zero or more **overflow pages**
- Hashing function **h**:
 - $h(r)$ = bucket in which (data entry for) record r belongs
 - h** looks at the **search key** fields of r
 - No need for "index entries" in this scheme

Duke CS, Fall 2019 CompSci 516: Database Systems 55

Example: Hash-based index

Index organized file hashed on AGE, with Auxiliary index on SAL

Employee File hashed on AGE

Alternative 2

Duke CS, Fall 2019 CompSci 516: Database Systems 56

Introduction

- Hash-based indexes are best for **equality selections**
 - Find all records with name = "Joe"
 - Cannot support range searches
 - But useful in implementing relational operators like join (later)
- Static and dynamic hashing techniques exist
 - trade-offs similar to ISAM vs. B+ trees

Duke CS, Fall 2019 CompSci 516: Database Systems 57

Static Hashing

- Pages containing data = a collection of **buckets**
 - each bucket has one primary page, also possibly overflow pages
 - buckets contain **data entries k^***

Primary bucket pages Overflow pages

Duke CS, Fall 2019 CompSci 516: Database Systems 58

Static Hashing

- # primary pages fixed
 - allocated sequentially, never de-allocated, overflow pages if needed.
- $h(k) \bmod N$ = bucket to which data entry with key k belongs
 - N = # of buckets

Primary bucket pages Overflow pages

Duke CS, Fall 2019 CompSci 516: Database Systems 59

Static Hashing

- Hash function works on search key field of record r
 - Must distribute values over range $0 \dots N-1$
 - $h(\text{key}) = (a * \text{key} + b)$ usually works well
 - bucket = $h(\text{key}) \bmod N$
 - a and b are constants – chosen to tune h
- Advantage:
 - #buckets known – pages can be allocated sequentially
 - search needs 1 I/O (if no overflow page)
 - insert/delete needs 2 I/O (if no overflow page) (why 2?)
- Disadvantage:
 - Long overflow chains can develop if file grows and degrade performance (**data skew**)
 - Or waste of space if file shrinks
- Solutions:
 - keep some pages say 80% full initially
 - Periodically **rehash** if overflow pages (can be expensive)
 - or use **Dynamic Hashing**

Duke CS, Fall 2019 CompSci 516: Database Systems 60

Dynamic Hashing Techniques

- Extendible Hashing
- Linear Hashing

Duke CS, Fall 2019

CompSci 516: Database Systems

61

Extendible Hashing

- Consider static hashing
- Bucket (primary page) becomes full
 - Why not re-organize file by doubling # of buckets?
 - Reading and writing (double #pages) all pages is expensive
- Idea: Use directory of pointers to buckets
 - double # of buckets by doubling the directory, splitting just the bucket that overflowed
 - Directory much smaller than file, so doubling it is much cheaper
 - Only one page of data entries is split
 - No overflow page (new bucket, no new overflow page)
 - Trick lies in how hash function is adjusted

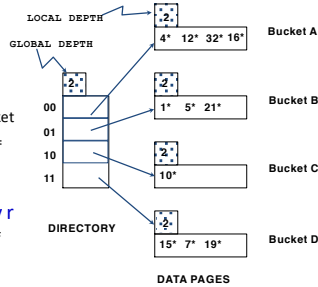
Duke CS, Fall 2019

CompSci 516: Database Systems

62

Example

- Directory is array of size 4
 - each element points to a bucket
 - #bits to represent = $\log 4 = 2 =$ global depth
- To find bucket for search key r
 - take last global depth # bits of $h(r)$
 - assume $h(r) = r$
 - If $h(r) = 5 =$ binary 101
 - it is in bucket pointed to by 01



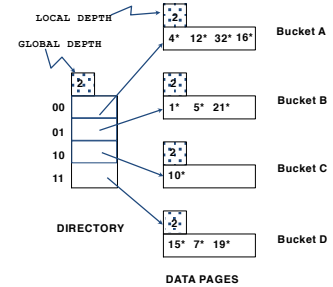
Duke CS, Spring 2016

CompSci 516: Data Intensive Computing Systems

13

Example

- Insert:
- If bucket is full, split it
 - allocate new page
 - re-distribute
- Suppose inserting 13*
- binary = 1101
 - bucket 01
 - Has space, insert



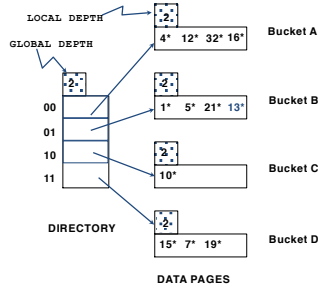
Duke CS, Spring 2016

CompSci 516: Data Intensive Computing Systems

14

Example

- Insert:
- If bucket is full, split it
 - allocate new page
 - re-distribute
- Suppose inserting 20*
- binary = 10100
 - bucket 00
 - Already full
 - To split, consider last three bits of 10100
 - Last two bits the same 00 – the data entry will belong to one of these buckets
 - Third bit to distinguish them



Duke CS, Spring 2016

CompSci 516: Data Intensive Computing Systems

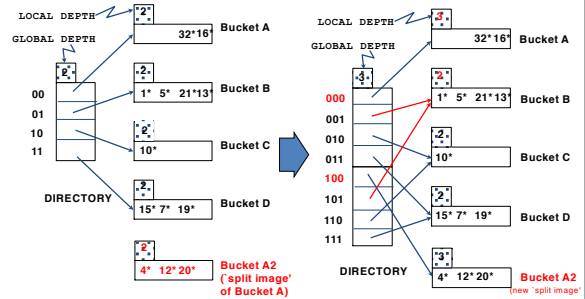
15

Example

Global depth: Max # of bits needed to tell which bucket an entry belongs to

Local depth: # of bits used to determine if an entry belongs to this bucket

- also denotes whether a directory doubling is needed while splitting
- no directory doubling needed when $9 = 1001$ is inserted (LD < GD)



Duke CS, Fall 2019

CompSci 516: Database Systems

66

When does bucket split cause directory doubling?

- Before insert, local depth of bucket = global depth
- Insert causes local depth to become > global depth
- directory is doubled by copying it over and 'fixing' pointer to split image page

Duke CS, Fall 2019 CompSci 516: Database Systems 67

Comments on Extendible Hashing

- If directory fits in memory, equality search answered with one disk access (to access the bucket); else two.
 - 100MB file, 100 bytes/rec, 4KB page size, contains 10^6 records (as data entries) and 25,000 directory elements; chances are high that directory will fit in memory.
 - Directory grows in spurts, and, if the distribution of hash values is skewed, directory can grow large
 - Multiple entries with same hash value cause problems
- Delete:
 - If removal of data entry makes bucket empty, can be merged with 'split image'
 - If each directory element points to same bucket as its split image, can halve directory.

Duke CS, Fall 2019 CompSci 516: Database Systems 68

Linear Hashing

- This is another dynamic hashing scheme
 - an alternative to Extendible Hashing
- LH handles the problem of long overflow chains
 - without using a directory
 - handles duplicates and collisions
 - very flexible w.r.t. timing of bucket splits

Duke CS, Fall 2019 CompSci 516: Database Systems 69

Linear Hashing: Basic Idea

- Use a family of hash functions h_0, h_1, h_2, \dots
 - $h_i(\text{key}) = h(\text{key}) \bmod(2^i N)$
 - $N =$ initial # buckets
 - h is some hash function (range is not 0 to $N-1$)
 - If $N = 2^{d_0}$, for some d_0 , h_i consists of applying h and looking at the last d_i bits, where $d_i = d_0 + i$
 - Note: $h_i(\text{key}) = h(\text{key}) \bmod(2^{d_0+i})$
 - h_{i+1} doubles the range of h_i
 - if h_i maps to M buckets, h_{i+1} maps to $2M$ buckets
 - similar to directory doubling
 - Suppose $N = 32, d_0 = 5$
 - $h_0 = h \bmod 32$ (last 5 bits)
 - $h_1 = h \bmod 64$ (last 6 bits)
 - $h_2 = h \bmod 128$ (last 7 bits) etc.

Duke CS, Fall 2019 CompSci 516: Database Systems 70

Linear Hashing: Rounds

- Directory avoided in LH by using overflow pages, and choosing bucket to split round-robin
- During round $Level$, only h_{Level} and $h_{Level+1}$ are in use
- The buckets from start to last are split sequentially
 - this doubles the no. of buckets
- Therefore, at any point in a round, we have
 - buckets that have been split
 - buckets that are yet to be split
 - buckets created by splits in this round

Duke CS, Fall 2019 CompSci 516: Database Systems 71

Overview of LH File

- In the middle of a round $Level$ – originally 0 to N_{Level}

The diagram illustrates the Linear Hashing File structure. It shows a vertical stack of buckets. At the top, buckets 0 to Next-1 are marked as 'Buckets split in this round:'. A red arrow points to a bucket labeled 'Bucket to be split Next'. Below this, a green bracket indicates the range of buckets that existed at the beginning of this round, labeled h_{Level} . A blue arrow points to a bucket labeled 'split image' bucket, with a note: 'if $h_{Level}(r)$ is in this range, must use $h_{Level+i}(r)$ to decide if entry is in 'split image' bucket.' Another blue arrow points to a bucket labeled 'split image' bucket, with a note: 'if $h_{Level}(r)$ is in this range, no need'. At the bottom, a blue arrow points to a bucket labeled 'split image' bucket, with a note: 'split image' buckets: created (through splitting of other buckets) in this round'. A legend at the bottom left explains:

- Buckets 0 to Next-1 have been split
- Next to N_{Level} yet to be split
- Round ends when all N_{Level} initial (for round $Level$) buckets are split

Duke CS, Fall 2019 CompSci 516: Database Systems 72

Overview of LH File

- In the middle of a round **Level** – originally 0 to N_{Level}

0
Next - 1
Bucket to be split
Next

Buckets split in this round:
if $h_{Level}(r)$ is in this range, must use $h_{Level+1}(r)$ to decide if entry is in 'split image' bucket.

Buckets that existed at the beginning of this round: this is the range of h_{Level}

if $h_{Level}(r)$ is in this range, no need

'split image' buckets: created (through splitting of other buckets) in this round

- Buckets 0 to Next-1 have been split
- Next to N_{Level} yet to be split
- Round ends when all N_{Level} initial (for round R) buckets are split

- Search: To find bucket for data entry r , find $h_{Level}(r)$:
- If $h_{Level}(r)$ in range 'Next to N_{Level} ', r belongs here.
- Else, r could belong to bucket $h_{Level}(r)$ or $h_{Level}(r)+N_x$
- Apply $h_{Level+1}(r)$ to find out

Duke CS, Fall 2019 CompSci 516: Database Systems 73

Linear Hashing: Insert

- Insert:** Find bucket by applying $h_{Level} / h_{Level+1}$:
 - If bucket to insert into is full:
 - Add overflow page and insert data entry
 - Split **Next** bucket and increment **Next**
- Note:** We are going to assume that a split is 'triggered' whenever an insert causes the creation of an overflow page, but in general, we could impose additional conditions for better space utilization ([RG], p.380)

Duke CS, Fall 2019 CompSci 516: Database Systems 74

Example of Linear Hashing

Level=0, $N_0 = 4 = 2^2$, $d=2$

h	h	PRIMARY PAGES
1	0	Next=0 32 44 36
000	00	
001	01	9 25 5 Data entry r with $h(r)=5$
010	10	14 18 10 30 Primary bucket page
011	11	31 35 7 11

- Insert $43^* = 101011$
- $h_0(43) = 11$
- Full
- Insert in an overflow page
- Need a split at Next (=0)
- Entries in 00 is distributed to 000 and 100

(This info is for illustration only!) (The actual contents of the linear hashed file)

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 26

Example of Linear Hashing

Level=0, $N_0 = 4 = 2^2$, $d=2$

h	h	PRIMARY PAGES	OVERFLOW PAGES
1	0	Next=0 32 44 36	
000	00		32
001	01	9 25 5 Data entry r with $h(r)=5$	Next=1 9 25 5
010	10	14 18 10 30 Primary bucket page	
011	11	31 35 7 11	43
100	00	44 36	

- Next is incremented after split
- Note the difference between overflow page of 11 and split image of 00 (000 and 100)

(This info is for illustration only!) (The actual contents of the linear hashed file)

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 27

Example of Linear Hashing

- Search for $18^* = 10010$
 - between Next (=1) and 4
 - this bucket has not been split
 - 18 should be here
- Search for $32^* = 10000$ or $44^* = 101100$
- Between 0 and Next-1
 - Need h_1
- Not all insertion triggers split
 - Insert $37^* = 100101$
 - Has space
- Splitting at Next?
 - No overflow bucket needed
 - Just copy at the image/original
- Next = $N_{Level}-1$ and a split?
 - Start a new round
 - Increment Level
 - Next reset to 0

Level=0, $N_0 = 4 = 2^2$, $d=2$

h	h	PRIMARY PAGES	OVERFLOW PAGES
1	0	Next=1 32	
000	00		
001	01	9 25 5	
010	10	14 18 10 30 Primary bucket page	
011	11	31 35 7 11	43
100	00	44 36	

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 28

Example of Linear Hashing

- Not all insertion triggers split
- Insert $37^* = 100101$
 - Has space

Level=0, $N_0 = 4 = 2^2$, $d=2$

h	h	PRIMARY PAGES	OVERFLOW PAGES
1	0	Next=1 32	
000	00		
001	01	9 25 5 37	
010	10	14 18 10 30 Primary bucket page	
011	11	31 35 7 11	43
100	00	44 36	

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 28

Example of Linear Hashing

- Splitting at Next?
 - No overflow bucket needed
 - Just copy at the image/original

insert $29^* = 11101$

Level=0, $N_0 = 4 = 2^0$, $d_0=2$

Level=0, $N_0 = 4 = 2^0$, $d_0=2$

Duke CS, Spring 2016 CompSci 516: Data Intensive Computing Systems 28

Example: End of a Round

insert $50^* = 110010$ Level=1, $N_1 = 8 = 2^1$, $d_1=3$

(after inserting 22^* , 66^* , 34^* - check yourself!)

Level=0, $N_0 = 4 = 2^0$, $d_0=2$

Level=1, $N_1 = 8 = 2^1$, $d_1=3$

Duke CS, Fall 2019 CompSci 516: Database Systems 80

LH vs. EH

- They are very similar
 - h_i to h_{i+1} is like doubling the directory
 - LH: avoid the explicit directory, clever choice of split
 - EH: always split – higher bucket occupancy
- Uniform distribution: LH has lower average cost
 - No directory level
- Skewed distribution
 - Many empty/nearly empty buckets in LH
 - EH may be better

Duke CS, Fall 2019 CompSci 516: Database Systems 81

System Catalogs

- For each index:
 - structure (e.g., B+ tree) and search key fields
- For each relation:
 - name, file name, file structure (e.g., Heap file)
 - attribute name and type, for each attribute
 - index name, for each index
 - integrity constraints
- For each view:
 - view name and definition
- Plus statistics, authorization, buffer pool size, etc.
- (described in [RG] 12.1)

Catalogs are themselves stored as relations!

Duke CS, Fall 2019 CompSci 516: Database Systems 82

Summary

- Hash-based indexes: best for equality searches, cannot support range searches.
- Static Hashing can lead to long overflow chains.
- Extendible Hashing avoids overflow pages by splitting a full bucket when a new data entry is to be added to it
 - Duplicates may still require overflow pages
 - Directory to keep track of buckets, doubles periodically
 - Can get large with skewed data; additional I/O if this does not fit in main memory

Duke CS, Fall 2019 CompSci 516: Database Systems 83

Summary

- Linear Hashing avoids directory by splitting buckets round-robin, and using overflow pages
 - Overflow pages not likely to be long
 - Duplicates handled easily
- For hash-based indexes, a skewed data distribution is one in which the hash values of data entries are not uniformly distributed
 - bad

Duke CS, Fall 2019 CompSci 516: Database Systems 84