

Relational Algebra

(π) Projection: retains only attributes in projection list
 ex: $S2 = \{name, phone\}$
 $\pi_{name}(S2)$

Selection (σ): select rows satisfying condition
 $\sigma_{age > 8}(S2)$

Cross product (\times): each row of $S1$ paired w/ each row of $P1$

Union (\cup): union-compatible... same # of fields, fields of same type

Set difference ($-$)

sid	sname	rating
22	A	1
23	B	3
24	C	3

$S1 - S2$
58, C, 3

$S2 - S1$
28, D, 5
8, Y, D

Compound operators: (intersect) \cap and Δ join

natural join example:

sid	bid	age
1	3	20
2	4	21
3	2	22
4	1	23

conditional join (\bowtie) theta join

$R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$

Examples:
 find names of sailors who've reserved red boat
 $\pi_{sname}(\sigma_{color='red'}(Boats) \bowtie Reserves \bowtie Sailors)$

find names of artists who have albums of genre of either "pop" or "rock"
 $\pi_{artist_name}(\sigma_{(genre='pop' \vee genre='rock')}(Albums))$

find id of artists who have albums of genre "pop" or have spent over 10 weeks in top 40.
 $\pi_{artist_id}(\sigma_{(Album.genre='pop' \vee (Album.genre='pop' \wedge Songs.wk_top_40 > 10))}(Artists \bowtie Album))$

find names of artists who do not have albums
 $\pi_{artists.name}(\sigma_{\neg(\exists Album.artist_id = Artists.artist_id)}(Artists))$

Info: Songs (song-id, song-name, album-id, top40)
 Artists (artist-id, artist-name, first-yr)
 Albums (album-id, album-name, artist-id, yr-released, genre)

Storing Data: Disks & Files

DBMS stores info. on disks
 lowest level manages space.

Unordered heap files: linked list or pg. directory
 allows us to retrieve records: scan sequentially

record formats: fixed length

variable length (2 options)

fields delimited by special symbols

array of field offsets

base address (b)

B+ L1+L2

OR

array of field offsets

fields delimited by special symbols

array of field offsets

Joins

$[R]$: # pages to store R
 P_R : # records/pg R
 $|R|$: # records in R

$([R] \times P_R) = |R|$

Simple Nested loops Join
 $(P_R \times [R]) \times [S] + [P]$

for each tuple, s can s

Page Oriented Nested Loop Join
 $([P] \times [S]) + R$ (doesn't exploit buffers)

for each page P, scan S

Chunk Nested Loop Join
 Say ... $B = 100 + 2$ memory buffers

Join cost: $[outer + [outer_chunks] * [inner]] < k, S \{items\}$

Hash Join
 2 pass hash join: $3([R] + [S])$ + output

partitioning phase:
 $2([R] + [S])$ + output

matching phase
 read both rel., write output

to read/write both relations

read both rel., write output

$([R] + [S]) + [output]$

Compare
 nested loops: nm
 Equi joins (hash doesn't work)

hash better if relation fits into mem, equivalent

min 50%
 node occupancy except root

each node: m entries where $d \leq m \leq 2d$

entries, where "d" is order of tree

B+ Discussion

Info: 2 million users in database
 each user entry 2kb, mainly performing range queries on user's age
 page size: 16kb

1) Clustered B+ on age field.
 fanout = 200, height = 3
 avg 50,000 users/query

How many I/Os per query? $3 + \frac{50000 \times 2}{16}$

2) unclustered. worst case?
 Index entry 3x smaller than entry.
 3 I/Os to descend to leaf pgs.
 read data entries:

$(50,000 \times \frac{2}{3}) = 2084$ I/Os.

50,000 to read unordered data pages.

Data vs Index Pg. Split

data page split: entry inserted in parent node appears in leaf (copy up)

Index pg. split: push up

30 still appears in leaf!!
 copied up b/c data node

final tree...

25 pushed up

But for Manager (Clib)

A	B	C	D	A	F	A	D	G	D	G	E	D	F
A			(1)	1	(1)	1	1	1	1	1	0	1	(1)
B			(1)	1	0	1	1	1	1	0	0	0	(1)
C			(1)	1	0	0	0	1	1	0	0	0	(1)
D			(1)	1	0	0	0	1	1	0	0	0	(1)

all equal here due to ref bit. A isn't min' b/c it already had a ref bit of 1!

b/c the J don't "add" if ref 1

Files & Access Management

Indexes: Disk based data structure for fast lookup by value.

contains a collection of data entries

Alternatives for Data Entry k in index

1) Actual data record (w/ key value k)

2) $< k, rid \text{ of matching data record} >$

3) $< k, list \text{ of rids of matching records} >$

All 1 \Rightarrow clustered

B+ trees

Create index on files to speed up selection on search by fields. Can have diff indexes on diff keys.

insert/delete at $\log_2 N$ cost

$f = \# \text{ entries/pg (fanout)}$
 $N = \# \text{ leaf pages}$

height balanced!, length of path from root to leaf node.

min 50%
 node occupancy except root

each node: m entries where $d \leq m \leq 2d$

entries, where "d" is order of tree

B+ Discussion

Info: 2 million users in database
 each user entry 2kb, mainly performing range queries on user's age
 page size: 16kb

1) Clustered B+ on age field.
 fanout = 200, height = 3
 avg 50,000 users/query

How many I/Os per query? $3 + \frac{50000 \times 2}{16}$

2) unclustered. worst case?
 Index entry 3x smaller than entry.
 3 I/Os to descend to leaf pgs.
 read data entries:

$(50,000 \times \frac{2}{3}) = 2084$ I/Os.

50,000 to read unordered data pages.

Data vs Index Pg. Split

data page split: entry inserted in parent node appears in leaf (copy up)

Index pg. split: push up

30 still appears in leaf!!
 copied up b/c data node

final tree...

When to change col's ref bit to 0? See ... when you have a col of 1's and LRM (list changed replaced)

Query Opt. & Execution

Relational Operators

Files & Access Methods

Buffer Management

Disk Space Management

DB

Out of Operations

assume single record insert/delete

exactly 1 match for equality sel.

heap files: insert always appends to end

sorted files: Also compacted after deletions & selection on search key.

B: no. data pages

R: no. records/page

D: (avg) time to read/write disk page

Heap file

Sorted file

Clustered File

Scan all records

equality search

range search

insert

delete

BD

BD

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

Query Opt. & Execution

Relational Operators

Files & Access Methods

Buffer Management

Disk Space Management

DB

Out of Operations

assume single record insert/delete

exactly 1 match for equality sel.

heap files: insert always appends to end

sorted files: Also compacted after deletions & selection on search key.

B: no. data pages

R: no. records/page

D: (avg) time to read/write disk page

Heap file

Sorted file

Clustered File

Scan all records

equality search

range search

insert

delete

BD

BD

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

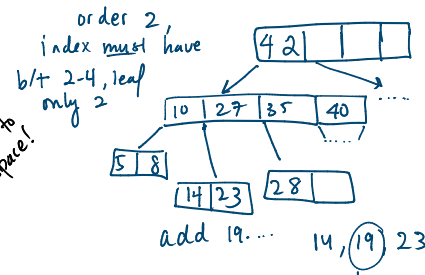
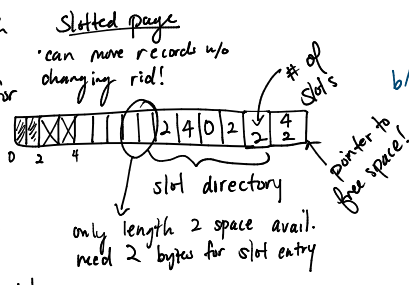
$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

$(\log_2 B)D$

Page formats: fixed length
 Rid = <page id, Slot #>
 • problematic: moving records for free space changes rid!
 Q: 80 bytes free space (slotted pg). Costs 4 bytes for directory entry.
 How many 1 byte records?
 $\frac{80}{1+4} = 16$ records!



Example: file org/indexes
 Consider 500 pgs (B), 6k tuples, query sid > 450 (sid unique, range 0-6k)
 ① i/o's on heap file? B = 500
 ② i/o if stored in sorted file by sid?
 $\log_2(500) + \frac{1}{4} \cdot 500$
 \uparrow
 $\frac{4500}{6000} = .25$

Hashing
 150 buffer pages, 4 tuples/page
 40,000 animals

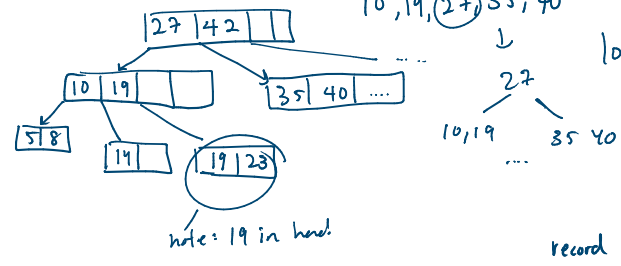
hash aggregation w/ hybrid hashing
 reserve k pages for in-mem table.

max k to hash all animals in 2 passes?
 (using $B-k$ buffers?)

1 pass ... $B-k$ pages for output partitions, each...
 no bigger than B . $N \leq (B-k) * B$

$B = 150$ and $N = 10,000$ ($\frac{40000}{4}$)
 $B - k = 66.67$
 $k = 150 - 66.67 = 83$

Sequential flooding: LRU + repeated seq. scans
 # buffer frames < pgs in file
 each pg request \Rightarrow i/o!



Buffer Manager:

record evictions, 3 page buffer

1	2	3	4	5	6	7	8	9	10
D	B	S	A	B	C	D	S	B	S

Mrn: D
 B (8) A (B) C (S) (S) (S) evicted (S)

clock: D(1) A(1) 1 1 D(1) 0 0 0
 B(1) 0 +1 1 1 S(1) 1 ✓
 S(1) 0 0 C(1) 1 D(1) B(1) 1
 1 2 3 4 5 6 7 8 9 10
 D B S A B C D S B S