

External Sorting:

Merge:

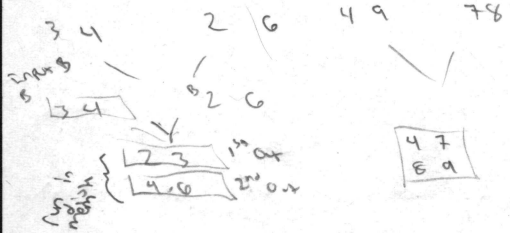
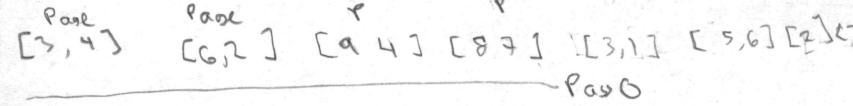
passes

2 way $\lceil \log_2 N \rceil + 1$

General $\lceil \log_{B-1} \lceil \frac{N}{B} \rceil \rceil + 1$

Total cost:

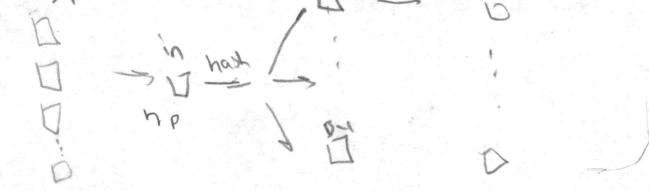
$2N \cdot [\text{# passes}]$



2 passes size $N = B(B-1)$
 space = $B \approx \sqrt{N}$
 Sort N pages of data
 in about \sqrt{N} space
 $\frac{N}{B}$ comes from
 Pass 0, creates $\frac{N}{B}$ runs
 of size B .

External Hashing: Good for removing duplicates or forming groups

1. Divide stage: use hashing function h_{B-1} to hash stream of incoming data into $B-1$ output buffers connected to $B-1$ partitions on disk.



2. Rerash/Conquer: use hashing func. to read $B-1$ partitions created in first stage into RAM hash table. Right out complete hash table to disk.

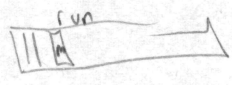
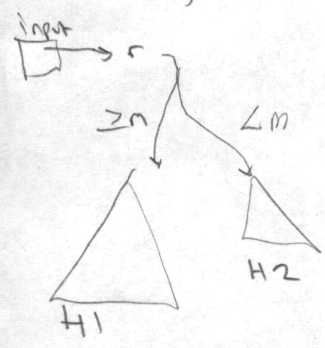
Runtime: $4N$

Biggest trade in two passes = $B(B-1)$, 1st stage creates $B-1$ partitions and partitions can't be more than B pages. Recursive partitioning can be used for data that

can't be hashed in 2 passes.

Replacement Sort/Tournament/Heap Sort:

- load in $B-2$ pages into "current run" heap H1
- keep H2 initially empty for next run.
- Pop off an m and write to run
- read in r from memory and if $\geq m$ add r to H1, if $< m$ add to H2.
- done when H1 empty. H2 starts next run, H2 guaranteed size $B-2$
- on average runs sized $2(B-2)$ nodes.



	Heap File	Sorted File	Clustered File
Scan all records	$B \cdot D$	$B \cdot D$	$1.5 \cdot B \cdot D$
Equality search	$\frac{1}{2} B \cdot D$	$(\log_2 B) \cdot D$	$(\log_2 1.5 B + 1) \cdot D$
Range Search	$B \cdot D$	$(\log_2 B + \# \text{match}) \cdot D$	$(\log_2 1.5 B + \# \text{match}) \cdot D$
Insert	$D + D$	$(\log_2 B + \frac{1}{2} B \cdot \frac{1}{2} B)$	$(\log_2 1.5 B + 2) \cdot D$
Delete	$\frac{1}{2} B \cdot D + D$	$\frac{1}{2} B \cdot D$	$(\log_2 1.5 B + 2) \cdot D$
average!		D	

Page Formats

Packed: store records in 1st N slots ($N = \text{records on page}$): when record deleted, move last record into vacated slot: all empty slots at end. Con: Doesn't work if use external ref. to main record (rid has slot # which moves)

Unpacked (Bitmap): handle deletions with array of bits locating record takes scanning bit map to find slots where bit on.

Slotted Page: used for variable length records and fixed length

Double Buffering

- CPU idle when have to read in input buffers and write out out buffer.
- allocate additional pages to input or output that are used when other set reading or writing so CPU not idle.

- Pro: need to move records on page for reasons other than tracking space freed by deletions
- when record deleted, more records to fill hole so free space contiguous
- maintain directory of slots/page with (offset, length) pair
- delete by setting offset to -1
- keep pointer to free space, when trying to put in record too large more records to reclaim freed space.

Terminology for Tree Index

- Fanout (F): # of children per non-leaf node
- Order (d): min # of entries per node
- N_i : # of leaf pages

- height of B+ tree depends on # data entries / size of index entries
- if search key values long not many index entries fit on page, fan-out low height large.
- high fan out more space efficient for duplicates.

Join Algorithms: $M = \text{pages in outer } (R)$

Nested Loops Join:

$N = \text{pages in inner } (S)$
 $P_r = \text{records/page}$
 Stream on outside.
 - works for continuous streams
 - tuple at a time
 - scan R and for each tuple in R scan all of S
 total cost = $M + P_r \cdot M \cdot N$

Refinement: page at a time: each page R retrieve page S

cost = $M + M \cdot N$

Optimize: choose R to be smaller

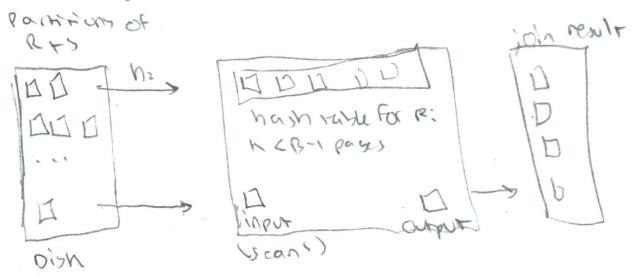
Index - Nested Loop Join

- index on on relations, on join attribute:
 make indexed relation be inner relation.
 - compare r only with tuples s in same partition. = have same value in join.
 - cost of scanning $R = M$, cost of retrieving matching S depends on kind of index / # of matches.
 1) if index on S in $B+$, cost 2-4/10. In hash 1-2
 2) once find leaf, cost of retrieving depends on clustering. If clustered cost +1, if not could be 1 per match S tuple. (on diff. pages in worst case)

cost = $M + (M \cdot P_r) (\text{cost to leaf}) + M P_r (\text{cost to retrieve})$

Hash - Join: works for equi-join, can't use if can't sort ie infinite stream

- uses hashing to identify partitions
 - partitioning phase similar to external hashing
 - probing phase illustrated:



- in partitioning, scan R and S both once and write once $\Rightarrow 2(M+N)$; 2nd phase scan each partition once for additional $M+N$

cost = $3(M+N)$ assuming length(partitions) $\leq B$ pages

Hybrid Hash Join: can handle infinite streams but bad choice
 - keep one of hash buckets in memory

select From Table A Left Join B on A.key=B.key
 select From A Inner Join B on A.key=B.key
 select From A Right Join B on A.key=B.key where A.key is null
 select From A Full outer Join B on A.key=B.key where A.key is null or B.key is null

Block Nested Loops Join: bad for streams unless outer fits in memory

- suppose enough mem to hold R with 2 extra buffer pages
 - read in smaller relation, use buffer scan S .

cost = $M + N$ - if R fits in buffer

refinement: build in memory hash for smaller relation R , minimizes CPU cost.

- when R can't be loaded completely in; load in $B-2$ pages of R , scan S .

cost = $M + N \lceil \frac{M}{B-2} \rceil$

Sort Merge Join: can't use for duplicates or continuous streams

- sort both relations on join attribute then look for qualifying tuples by merging the two.
 - external sorting alg. used to sort, if already sorted then no need.

Assuming no duplicates: cost = cost sort R + cost sort S + $M+N$

Refinement: combine the merging phase of sorting with merging phase of join.

- produce sorted runs of size B for R and S , $B \geq \sqrt{L}$ where L is size of larger, # runs per relation $\leq \sqrt{L}$
 - suppose # buffers available for merging at least $2\sqrt{L}$; more than total # of runs of R and S . allocate 1 buffer per run of R and 1 for S .

- merge runs of R , runs of S , and resulting R and S streams.

- increases # of buffers to $2\sqrt{L}$, unless use tech. to produce runs of size $2B$, make it back to \sqrt{L}

cost = $3(M+N)$

- enough memory for refinement: $B \geq 1 + \frac{R}{B} + \frac{S}{B}$
 basically $N = R+S$
 - mem requirement by tournament sort $\lceil \frac{S}{2(B-2)} \rceil + \lceil \frac{R}{2(B-2)} \rceil \leq B-1$
Clustered index

- index data entries stored in approx order by value of search keys
 - file can be clustered on at most one search key
 Pros: good for RANGE queries, potential locality
 cons: more expensive to maintain.

Indexes

Index - disk based data structure for fast lookup by value

Search key: any subset of columns in relation, need not be a key of relation (can be multiple items)

3 Alternatives

- 1) Actual data record (with key values)
- 2) $\langle k, rid \text{ of matching data record} \rangle$
- 3) $\langle k, \text{list of rids of matching data records} \rangle$

CANNOT SORT INFINITE STREAMS