

Announcements (Mon., Mar. 20)

- 3.1 and 3.2 are due on Wednesday March 22
- Milestone 2 due next Monday March 27
- · Feedback posted on private piazza threads

Where are we?

- We are covering DB internals and query processing • So far:
- Index: mostly B+ tree

• Today:

· finish query processing and join algorithms

Query Processing

Overview

- Many different ways of processing the same query • Scan? Sort? Hash? Use an index?
 - All have different performance characteristics and/or make different assumptions about data
- Best choice depends on the situation
 - Implement all alternatives
 - Let the query optimizer choose at run-time
 - Often not the "best choice"
 - · Optimizer tries NOT to select a "bad choice"

Notation

- Relations: R, S
- Tuples: r, s
- Number of tuples: |R|, |S|
- Number of disk blocks: *B*(*R*), *B*(*S*)
- Number of memory blocks available: M
- Cost metric • Number of I/O's
 - Memory requirement
- We do not count the cost of final write to disk
- Do not try to memorize the formulas for cost estimation!
 - understand the logic
 - · recall the diagram of disk and memory on whiteboard



Table scan

- Scan table R and process the query
 - Selection over R
 - Projection of R without duplicate elimination
- I/O's: **B(R)**
 - Trick for selection: stop early if it is a lookup by key
- Memory requirement: 2
- Not counting the cost of writing the result out • Same for any algorithm!
 - Maybe not needed—results may be pipelined into another operator







Toy example

• 3 memory blocks available; each holds one number

```
• Input: 1, 7, 4, 5, 2, 8, 3, 6, 9

    Pass o

    1, 7, 4 → 1, 4, 7

    • 5, 2, 8 \rightarrow 2, 5, 8

    9, 6, 3 → 3, 6, 9

• Pass 1

    1, 4, 7 + 2, 5, 8 → 1, 2, 4, 5, 7, 8

    • 3, 6, 9
• Pass 2 (final)

    1, 2, 4, 5, 7, 8 + 3, 6, 9 → 1, 2, 3, 4, 5, 6, 7, 8, 9
```





Some tricks for sorting

Double buffering

- · Allocate an additional block for each run
- Overlap I/O with processing
- Trade-off: smaller fan-in (more passes)
- Instead of reading/writing one disk block at time, read/write a bunch ("cluster")
- More sequential I/O's
- Trade-off: larger cluster → smaller fan-in (more passes)

Sort-merge join

$R \bowtie_{R.A=S.B} S$

• Sort *R* and *S* by their join attributes; then merge r, s = the first tuples in sorted R and S Repeat until one of *R* and *S* is exhausted: If r. A > s. B then s = next tuple in Selse if r.A < s.B then r = next tuple in Relse output all matching tuples, and r, s = next in R and S

• I/O's: sorting + 2B(R) + 2B(S)

- In most cases (e.g., join of key and foreign key)
- Worst case is $B(R) \cdot B(S)$: everything joins

Example of merge join R: S: $R \bowtie_{R.A=S.B} S$: \implies $s_1.B = 1$ $\rightarrow r_1 \cdot A = 1$ r_1s_1 $\rightarrow s_2 \cdot B = 2$ $r_2 \cdot A = 3$ r_2s_3 $r_3.A = 3$ \implies $s_3.B = 3$ r_2s_4 $\rightarrow r_4.A = 5$ $s_4.B = 3$ $\rightarrow r_5.A = 7$ \implies $s_5.B = 8$ r_3s_3 \rightarrow $r_6.A = 7$ r_3s_4 $\rightarrow r_7.A = 8$ $r_{7}s_{5}$





Other sort-based algorithms

- Union (set), difference, intersection
 More or less like SMJ
- Duplication elimination
 - External merge sort
 - Eliminate duplicates in sort and merge
- Grouping and aggregation
 - External merge sort, by group-by columns
 - Trick: produce "partial" aggregate values in each run, and combine them during merge
 - This trick doesn't always work though
 - Examples: SUM(DISTINCT ...), MEDIAN(...)

Hashing-based algorithms



Hash join $R \bowtie_{R,A=S,B} S$ • Main idea • Partition *R* and *S* by hashing their join attributes, and then consider corresponding partitions of *R* and *S*. • If *r*. *A* and *s*. *B* get hashed to different partitions, they con't join • State - Consider solutions •





