CSE 444: Database Internals

Lecture 7 Query Execution and Operator Algorithms (part 1)

CSE 444 - Spring 2015

Announcements

• Lab 2 / part 1 due Friday, 11pm

• CSE544M: review 2 due today, 11pm

What We Have Learned So Far

- Overview of the architecture of a DBMS
- Access methods
 - Heap files, sequential files, Indexes (hash or B+ trees)
- Role of buffer manager
- Practiced the concepts in hw1 and lab1

DBMS Architecture



Next Lectures

- How to answer queries efficiently!
 - Physical query plans and operator algorithms
- How to automatically find good query plans
 - How to compute the cost of a complete plan
 - How to pick a good query plan for a query
 - i.e., Query optimization

Query Evaluation Steps Review







Physical Query Plan

• Access path selection for each relation:

- File scan, or
- Index lookup with a predicate
- Implementation choice for each operator
 - We will learn different algorithms
- Scheduling decisions for operators
 - Pipelined execution, or
 - Intermediate tuple materialization

Iterator Interface

- open()
 - Initializes operator state
 - Sets parameters such as selection condition
- next()
 - Operator invokes get_next() recursively on its inputs
 - Performs processing and produces an output tuple
- close(): clean-up state

Pipelined Query Execution



Pipelined Query Execution



Pipelined Execution

- Tuples generated by an operator are immediately sent to the parent
- Benefits:
 - No operator synchronization issues
 - Saves cost of writing intermediate data to disk
 - Saves cost of reading intermediate data from disk
- This approach is used whenever possible

Intermediate Tuple Materialization

• Tuples generated by an operator are written to disk an in intermediate table

- No direct benefit
- Necessary:
 - For certain operator implementations
 - When we don't have enough memory

Intermediate Tuple Materialization



CSE 444 - Spring 2015

Memory Management

Each operator:

- Pre-allocates heap space for tuples
 - Pointers to base data in buffer pool
 - Or new tuples on the heap
- Allocates memory for its internal state
 - Either on heap or buffer pool (depends on system)

DMBS may **limit** how much memory each operator, or each query can use

CSE 444 - Spring 2015

Query Execution Bottom Line

- SQL query transformed into physical plan
 - Access path selection for each relation
 - Implementation choice for each operator
 - Scheduling decisions for operators
- Execution of the physical plan is pull-based
- Operators given a limited amount of memory

Operator Algorithms

Operator Algorithms

Design criteria

- Cost: IO, CPU, Network
- Memory utilization
- Load balance (for parallel operators)

Cost Parameters

Cost = total number of I/Os

– This is a simplification that ignores CPU, network

- Parameters:
 - B(R) = # of blocks (i.e., pages) for relation R
 - T(R) = # of tuples in relation R
 - V(R, a) = # of distinct values of attribute a
 - When a is a key, V(R,a) = T(R)
 - When a is not a key, V(R,a) can be anything < T(R)

Convention

- Cost = the cost of reading operands from disk
- Cost of writing the result to disk is not included; need to count it separately when applicable

Example: Cost of Scanning a Table

- Result may be unsorted: B(R)
- Result needs to be sorted: 3B(R)
 - We will discuss sorting later

Outline

- Join operator algorithms
 - One-pass algorithms (Sec. 15.2 and 15.3)
 - Index-based algorithms (Sec 15.6)
 - Two-pass algorithms (Sec 15.4 and 15.5)
- Note about readings:
 - In class, we discuss only algorithms for joins
 - Other operators are easier: read the book

Join Algorithms

- Hash join
- Nested loop join
- Sort-merge join

Hash Join

Hash join: $R \bowtie S$

- Scan R, build buckets in main memory
- Then scan S and join
- Cost: B(R) + B(S)
- One-pass algorithm when $B(R) \le M$





Step 1: Scan Patient and build hash table in memoryCan be done in
method open()Memory M = 21 pages



Step 2: Scan Insurance and probe into hash table



Step 2: Scan Insurance and probe into hash table



Step 2: Scan Insurance and probe into hash table



Nested Loop Joins

- Tuple-based nested loop $\mathsf{R} \bowtie \mathsf{S}$
- R is the outer relation, S is the inner relation

```
\begin{array}{l} \label{eq:total_for_each} \hline \begin{tabular}{l} for each tuple $t_1$ in $R$ $do$ \\ \hline \begin{tabular}{l} for each tuple $t_2$ in $S$ $do$ \\ \hline \begin{tabular}{l} if $t_1$ and $t_2$ join $\underline{then}$ output $(t_1,t_2)$ \\ \end{array}
```

Nested Loop Joins

- Tuple-based nested loop $\mathsf{R} \bowtie \mathsf{S}$
- R is the outer relation, S is the inner relation

for each tuple t_1 in R do for each tuple t_2 in S do if t_1 and t_2 join then output (t_1, t_2)

- Cost: B(R) + T(R) B(S)
- Multiple-pass since S is read many times

 $\begin{array}{l} \label{eq:for} \mbox{for each page of tuples r in R do} \\ \mbox{for each page of tuples s in S do} \\ \mbox{for all pairs of tuples } t_1 \mbox{ in r, } t_2 \mbox{ in s} \\ \mbox{if } t_1 \mbox{ and } t_2 \mbox{ join } \mbox{then} \mbox{ output } (t_1,t_2) \end{array}$

• Cost: B(R) + B(R)B(S)







Block-Nested-Loop Refinement

Block-Nested-Loop Refinement

 $\begin{array}{l} \label{eq:starses} \begin{array}{l} \mbox{for each group of M-1 pages r in R } \mbox{do} \\ \mbox{for each page of tuples s in S } \mbox{do} \\ \mbox{for all pairs of tuples } t_1 \mbox{ in r, } t_2 \mbox{ in s} \\ \mbox{if } t_1 \mbox{ and } t_2 \mbox{ join } \mbox{then} \mbox{ output } (t_1,t_2) \end{array}$

• Cost: B(R) + B(R)B(S)/(M-1)

Sort-Merge Join

Sort-merge join: $R \bowtie S$

- Scan R and sort in main memory
- Scan S and sort in main memory
- Merge R and S
- Cost: B(R) + B(S)
- One pass algorithm when B(S) + B(R) <= M
- Typically, this is NOT a one pass algorithm

Step 1: Scan Patient and sort in memory



Step 2: Scan Insurance and sort in memory





Step 3: Merge Patient and Insurance

Memory M = 21 pages



Step 3: Merge Patient and Insurance



