

CSE 444: Database Internals

Lecture 7 Query Execution and Operator Algorithms (part 1)

Magda Balazinska - CSE 444, Spring 2013

1

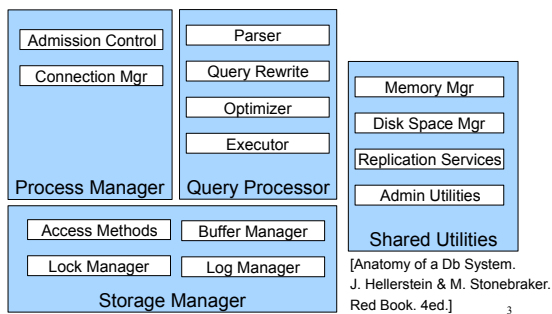
What We Have Learned So Far

- Overview of the architecture of a DBMS
 - Main components (summarized on next slide)
 - Steps involved in query evaluation
- Access methods
 - Storing data in heap files or sequential files
 - Indexes (hash or B+ trees)
- Role of buffer manager
- Interaction of components during query execution

Magda Balazinska - CSE 444, Spring 2013

2

DBMS Architecture



3

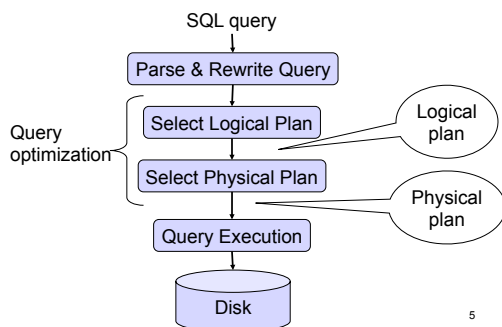
What We Will Learn Next

- How to answer queries **efficiently!**
 - Operator algorithms, especially for joins
 - How to leverage indexes for selections and joins
 - How to compute costs of individual operations
- 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?
- Start by reviewing physical plans

Magda Balazinska - CSE 444, Spring 2013

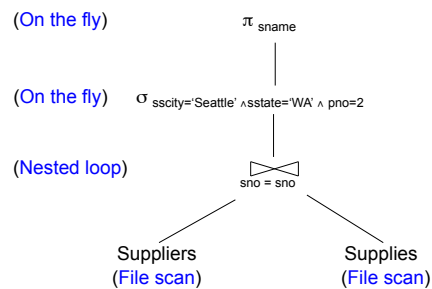
4

Query Evaluation Steps Review



5

Physical Query Plan



Magda Balazinska - CSE 444, Spring 2013

6

Physical Query Plan

- Logical query plan with extra annotations
- **Access path selection** for each relation
 - Use a file scan or use an index with a predicate
 - We learned various alternatives in past few lectures
- **Implementation choice** for each operator
 - We will learn about operator algorithms
- **Scheduling decisions** for operators
 - Pipelined execution
 - Or intermediate tuple materialization

Magda Balazinska - CSE 444, Spring 2013

7

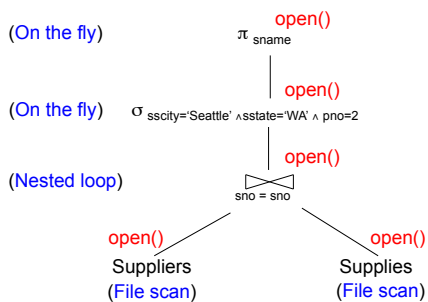
Iterator Interface

- Each **operator implements this interface**
- Interface has only three methods
- **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

Magda Balazinska - CSE 444, Spring 2013

8

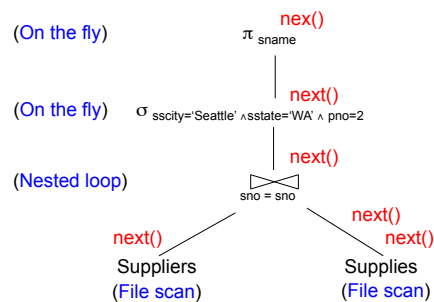
Pipelined Query Execution



Magda Balazinska - CSE 444, Spring 2013

9

Pipelined Query Execution



Magda Balazinska - CSE 444, Spring 2013

10

Pipelined Execution

- Applies parent operator to tuples directly as they are produced by child operators
- Benefits
 - No operator synchronization issues
 - Saves cost of writing intermediate data to disk
 - Saves cost of reading intermediate data from disk
 - Good resource utilizations on single processor
- This approach is used whenever possible

Magda Balazinska - CSE 444, Spring 2013

11

Intermediate Tuple Materialization

- Writes the results of an operator to an intermediate table on disk
- No direct benefit but
- Necessary for some operator implementations
- When operator needs to examine the same tuples multiple times

Magda Balazinska - CSE 444, Spring 2013

12

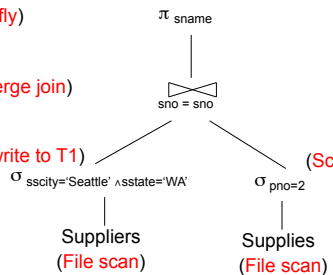
Intermediate Tuple Materialization

(On the fly)

(Sort-merge join)

(Scan: write to T1)

(Scan: write to T2)



Magda Balazinska - CSE 444, Spring 2013

13

Memory Management

- Each operator has some pre-allocated space for input and output tuples
 - These are tuples in-flight between operators
 - These input/output queues either point to base data in buffer pool
 - Or point to new tuples on the heap
- Each operator can also allocate separate memory for its internal state
 - Either on heap or buffer pool (depends on system)

Magda Balazinska - CSE 444, Spring 2013

14

Memory Management (cont.)

- DBMS limits how much memory each operator can use
 - This limit can be configurable
- DBMS limits how much memory each query can use
 - This limit can be configurable

Magda Balazinska - CSE 444, Spring 2013

15

Query Execution Bottom Line

- SQL query transformed into **physical plan**
 - Logical query plan with extra annotations
 - **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

Magda Balazinska - CSE 444, Spring 2013

16

Operator Algorithms

Magda Balazinska - CSE 444, Spring 2013

17

Why Learn About Op Algos?

- Good algorithms can greatly improve performance
 - Need to know operator algorithms to understand query plans
 - Need to understand query plans to tune a DBMS
- Implemented in commercial DBMSs
 - DBMSs implement different subsets of known algorithms
- Operator costs are first step toward query optimization
- Basic ideas to achieve high performance in operators go beyond relational operators

Magda Balazinska - CSE 444, Spring 2013

18

Operator Algorithms

- How to compare implementations/algorithms?
 - Using a cost model: IO, CPU, (and network bw)
 - Later, will see how this plays a role in optimization
- Some key design criteria
 - Cost: Different algorithms have different costs
 - Cost depends on input data and other parameters
 - Memory utilization
 - Operators only have access to limited amount of memory
 - Load balance (for parallel operators)

Magda Balazinska - CSE 444, Spring 2013

19

Cost Parameters

- In database systems the data is on disk
- **Cost = total number of I/Os**
 - This is a simplification
 - Normally, need to consider IO, CPU, and 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 $< T(R)$

20

Cost

- Cost of an operation = number of disk I/Os to
 - Read the operands
 - Compute the result
- Cost of writing the result to disk is *not included*
 - Need to count it separately when applicable

Magda Balazinska - CSE 444, Spring 2013

21

Cost of Scanning a Table

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

Magda Balazinska - CSE 444, Spring 2013

22

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 will discuss only algorithms for join operator (because other operators are easier)
 - Read the book to get more details about these algos
 - Read the book to learn about algos for other operators

Magda Balazinska - CSE 444, Spring 2013

23

Basic Join Algorithms

- Logical operator:
 - $\text{Product}(\text{pname, cname}) \bowtie \text{Company}(\text{cname, city})$
- Propose three physical operators for the join, assuming the tables are in main memory:
 - **Hash join**
 - **Nested loop join**
 - **Sort-merge join**

Magda Balazinska - CSE 444, Spring 2013

24

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 $B(R) \leq M$
 - By "one pass", we mean that the operator reads its operands only once. It does not write intermediate results back to disk.

Magda Balazinska - CSE 444, Spring 2013

25

Hash Join Example

Patient(pid, name, address)

Insurance(pid, provider, policy_nb)

Patient \bowtie Insurance

Patient

1	'Bob'	'Seattle'
2	'Ela'	'Everett'
3	'Jill'	'Kent'
4	'Joe'	'Seattle'

Insurance

2	'Blue'	123
4	'Prem'	432
4	'Prem'	343
3	'GrpH'	554

Two tuples per page

26

Hash Join Example

Patient \bowtie Insurance

Some large-enough nb

Memory M = 21 pages

Showing pid only

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

This is one page with two tuples

27

Hash Join Example

Step 1: Scan Patient and create hash table in memory

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------



Input buffer

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

28

Hash Join Example

Step 2: Scan Insurance and probe into hash table

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------

Input buffer

Output buffer

Write to disk or pass to next operator

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

29

Hash Join Example

Step 2: Scan Insurance and probe into hash table

Memory M = 21 pages

Hash h: pid % 5

5	1 6 2	3 8 4 9
---	-------	---------

Input buffer

Output buffer

Disk

Patient	Insurance
1 2	2 4 6 6
3 4	4 3 1 3
9 6	2 8
8 5	8 9

30

Hash Join Example

Step 2: Scan Insurance and probe into hash table

Memory M = 21 pages
Hash h: pid % 5

5	1 6	2	3 8	4 9
---	-----	---	-----	-----

Input buffer: 4 3 Output buffer: 4 4

Keep going until read all of Insurance

Cost: $B(R) + B(S)$

31

Hash Join Details

```

Open() {
  H = newHashTable( );
  R.Open( );
  x = R.GetNext( );
  while (x != null) {
    H.insert(x); x = R.GetNext( );
  }
  R.Close( );
  S.Open( );
  buffer = [ ];
}

```

32

Hash Join Details

```

GetNext() {
  while (buffer == [ ]) {
    x = S.GetNext( );
    if (x==Null) return NULL;
    buffer = H.find(x);
  }
  z = buffer.first( );
  buffer = buffer.reset( );
  return z;
}

```

Magda Balazinska - CSE 444, Spring 2013

33

Hash Join Details

```

Close() {
  release memory (H, buffer, etc.);
  S.Close( );
}

```

Magda Balazinska - CSE 444, Spring 2013

34

Nested Loop Joins

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

```

for each tuple r in R do
  for each tuple s in S do
    if r and s join then output (r,s)

```

- Cost: $B(R) + T(R)B(S)$
- Not quite one-pass since S is read many times

Magda Balazinska - CSE 444, Spring 2013

35

Page-at-a-time Refinement

```

for each page of tuples r in R do
  for each page of tuples s in S do
    for all pairs of tuples
      if r and s join then output (r,s)

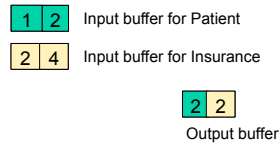
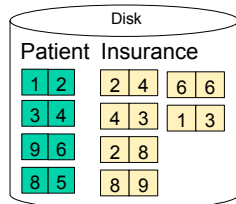
```

- Cost: $B(R) + B(R)B(S)$

Magda Balazinska - CSE 444, Spring 2013

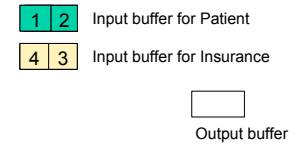
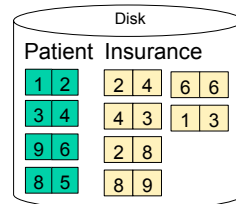
36

Nested Loop Example



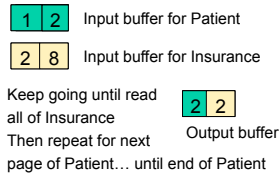
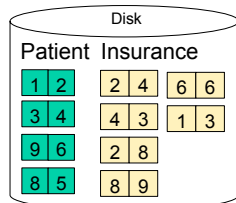
37

Nested Loop Example



38

Nested Loop Example



Cost: $B(R) + B(R)B(S)$

39

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 $B(S) + B(R) \leq M$
- Typically, this is NOT a one pass algorithm

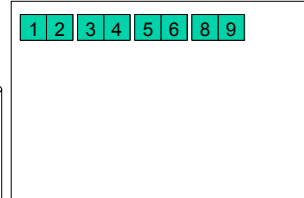
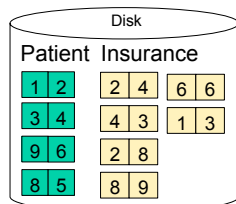
Magda Balazinska - CSE 444, Spring 2013

40

Sort-Merge Join Example

Step 1: Scan Patient and sort in memory

Memory M = 21 pages

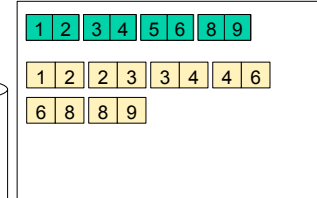
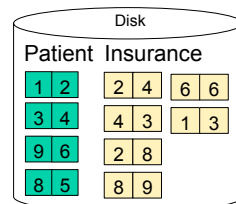


41

Sort-Merge Join Example

Step 2: Scan Insurance and sort in memory

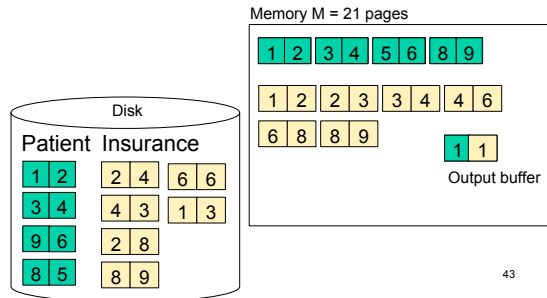
Memory M = 21 pages



42

Sort-Merge Join Example

Step 3: Merge Patient and Insurance



Sort-Merge Join Example

Step 3: Merge Patient and Insurance

