

# Chapter 15

## Algorithms for Query Processing and Optimization



5th Edition

Elmasri / Navathe

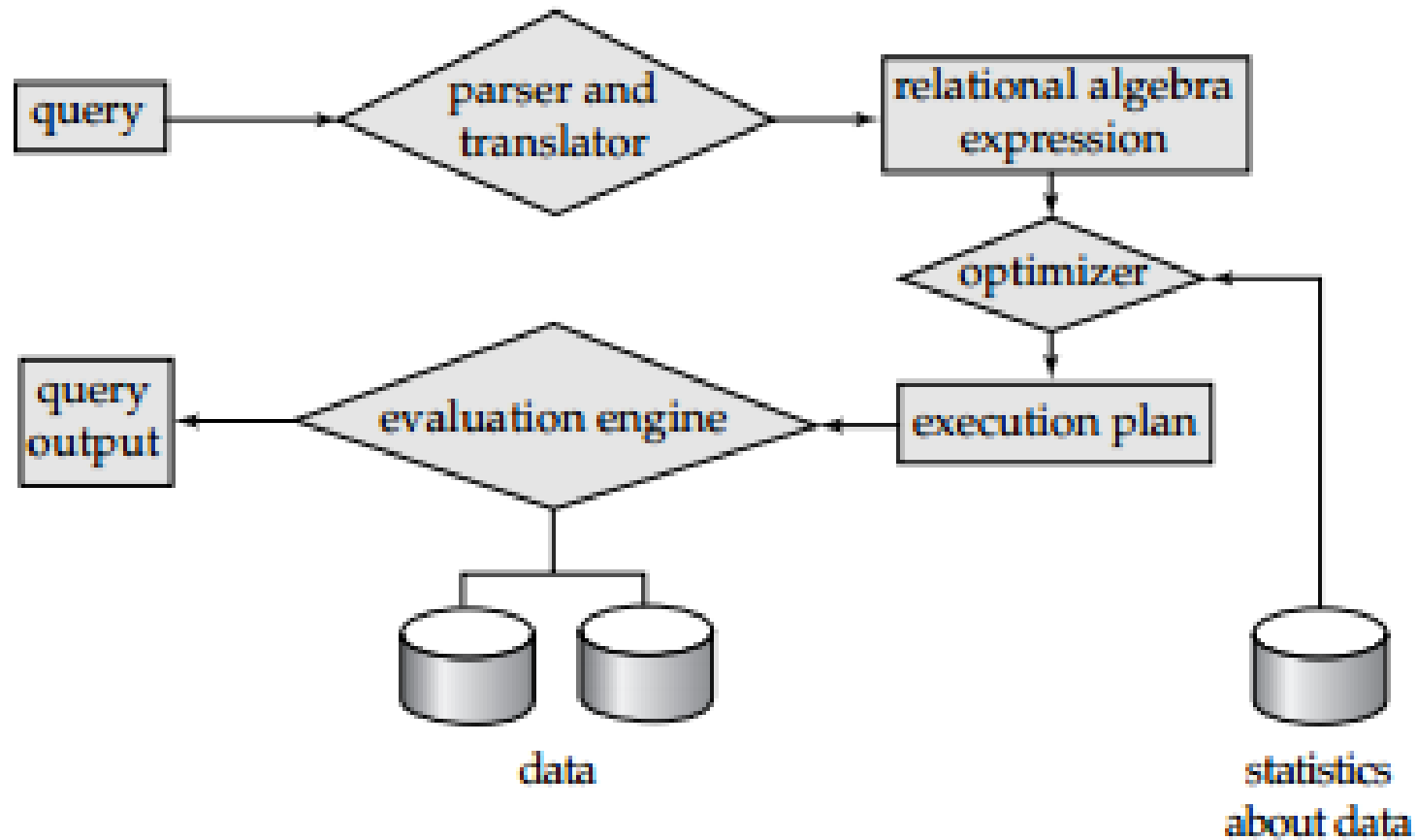
# Chapter Outline (1)

- 0. Introduction to Query Processing
- 1. Translating SQL Queries into Relational Algebra
- 2. Algorithms for External Sorting
- 3. Algorithms for SELECT and JOIN Operations
- 4. Algorithms for PROJECT and SET Operations
- 5. Implementing Aggregate Operations and Outer Joins

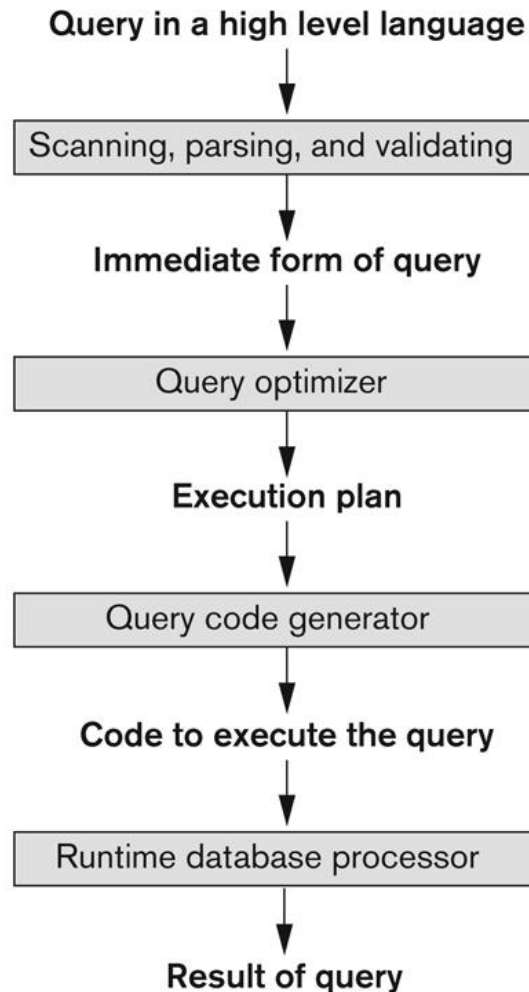
# 0. Introduction to Query Processing (1)

- **Query processing:**
  - Is the list of activities that are performed to obtain the required tuples that satisfy a given query.
- **Query optimization:**
  - The process of choosing a suitable execution strategy for processing a query.
- **Two internal representations of a query:**
  - **Query Tree**
  - **Query Graph**

# Introduction to Query Processing (2)



# Introduction to Query Processing (2)



**Figure 15.1**

Typical steps when processing a high-level query.

**Code can be:**

Executed directly (interpreted mode)  
Stored and executed later whenever needed (compiled mode)

# 1. Translating SQL Queries into Relational Algebra

- **Parser & Translator:**

- Syntax
- Schema element
- Converts the query into R.A expression.

- **Optimizer**

- Find all equivalent R.A expressions
- Find the R.A expression with least cost
- Cost(CPU, Block access, time spent)
- Will create query evaluation plan which tell what R.A and what algorithm is used.

- **Query evaluation plan:**

- Evaluate the above plan and get the result

# 1. Translating SQL Queries into Relational Algebra

- **Query block:**
  - The basic unit that can be translated into the algebraic operators and optimized.
- A query block contains a single SELECT-FROM-WHERE expression, as well as GROUP BY and HAVING clause if these are part of the block.
- **Nested queries** within a query are identified as separate query blocks.
- Aggregate operators in SQL must be included in the extended algebra.

# Translating SQL Queries into Relational Algebra

<b>SELECT</b>	LNAME, FNAME	<b>SELECT</b>	MAX (SALARY)
<b>FROM</b>	EMPLOYEE	<b>FROM</b>	EMPLOYEE
<b>WHERE</b>	SALARY > (	<b>WHERE</b>	DNO = 5);

SELECT LNAME, FNAME  
FROM EMPLOYEE  
WHERE SALARY > C

$\pi_{\text{LNAME, FNAME}} (\sigma_{\text{SALARY} > C} (\text{EMPLOYEE}))$

SELECT MAX (SALARY)  
FROM EMPLOYEE  
WHERE DNO = 5

$\mathcal{F}_{\text{MAX SALARY}} (\sigma_{\text{DNO}=5} (\text{EMPLOYEE}))$



## 2. Algorithms for External Sorting (1)

### ■ External sorting:

- Refers to sorting algorithms that are suitable for large files of records stored on disk that do not fit entirely in main memory, such as most database files.

### ■ Sort-Merge strategy:

- Starts by sorting small subfiles (**runs**) of the main file and then merges the sorted runs, creating larger sorted subfiles that are merged in turn.
- Sorting phase:  $n_R = \lceil (b/n_B) \rceil$
- Merging phase:  $d_M = \text{Min}(n_B-1, n_R)$ ;  $n_P = \lceil (\log_{d_M}(n_R)) \rceil$
- $n_R$ : number of initial runs;  $b$ : number of file blocks;
- $n_B$ : available buffer space;  $d_M$ : degree of merging;
- $n_P$ : number of passes.

# Algorithms for External Sorting (2)

```
set   $i \leftarrow 1$ ;  
      $j \leftarrow b$ ;           {size of the file in blocks}  
      $k \leftarrow n_B$ ;       {size of buffer in blocks}  
      $m \leftarrow \lceil (j/k) \rceil$ ;  
  
{Sort Phase}  
while ( $i \leq m$ )  
do {  
    read next  $k$  blocks of the file into the buffer or if there are less than  $k$  blocks  
    remaining, then read in the remaining blocks;  
    sort the records in the buffer and write as a temporary subfile;  
     $i \leftarrow i + 1$ ;  
}  
  
{Merge Phase: merge subfiles until only 1 remains}  
set   $i \leftarrow 1$ ;  
      $p \leftarrow \lceil \log_{k-1} m \rceil$ ; { $p$  is the number of passes for the merging phase}  
      $j \leftarrow m$ ;  
while ( $i \leq p$ )  
do {  
     $n \leftarrow 1$ ;  
     $q \leftarrow \lceil (j/(k-1)) \rceil$ ; {number of subfiles to write in this pass}  
    while ( $n \leq q$ )  
    do {  
        read next  $k-1$  subfiles or remaining subfiles (from previous pass)  
        one block at a time;  
        merge and write as new subfile one block at a time;  
         $n \leftarrow n + 1$ ;  
    }  
     $j \leftarrow q$ ;  
     $i \leftarrow i + 1$ ;  
}
```

**Figure 15.2**

Outline of the sort-merge algorithm for external sorting.

# 3. Algorithms for SELECT and JOIN Operations (1)

- Implementing the SELECT Operation
- Examples:
  - (OP1):  $\sigma_{SSN='123456789'}(EMPLOYEE)$
  - (OP2):  $\sigma_{DNUMBER>5}(DEPARTMENT)$
  - (OP3):  $\sigma_{DNO=5}(EMPLOYEE)$
  - (OP4):  $\sigma_{DNO=5 \text{ AND } SALARY>30000 \text{ AND } SEX=F}(EMPLOYEE)$
  - (OP5):  $\sigma_{ESSN=123456789 \text{ AND } PNO=10}(WORKS\_ON)$

# Algorithms for SELECT and JOIN Operations (2)

- Implementing the SELECT Operation (contd.):
- Search Methods for Simple Selection:
  - **S1 Linear search** (brute force):
    - Retrieve every record in the file, and test whether its attribute values satisfy the selection condition.
  - **S2 Binary search:**
    - If the selection condition involves an equality comparison on a key attribute on which the file is ordered, binary search (which is more efficient than linear search) can be used. (See OP1).
  - **S3 Using a primary index or hash key to retrieve a single record:**
    - If the selection condition involves an equality comparison on a key attribute with a primary index (or a hash key), use the primary index (or the hash key) to retrieve the record.

# Algorithms for SELECT and JOIN Operations (3)

- Implementing the SELECT Operation (contd.):
- Search Methods for Simple Selection:
  - **S4 Using a primary index to retrieve multiple records:**
    - If the comparison condition is  $>$ ,  $\geq$ ,  $<$ , or  $\leq$  on a key field with a primary index, use the index to find the record satisfying the corresponding equality condition, then retrieve all subsequent records in the (ordered) file.
  - **S5 Using a clustering index to retrieve multiple records:**
    - If the selection condition involves an equality comparison on a non-key attribute with a clustering index, use the clustering index to retrieve all the records satisfying the selection condition.
  - **S6 Using a secondary (B+-tree) index:**
    - On an equality comparison, this search method can be used to retrieve a single record if the indexing field has unique values (is a key) or to retrieve multiple records if the indexing field is not a key.
    - In addition, it can be used to retrieve records on conditions involving  $>$ ,  $\geq$ ,  $<$ , or  $\leq$ . (FOR RANGE QUERIES)

# Algorithms for SELECT and JOIN Operations (4)

- Implementing the SELECT Operation (contd.):
- Search Methods for Simple Selection:
  - **S7 Conjunctive selection using a individual index :**
    - If an attribute involved in any single simple condition in the conjunctive condition has an access path that permits the use of one of the methods S2 to S6, use that condition to retrieve the records and then check whether each retrieved record satisfies the remaining simple conditions in the conjunctive condition.
  - **S8 Conjunctive selection using a composite index**
    - If two or more attributes are involved in equality conditions in the conjunctive condition and a composite index (or hash structure) exists on the combined field, we can use the index directly. E.g index has been created on the composite key (Essn, Pno) of the relational R( Works\_on)

# Algorithms for SELECT and JOIN Operations (5)

- Implementing the SELECT Operation (contd.):
- Search Methods for Complex Selection:
  - **S9 Conjunctive selection by intersection of record pointers:**
    - This method is possible if secondary indexes are available on all (or some of) the fields involved in equality comparison conditions in the conjunctive condition and if the indexes include record pointers (rather than block pointers).
    - Each index can be used to retrieve the record pointers that satisfy the individual condition.
    - The intersection of these sets of record pointers gives the record pointers that satisfy the conjunctive condition, which are then used to retrieve those records directly.
    - If only some of the conditions have secondary indexes, each retrieved record is further tested to determine whether it satisfies the remaining conditions.

# Algorithms for SELECT and JOIN Operations (7)

- Implementing the SELECT Operation (contd.):
  - Whenever a **single condition** specifies the selection, we can only check whether an access path exists on the attribute involved in that condition.
    - If an access path exists, the method corresponding to that access path is used; otherwise, the “brute force” linear search approach of method S1 is used. (See OP1, OP2 and OP3)
  - For **conjunctive selection conditions**, whenever *more than one* of the attributes involved in the conditions have an access path, query optimization should be done to choose the access path that *retrieves the fewest records* in the most efficient way.
  - **Disjunctive selection conditions**



# Algorithms for SELECT and JOIN Operations (8)

- Implementing the JOIN Operation:
  - Join (EQUIJOIN, NATURAL JOIN)
    - two-way join: a join on two files
    - e.g.  $R \bowtie_{A=B} S$
    - multi-way joins: joins involving more than two files.
    - e.g.  $R \bowtie_{A=B} S \bowtie_{C=D} T$
- Examples
  - (OP6):  $\text{EMPLOYEE} \bowtie_{\text{DNO}=\text{DNUMBER}} \text{DEPARTMENT}$
  - (OP7):  $\text{DEPARTMENT} \bowtie_{\text{MGRSSN}=\text{SSN}} \text{EMPLOYEE}$

# Algorithms for SELECT and JOIN Operations (9)

- Implementing the JOIN Operation (contd.):
- Methods for implementing joins:
  - **J1 Nested-loop join** (brute force):
    - For each record  $t$  in  $R$  (outer loop), retrieve every record  $s$  from  $S$  (inner loop) and test whether the two records satisfy the join condition  $t[A] = s[B]$ .
  - **J2 Single-loop join** (Using an access structure to retrieve the matching records):
    - If an index (or hash key) exists for one of the two join attributes — say,  $B$  of  $S$  — retrieve each record  $t$  in  $R$ , one at a time, and then use the access structure to retrieve directly all matching records  $s$  from  $S$  that satisfy  $s[B] = t[A]$ .

# Algorithms for SELECT and JOIN Operations (10)

- Implementing the JOIN Operation (contd.):
- Methods for implementing joins:
  - **J3 Sort-merge join:**
    - If the records of R and S are *physically sorted (ordered)* by value of the join attributes A and B, respectively, we can implement the join in the most efficient way possible.
    - Both files are scanned in order of the join attributes, matching the records that have the same values for A and B.
    - In this method, the records of each file are scanned only once each for matching with the other file—unless both A and B are non-key attributes, in which case the method needs to be modified slightly.

# Algorithms for SELECT and JOIN Operations (11)

- Implementing the JOIN Operation (contd.):
- Methods for implementing joins:
  - **J4 Hash-join:**
    - The records of files R and S are both hashed to the *same hash file*, using the *same hashing function* on the join attributes A of R and B of S as hash keys.
    - A single pass through the file with fewer records (say, R) hashes its records to the hash file buckets.
    - A single pass through the other file (S) then hashes each of its records to the appropriate bucket, where the record is combined with all matching records from R.

# Algorithms for SELECT and JOIN Operations (14)

- Implementing the JOIN Operation (contd.):
- Factors affecting JOIN performance
  - Available buffer space
  - Join selection factor
  - Choice of inner VS outer relation

# Algorithms for SELECT and JOIN Operations (15)

- Implementing the JOIN Operation (contd.):
- Other types of JOIN algorithms
- Partition hash join
  - Partitioning phase:
    - Each file (R and S) is first partitioned into M partitions using a partitioning hash function on the join attributes:
      - $R_1, R_2, R_3, \dots, R_m$  and  $S_1, S_2, S_3, \dots, S_m$
    - Minimum number of in-memory buffers needed for the partitioning phase:  $M+1$ .
    - A disk sub-file is created per partition to store the tuples for that partition.
  - Joining or probing phase:
    - Involves M iterations, one per partitioned file.
    - Iteration i involves joining partitions  $R_i$  and  $S_i$ .

# Algorithms for SELECT and JOIN Operations (16)

- Implementing the JOIN Operation (contd.):
- Partitioned Hash Join Procedure:
  - Assume  $R_i$  is smaller than  $S_i$ .
    1. Copy records from  $R_i$  into memory buffers.
    2. Read all blocks from  $S_i$ , one at a time and each record from  $S_i$  is used to *probe* for a matching record(s) from partition  $S_i$ .
    3. Write matching record from  $R_i$  after joining to the record from  $S_i$  into the result file.

# Algorithms for SELECT and JOIN Operations (17)

- Implementing the JOIN Operation (contd.):
- Cost analysis of partition hash join:
  1. Reading and writing each record from R and S during the partitioning phase:
$$(b_R + b_S), (b_R + b_S)$$
  2. Reading each record during the joining phase:
$$(b_R + b_S)$$
  3. Writing the result of join:
$$b_{RES}$$
- Total Cost:
  - $3 * (b_R + b_S) + b_{RES}$



# Algorithms for SELECT and JOIN Operations (18)

- Implementing the JOIN Operation (contd.):
- **Hybrid hash join:**
  - Same as partitioned hash join except:
    - Joining phase of one of the partitions is included during the partitioning phase.
  - **Partitioning phase:**
    - Allocate buffers for smaller relation- one block for each of the M-1 partitions, remaining blocks to partition 1.
    - Repeat for the larger relation in the pass through S.)
  - **Joining phase:**
    - M-1 iterations are needed for the partitions R2 , R3 , R4 , .....Rm and S2 , S3 , S4 , .....Sm. R1 and S1 are joined during the partitioning of S1, and results of joining R1 and S1 are already written to the disk by the end of partitioning phase.

# 4. Algorithms for PROJECT and SET Operations (1)

- Algorithm for PROJECT operations (Figure 15.3b)

$\pi_{\langle \text{attribute list} \rangle}(R)$

1. If  $\langle \text{attribute list} \rangle$  has a key of relation  $R$ , extract all tuples from  $R$  with only the values for the attributes in  $\langle \text{attribute list} \rangle$ .
2. If  $\langle \text{attribute list} \rangle$  does NOT include a key of relation  $R$ , duplicated tuples must be removed from the results.

- Methods to remove duplicate tuples

1. Sorting
2. Hashing

# Algorithms for PROJECT and SET Operations (2)

- **Algorithm for SET operations**
- **Set operations:**
  - UNION, INTERSECTION, SET DIFFERENCE and CARTESIAN PRODUCT
- **CARTESIAN PRODUCT** of relations R and S include all possible combinations of records from R and S. The attribute of the result include all attributes of R and S.
- **Cost analysis of CARTESIAN PRODUCT**
  - If R has n records and j attributes and S has m records and k attributes, the result relation will have  $n*m$  records and  $j+k$  attributes.
- **CARTESIAN PRODUCT operation is very expensive and should be avoided if possible.**

# Algorithms for PROJECT and SET Operations (3)

- **Algorithm for SET operations (contd.)**
- **UNION** (See Figure 15.3c)
  - Sort the two relations on the same attributes.
  - Scan and merge both sorted files concurrently, whenever the same tuple exists in both relations, only one is kept in the merged results.
- **INTERSECTION** (See Figure 15.3d)
  - Sort the two relations on the same attributes.
  - Scan and merge both sorted files concurrently, keep in the merged results only those tuples that appear in both relations.
- **SET DIFFERENCE R-S** (See Figure 15.3e)
  - Keep in the merged results only those tuples that appear in relation R but not in relation S.

# 5. Implementing Aggregate Operations and Outer Joins (1)

- Implementing Aggregate Operations:
- **Aggregate operators:**
  - **MIN, MAX, SUM, COUNT and AVG**
- Options to implement aggregate operators:
  - **Table Scan**
  - **Index**
- Example
  - **SELECT MAX (SALARY)**
  - **FROM EMPLOYEE;**
- If an (ascending) index on SALARY exists for the employee relation, then the optimizer could decide on traversing the index for the largest value, which would entail following the right most pointer in each index node from the root to a leaf.

# Implementing Aggregate Operations and Outer Joins (2)

- Implementing Aggregate Operations (contd.):
- **SUM, COUNT and AVG**
- For a **dense index** (each record has one index entry):
  - Apply the associated computation to the values in the index.
- For a **non-dense index**:
  - Actual number of records associated with each index entry must be accounted for
- With **GROUP BY**: the aggregate operator must be applied separately to each group of tuples.
  - Use sorting or hashing on the group attributes to partition the file into the appropriate groups;
  - Computes the aggregate function for the tuples in each group.
- What if we have **Clustering index** on the grouping attributes?

# Implementing Aggregate Operations and Outer Joins (3)

- Implementing Outer Join:
- **Outer Join Operators:**
  - **LEFT OUTER JOIN**
  - **RIGHT OUTER JOIN**
  - **FULL OUTER JOIN.**
- The full outer join produces a result which is equivalent to the union of the results of the left and right outer joins.
- Example:

```
SELECT      FNAME, DNAME
FROM        (EMPLOYEE LEFT OUTER JOIN DEPARTMENT
              ON DNO = DNUMBER);
```
- Note: The result of this query is a table of employee names and their associated departments. It is similar to a regular join result, with the exception that if an employee does not have an associated department, the employee's name will still appear in the resulting table, although the department name would be indicated as null.

# Implementing Aggregate Operations and Outer Joins (4)

- Implementing Outer Join (contd.):
- **Modifying Join Algorithms:**
  - Nested Loop or Sort-Merge joins can be modified to implement outer join. E.g.,
    - For left outer join, use the left relation as outer relation and construct result from every tuple in the left relation.
    - If there is a match, the concatenated tuple is saved in the result.
    - However, if an outer tuple does not match, then the tuple is still included in the result but is padded with a null value(s).



# Implementing Aggregate Operations and Outer Joins (5)

- Implementing Outer Join (contd.):
- Executing a combination of relational algebra operators.
- Implement the previous left outer join example
  - {Compute the JOIN of the EMPLOYEE and DEPARTMENT tables}
    - $TEMP1 \leftarrow \pi_{FNAME, DNAME} (EMPLOYEE \bowtie_{DNO=DNUMBER} DEPARTMENT)$
  - {Find the EMPLOYEEs that do not appear in the JOIN}
    - $TEMP2 \leftarrow \pi_{FNAME} (EMPLOYEE) - \pi_{FNAME} (Temp1)$
  - {Pad each tuple in TEMP2 with a null DNAME field}
    - $TEMP2 \leftarrow TEMP2 \times 'null'$
  - {UNION the temporary tables to produce the LEFT OUTER JOIN}
    - $RESULT \leftarrow TEMP1 \cup TEMP2$
- The cost of the outer join, as computed above, would include the cost of the associated steps (i.e., join, projections and union).