



# Chapter 10: Query Processing

**Database System Concepts, 7<sup>th</sup> Ed.**

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



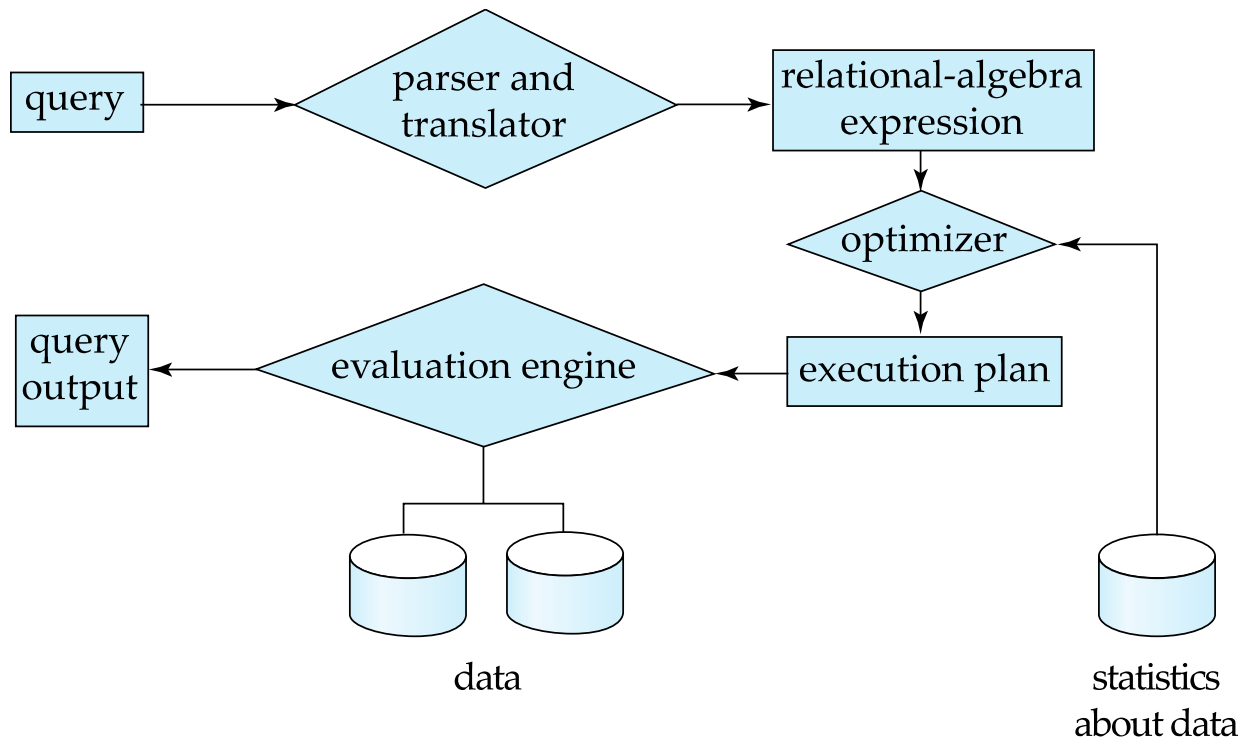
# Chapter 10: Query Processing

- Overview
- Measures of Query Cost
- Selection Operation
- Sorting
- Join Operation
- Other Operations
- Evaluation of Expressions



# Basic Steps in Query Processing

1. Parsing and translation
2. Optimization
3. Evaluation





# Basic Steps in Query Processing (Cont.)

- Parsing and translation
  - translate the query into its internal form. This is then translated into relational algebra expression.
  - Parser checks syntax and verifies validity relations.
- Evaluation
  - The query-execution engine takes a query-evaluation plan, executes that plan, and returns the answers to the query.



# Basic Steps in Query Processing: Optimization

- A relational algebra expression may have many equivalent forms
  - E.g.,  $\sigma_{salary < 75000}(\Pi_{salary}(instructor))$  is equivalent to  $\Pi_{salary}(\sigma_{salary < 75000}(instructor))$
- Each relational algebra operation can be evaluated using one of several different algorithms
  - Correspondingly, a relational-algebra expression can be evaluated in many ways.
- Annotated expression specifying a detailed evaluation strategy is called an **evaluation plan**. E.g.:
  - Use an index on *salary* to find instructors with a salary < 75000,
  - Or perform a complete relation scan and discard instructors with salary  $\geq 75000$



# Basic Steps: Optimization (Cont.)

- **Query Optimization:** Amongst all equivalent evaluation plans choose the one with the lowest execution cost.
  - Cost is estimated using statistical information from the database catalog
    - e.g. number of tuples in each relation, size of tuples, etc.
- In this chapter we study
  - How to measure query costs
  - Algorithms for evaluating relational algebra operations
  - How to combine algorithms for individual operations in order to evaluate a complete expression
- In the next chapter
  - We study how to optimize queries, that is, how to find an evaluation plan with the lowest estimated cost



# Measures of Query Cost

- Many factors contribute to time cost
  - *disk access, CPU, and network communication*
- Cost can be measured based on
  - **response time**, i.e., total elapsed time for answering a query, or
  - total **resource consumption**
- We use total resource consumption as a cost metric
  - Response time is harder to estimate, and minimizing resource consumption is a good idea in a shared database
- We ignore CPU costs for simplicity
  - Real systems do take CPU cost into account
  - Network costs must be considered for parallel systems
- We describe how to estimate the cost of each operation
  - We do not include the cost of writing output to disk



# Measures of Query Cost

- Disk cost can be estimated as:
  - Number of seeks \* average-seek-cost
  - Number of blocks read \* average-block-read-cost
  - Number of blocks written \* average-block-write-cost
- For simplicity we just use the **number of block transfers** *from disk and the number of seeks* as the cost measures
  - $t_T$  – time to transfer one block
    - Assuming for simplicity that the write cost is the same as the cost to read
  - $t_S$  – time for one seek
  - Cost for b block transfers plus S seeks
$$b * t_T + S * t_S$$
- $t_S$  and  $t_T$  depend on where data is stored; with 4 KB blocks:
  - High end magnetic disk:  $t_S = 4$  msec and  $t_T = 0.1$  msec
  - SSD:  $t_S = 20-90$  microsec and  $t_T = 2-10$  microsec for 4KB





# Measures of Query Cost (Cont.)

- Required data may be buffer resident already, avoiding disk I/O
  - But hard to consider for cost estimation
- Several algorithms can reduce disk IO by using extra buffer space
  - Amount of real memory available to buffer depends on other concurrent queries and OS processes, known only during execution
- **Worst case estimates** assume that no data is initially in the buffer and only the minimum amount of memory needed for the operation is available
  - But more optimistic estimates are used in practice



# Selection Operation

- **File scan**
- Algorithm **A1 (linear search)**. Scan each file block and test all records to see whether they satisfy the selection condition.
  - Cost estimate =  $b_r$  block transfers + 1 seek
    - $b_r$  denotes the number of blocks containing records of relation  $r$
  - If selection is on a key attribute, we can stop finding the record
    - cost =  $(b_r/2)$  block transfers + 1 seek
  - Linear search can be applied regardless of
    - selection condition or
    - ordering of records in the file, or
    - availability of indices
- Note: binary search generally does not make sense since data is not stored consecutively
  - except when there is an index available,
  - and binary search requires more seeks than index search



# Selections Using Indices

- **Index scan** – search algorithms that use an index
  - selection condition must be on search-key of the index.
- **A2 (clustering index, equality on key)**. Retrieve a single record that satisfies the corresponding equality condition
  - $Cost = (h_i + 1) * (t_T + t_S)$
- **A3 (clustering index, equality on a non-key)** Retrieve multiple records.
  - Records will be on consecutive blocks
    - Let  $b$  = number of blocks containing matching records
  - $Cost = h_i * (t_T + t_S) + t_S + t_T * b$



# Selections Using Indices

- **A4 (secondary index, equality on key/non-key).**
  - Retrieve a single record if the search-key is a candidate key
    - $Cost = (h_i + 1) * (t_T + t_S)$
  - Retrieve multiple records if search-key is not a candidate key
    - each of  $n$  matching records may be on a different block
    - $Cost = (h_i + n) * (t_T + t_S)$ 
      - Can be very expensive!



# Selections Involving Comparisons

- Can implement selections of the form  $\sigma_{A \leq V}(r)$  or  $\sigma_{A \geq V}(r)$  by using
  - a linear file scan,
  - or by using indices in the following ways:
- **A5 (clustering index, comparison)**. (Relation is sorted on  $A$ )
  - For  $\sigma_{A \geq V}(r)$  use index to find the first tuple  $\geq v$  and scan the relation sequentially from there
  - For  $\sigma_{A \leq V}^{\textcircled{R}}$  just scan the relation sequentially till first tuple  $> v$ ; do not use index
- **A6 (non-clustering index, comparison)**. (Relation is not sorted on  $A$ )
  - For  $\sigma_{A \geq V}(r)$  use index to find the first index entry  $\geq v$  and scan index sequentially from there, to find pointers to records.
  - For  $\sigma_{A \leq V}^{\textcircled{R}}$  just scan leaf pages of the index finding pointers to records, till the first entry  $> v$
  - In either case, retrieve records that are pointed to
  - requires an I/O per record; Linear file scan may be cheaper!



# Implementation of Complex Selections

- **Conjunction:**  $\sigma_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n}(r)$
- **A7 (conjunctive selection using one index).**
  - Select a combination of  $\theta_i$  and algorithms A1 through A7 that results in the least cost for  $\sigma_{\theta_i}(r)$ .
  - Test other conditions on the tuples fetched into the memory buffer.
- **A8 (conjunctive selection using a composite index).**
  - Use appropriate composite (multiple-key) index if available.
- **A9 (conjunctive selection by intersection of identifiers).**
  - Requires indices with record pointers.
  - Use the corresponding index for each condition and take the intersection of all the obtained sets of record pointers.
  - Then fetch records from the file
  - If some conditions do not have appropriate indices, apply the test in memory.



# Algorithms for Complex Selections

- **Disjunction:**  $\sigma_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n}(r)$ .
- **A10 (disjunctive selection by union of identifiers).**
  - Applicable if *all* conditions have available indices.
    - Otherwise use linear scan.
  - Use the corresponding index for each condition, and take the union of all the obtained sets of record pointers.
  - Then fetch records from the file
- **Negation:**  $\sigma_{\neg\theta}(r)$ 
  - Use linear scan on file
  - If very few records satisfy  $\neg\theta$ , and an index is applicable to  $\theta$ 
    - Find satisfying records using index and fetch from file



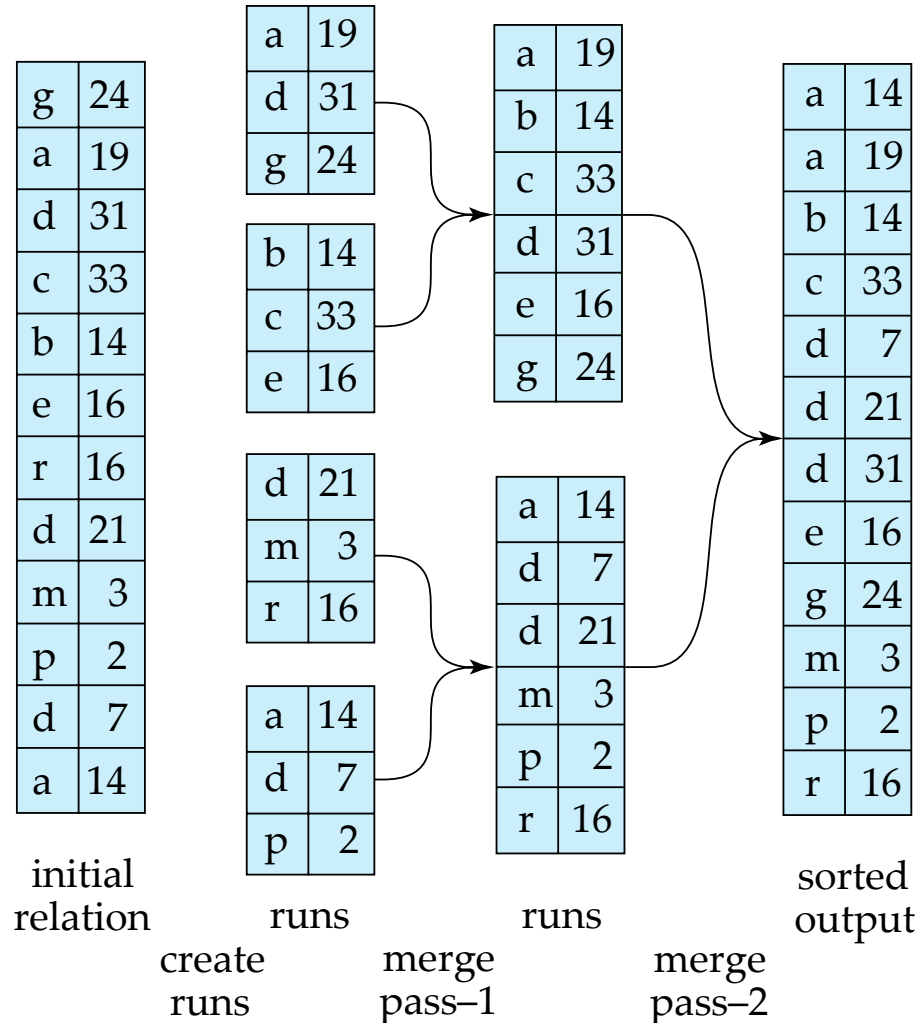
# Sorting

- We may build an index on the relation, and then use the index to read the relation in sorted order. This may lead to one disk block access for each tuple.
- For relations that fit in memory, techniques like quicksort can be used.
  - For relations that don't fit in memory, **external sort-merge** is a good choice.





# Example: External Sorting Using Sort-Merge





# External Sort-Merge

Let  $M$  denote memory size (in pages).

1. **Create sorted runs.** Let  $i$  be 0 initially.

Repeatedly do the following till the end of the relation:

- (a) Read  $M$  blocks of relation into memory
- (b) Sort the in-memory blocks
- (c) Write sorted data to run  $R_i$ ; increment  $i$ .

Let the final value of  $i$  be  $N$

2. *Merge the runs (next slide).....*



# External Sort-Merge (Cont.)

2. **Merge the runs ( $N$ -way merge).** We assume (for now) that  $N < M$ .
  1. Use  $N$  blocks of memory to buffer input runs, and 1 block to buffer the output. Read the first block of each run into its buffer page
  2. **repeat**
    1. Select the first record (in sort order) among all buffer pages
    2. Write the record to the output buffer. If the output buffer is full write it to disk.
    3. Delete the record from its input buffer page.  
**If** the buffer page becomes empty, **then**  
    read the following block (if any) of the run into the buffer.
  3. **until** all input buffer pages are empty:



# External Sort-Merge (Cont.)

- If  $N \geq M$ , several merge *passes* are required.
  - In each pass, contiguous groups of  $M - 1$  runs are merged.
  - A pass reduces the number of runs by a factor of  $M - 1$  and creates runs longer by the same factor.
    - E.g. If  $M=11$ , and there are 90 runs, one pass reduces the number of runs to 9, each of them 10 times the size of the initial runs
  - Repeated passes are performed till all runs have been merged into one.



# External Merge Sort (Cont.)

- Cost analysis:
  - 1 block per run leads to too many seeks during the merge
    - Instead use  $b_b$  buffer blocks per run
      - ➔ read/write  $b_b$  blocks at a time
    - Can merge  $\lfloor M/b_b \rfloor - 1$  runs in one pass
  - Total number of merge passes required:  $\lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil$ .
  - Block transfers for initial run creation as well as in each pass is  $2b_r$ 
    - For the final pass, we don't count the write cost
      - We ignore the final write cost for all operations since the output of an operation may be sent to the parent operation without being written to disk
    - Thus, the total number of block transfers for external sorting is:
$$b_r ( 2 \lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r/M) \rceil + 1 ) \lceil$$
  - Seeks: next slide



# External Merge Sort (Cont.)

- Cost of seeks
  - During run generation: one seek to read each run and one seek to write each run
    - $2 \lceil b_r / M \rceil$
  - During the merge phase
    - Need  $2 \lceil b_r / b_b \rceil$  seeks for each merge pass
      - except the final one which does not require a write
    - Total number of seeks:  
$$2 \lceil b_r / M \rceil + \lceil b_r / b_b \rceil (2 \lceil \log_{\lfloor M/b_b \rfloor - 1} (b_r / M) \rceil - 1)$$



# Join Operation

- Several different algorithms to implement joins
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- Choice based on the cost estimate
- Examples use the following information
  - Number of records of *student*: 5,000      *takes*: 10,000
  - Number of blocks of *student*:      100      *takes*:      400



# Nested-Loop Join

- To compute the theta join  $r \bowtie_{\theta} s$ 
  - for each** tuple  $t_r$  **in**  $r$  **do begin**
  - for each** tuple  $t_s$  **in**  $s$  **do begin**
  - test pair  $(t_r, t_s)$  to see if they satisfy the join condition  $\theta$
  - if they do, add  $t_r \cdot t_s$  to the result.
  - end**
  - end**
- $r$  is called the **outer relation** and  $s$  is the **inner relation** of the join.
- Requires no indices and can be used with any kind of join condition.
- Expensive since it examines every pair of tuples in the two relations.





## Nested-Loop Join (Cont.)

- In the worst case, if there is enough memory only to hold one block of each relation, the estimated cost is
$$n_r * b_s + b_r$$
 block transfers, plus  $n_r + b_r$  seeks
- If the smaller relation fits entirely in memory, use that as the inner relation.
  - Reduces cost to  $b_r + b_s$  block transfers and 2 seeks
- Assuming the worst-case memory availability, the cost estimate is
  - with *student* as outer relation:
    - $5000 * 400 + 100 = 2,000,100$  block transfers,
    - $5000 + 100 = 5100$  seeks
  - with *takes* as the outer relation
    - $10000 * 100 + 400 = 1,000,400$  block transfers and 10,400 seeks
- If smaller relation (*student*) fits entirely in memory, the cost estimate will be 500 block transfers.
- Block nested-loops algorithm (next slide) is preferable.



# Block Nested-Loop Join

- Variant of nested-loop join in which every block of inner relation is paired with every block of outer relation.

```
for each block  $B_r$  of  $r$  do begin  
  for each block  $B_s$  of  $s$  do begin  
    for each tuple  $t_r$  in  $B_r$  do begin  
      for each tuple  $t_s$  in  $B_s$  do begin  
        Check if  $(t_r, t_s)$  satisfy the join condition  
        if they do, add  $t_r \cdot t_s$  to the result.  
      end  
    end  
  end  
end  
end
```



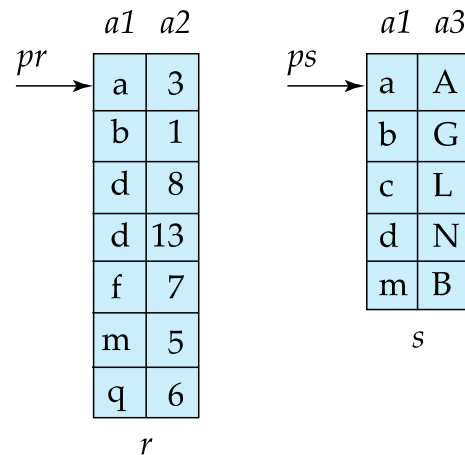
# Indexed Nested-Loop Join

- Index lookups can replace file scans if
  - join is the equi or natural join and
  - an index is available on the inner relation's join attribute
    - Can construct an index just to compute a join.
- For each tuple  $t_r$  in the outer relation  $r$ , use the index to look up tuples in  $s$  that satisfy the join condition with tuple  $t_r$ .
- Worst case: buffer has space for only one page of  $r$ , and, for each tuple in  $r$ , we perform an index lookup on  $s$ .
- Cost of the join:  $b_r(t_T + t_S) + n_r * c$ 
  - Where  $c$  is the cost of traversing the index and fetching all matching  $s$  tuples for one tuple of  $r$
  - $C$  can be estimated as the cost of a single selection on  $s$  using the join condition.
- If indices are available on join attributes of both  $r$  and  $s$ , use the relation with fewer tuples as the outer relation.



# Merge-Join

1. Sort both relations on their join attribute (if not already sorted on the join attributes).
2. Merge the sorted relations to join them
  1. Join step is similar to the merge stage of the sort-merge algorithm.
  2. Main difference is the handling of duplicate values in the join attribute — every pair with the same value on the join attribute must be matched
  3. Detailed algorithm in the book





# Merge-Join (Cont.)

- Can be used only for the equi and natural joins
- Each block needs to be read only once - assuming all tuples for any given value of the join attributes fit in memory
- Thus, the cost of merge join is:  
$$b_r + b_s \text{ block transfers} + \lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil \text{ seeks}$$

+ the cost of sorting if relations are unsorted.
- **hybrid merge-join:** If one relation is sorted, and the other has a secondary B<sup>+</sup>-tree index on the join attribute
  - Merge the sorted relation with the leaf entries of the B<sup>+</sup>-tree .
  - Sort the result on the addresses of the unsorted relation's tuples
  - Scan the unsorted relation in physical address order and merge with previous result, to replace addresses by the actual tuples
    - Sequential scan more efficient than random lookup

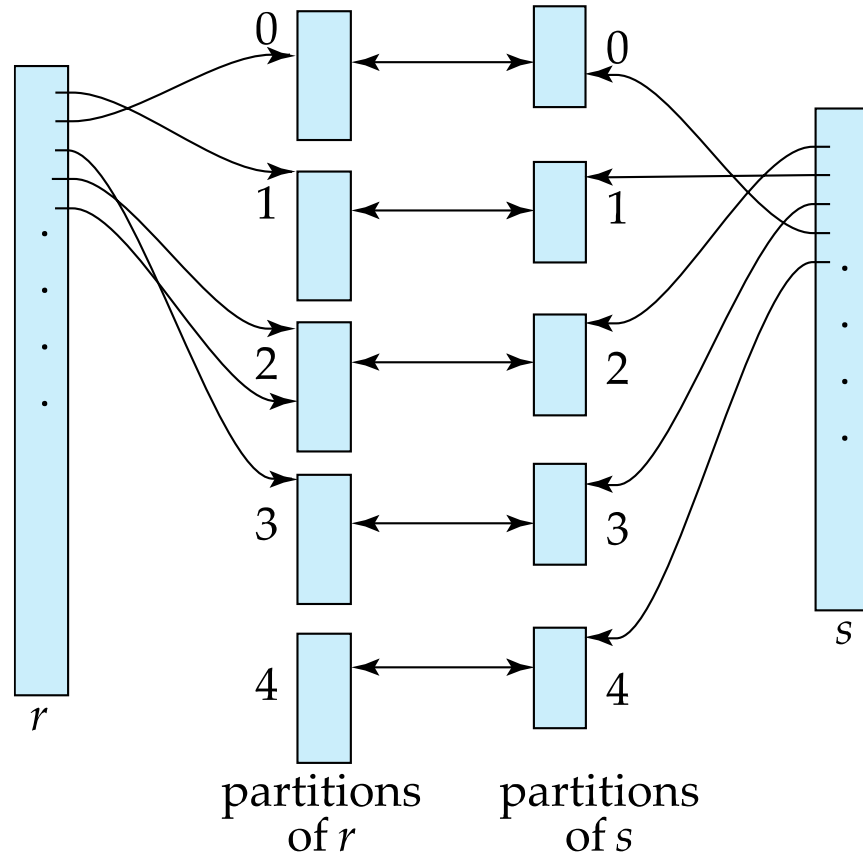


# Hash-Join

- Applicable for the equi and natural joins.
- A hash function  $h$  is used to partition tuples of both relations
- $h$  maps *JoinAttrs* values to  $\{0, 1, \dots, n\}$ , where *JoinAttrs* denotes the common attributes of  $r$  and  $s$  used in the natural join.
  - $r_0, r_1, \dots, r_n$  denote partitions of  $r$  tuples
    - Each tuple  $t_r \in r$  is put in partition  $r_i$  where  $i = h(t_r[\text{JoinAttrs}])$ .
  - $s_0, s_1, \dots, s_n$  denotes partitions of  $s$  tuples
    - Each tuple  $t_s \in s$  is put in partition  $s_i$ , where  $i = h(t_s[\text{JoinAttrs}])$ .
- *Note:* In book, Figure 12.10  $r_i$  is denoted as  $H_{r_i}$ ,  $s_i$  is denoted as  $H_{s_i}$  and  $n$  is denoted as  $n_h$ .



# Hash-Join (Cont.)





# Hash-Join Algorithm

The hash-join of  $r$  and  $s$  is computed as follows.

1. Partition the relation  $s$  using hashing function  $h$ . When partitioning a relation, one block of memory is reserved as the output buffer for each partition.
2. Partition  $r$  similarly.
3. For each  $i$ :
  - (a) Load  $s_i$  into memory and build an in-memory hash index on it using the join attribute. This hash index uses a different hash function than the earlier one  $h$ .
  - (b) Read the tuples in  $r_i$  from the disk one by one. For each tuple  $t_r$  locate each matching tuple  $t_s$  in  $s_i$  using the in-memory hash index. Output the concatenation of their attributes.

Relation  $s$  is called the **build input** and  $r$  is called the **probe input**.





# Complex Joins

- Join with a conjunctive condition:

$$r \bowtie_{\theta_1 \wedge \theta_2 \wedge \dots \wedge \theta_n} s$$

- Either use nested loops/block nested loops, or
- Compute the result of one of the simpler joins  $r \bowtie_{\theta_i} s$ 
  - final result comprises those tuples in the intermediate result that satisfy the remaining conditions

$$\theta_1 \wedge \dots \wedge \theta_{i-1} \wedge \theta_{i+1} \wedge \dots \wedge \theta_n$$

- Join with a disjunctive condition

$$r \bowtie_{\theta_1 \vee \theta_2 \vee \dots \vee \theta_n} s$$

- Either use nested loops/block nested loops, or
- Compute as the union of the records in individual joins  $r \bowtie_{\theta_i} s$ :

$$(r \bowtie_{\theta_1} s) \cup (r \bowtie_{\theta_2} s) \cup \dots \cup (r \bowtie_{\theta_n} s)$$



# Other Operations

- **Duplicate elimination** can be implemented via hashing or sorting.
  - On sorting, duplicates will come adjacent to each other, and all but one set of duplicates can be deleted.
  - *Optimization:* duplicates can be deleted during run generation as well as at intermediate merge steps in external sort-merge.
  - Hashing is similar – duplicates will come into the same bucket.
- **Projection:**
  - perform projection on each tuple
  - followed by duplicate elimination.



# Other Operations : Aggregation

- **Aggregation** can be implemented in a manner similar to duplicate elimination.
  - **Sorting** or **hashing** can be used to bring tuples in the same group together, and then the aggregate functions can be applied to each group.
  - Optimization: **partial aggregation**
    - combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values
    - For count, min, max, sum: keep aggregate values on tuples found so far in the group.
      - When combining partial aggregate for the count function, add up the partial aggregates
    - For avg, keep sum and count, and divide the sum by count at the end



# Other Operations : Set Operations

- **Set operations** ( $\cup$ ,  $\cap$  and  $-$ ): can either use variant of merge-join after sorting, or variant of hash-join.
- E.g., Set operations using hashing:
  1. Partition both relations using the same hash function
  2. Process each partition  $i$  as follows.
    1. Using a different hashing function, build an in-memory hash index on  $r_i$ .
    2. Process  $s_i$  as follows
      - $r \cup s$ :
        1. Add tuples in  $s_i$  to the hash index if they are not already in it.
        2. At end of  $s_i$  add the tuples in the hash index to the result.



# Other Operations : Set Operations

- E.g., Set operations using hashing:
  1. as before partition  $r$  and  $s$ ,
  2. as before, process each partition  $i$  as follows
    1. build a hash index on  $r_i$
    2. Process  $s_i$  as follows
      - $r \cap s$ :
        1. output tuples in  $s_i$  to the result if they are already there in the hash index
      - $r - s$ :
        1. for each tuple in  $s_i$ , if it is there in the hash index, delete it from the index.
        2. At the end of  $s_i$  add the remaining tuples in the hash index to the result.



# Answering Keyword Queries

- Indices mapping keywords to documents
  - For each keyword, store a sorted list of document IDs that contain the keyword
    - Commonly referred to as an **inverted index**
    - E.g.,: database: d1, d4, d11, d45, d77, d123  
distributed: d4, d8, d11, d56, d77, d121, d333
  - To answer a query with several keywords, compute the intersection of lists corresponding to those keywords
- To support ranking, inverted lists store extra information
  - “**Term frequency**” of the keyword in the document
  - “**Inverse document frequency**” of the keyword
  - **Page rank** of the document/web page



# Other Operations : Outer Join

- **Outer join** can be computed either as
  - A join followed by the addition of null-padded non-participating tuples.
  - by modifying the join algorithms.
- Modifying merge join to compute  $r \bowtie s$ 
  - In  $r \bowtie s$ , nonparticipating tuples are those in  $r - \Pi_R(r \bowtie s)$
  - Modify merge-join to compute  $r \bowtie s$ :
    - During merging, for every tuple  $t_r$  from  $r$  that does not match any tuple in  $s$ , output  $t_r$  padded with nulls.
  - Right outer-join and full outer-join can be computed similarly.



# Other Operations : Outer Join

- Modifying hash join to compute  $r \bowtie s$ 
  - If  $r$  is probe relation, output non-matching  $r$  tuples padded with nulls
  - If  $r$  is the build relation when probing keep track of which  $r$  tuples matched  $s$  tuples. At the end of  $s$ , output non-matched  $r$  tuples padded with nulls





# Evaluation of Expressions

- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree are:
  - **Materialization**: generate results of an expression whose inputs are relations or are already computed, **materialize** (store) it on disk. Repeat.
  - **Pipelining**: pass on tuples to parent operations even as an operation is being executed
- We study the above alternatives in more detail in the following

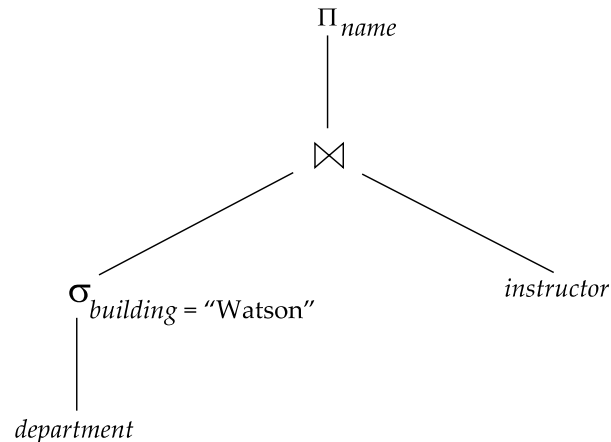


# Materialization

- **Materialized evaluation:** evaluate one operation at a time, starting at the lowest level. Use intermediate results materialized into temporary relations to evaluate next-level operations.
- E.g., in the figure below, compute and store

$$\sigma_{building="Watson"}(department)$$

then compute the store its join with *instructor*, and finally compute the projection on *name*.





# Materialization (Cont.)

- Materialized evaluation is always applicable
- Cost of writing results to disk and reading them back can be quite high
  - Our cost formulas for operations ignore the cost of writing results to disk, so
    - Overall cost = Sum of costs of individual operations + cost of writing intermediate results to disk
- **Double buffering**: use two output buffers for each operation, when one is full write it to disk while the other is getting filled
  - Allows overlap of disk writes with computation and reduces execution time



# Pipelining

- **Pipelined evaluation:** evaluate several operations simultaneously, passing the results of one operation on to the next.

- E.g., in the previous expression tree, don't store the result of

$$\sigma_{building = \text{"Watson"}}(department)$$

- instead, pass tuples directly to the join. Similarly, don't store the result of join, pass tuples directly to projection.
- Much cheaper than materialization: no need to store a temporary relation to disk.
- Pipelining may not always be possible – e.g., sort, hash-join.
- For pipelining to be effective, use evaluation algorithms that generate output tuples even as tuples are received for inputs to the operation.
- Pipelines can be executed in two ways: **demand-driven**, and **producer-driven**



# Pipelining (Cont.)

- In **demand-driven** or **lazy** evaluation
  - System repeatedly requests the next tuple from the top-level operation
  - Each operation requests the next tuple from children operations as required, in order to output its next tuple
  - In between calls, the operation has to maintain “**state**” so it knows what to return next
- In **producer-driven** or **eager** pipelining
  - Operators produce tuples eagerly and pass them up to their parents
    - Buffer maintained between operators, the child puts tuples in the buffer, parent removes tuples from the buffer
    - If the buffer is full, the child waits till there is space in the buffer, and then generates more tuples
  - System schedules operations that have space in the output buffer and can process more input tuples
- Alternative names: **pull** and **push** models of pipelining



# Pipelining (Cont.)

- Implementation of demand-driven pipelining
  - Each operation is implemented as an **iterator** implementing the following operations
    - **open()**
      - E.g., file scan: initialize file scan
        - state: pointer to the beginning of the file
      - E.g., merge join: sort relations;
        - state: pointers to the beginning of sorted relations
    - **next()**
      - E.g., for file scan: Output next tuple, and advance and store file pointer
      - E.g., for merge join: continue with the merge from the earlier state till the next output tuple is found. Save pointers as iterator state.
    - **close()**



# End of Chapter 10