

Chapter 15

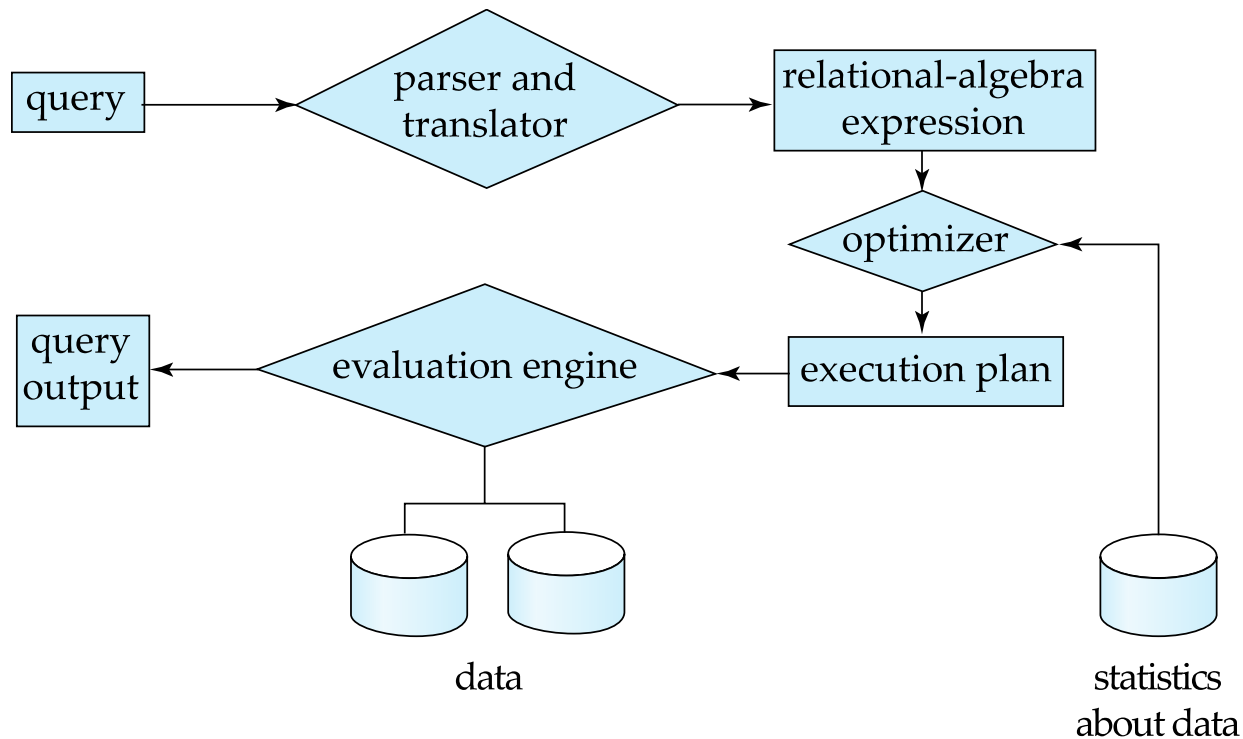
Query Processing

Chapter 15: 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.
 - Parser checks syntax, verifies 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 expressions
 - 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 detailed evaluation strategy is called an **evaluation-plan**. E.g.,:
 - Use an index on *salary* to find instructors with salary < 75000,
 - Or perform complete relation scan and discard instructors with salary ≥ 75000

Basic Steps: Optimization (Cont.)

- **Query Optimization:** Amongst all equivalent evaluation plans choose the one with lowest 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 Chapter 16
 - We study how to optimize queries, that is, how to find an evaluation plan with 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 query, or
 - total **resource consumption**
- We use total resource consumption as cost metric
 - Response time 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 distributed 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 write cost is same as read cost
 - 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 the storage where data is stored
 - High end magnetic disk: $t_S = 4$ msec and $t_T = 0.1$ msec with 4 KB blocks
 - SSD: $t_S = 20\text{-}90$ microsec and $t_T = 2\text{-}10$ microsec for 4KB blocks

Measures of Query Cost (Cont.)

- Required data may be buffer resident already, avoiding disk I/O
 - But hard to take into account 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 buffer and only the minimum amount of memory needed for the operation is available
 - But more optimistic estimates are used in practice

Selection (σ) 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 number of blocks containing records from relation r
 - If selection is on a key attribute, can stop on finding 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 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 nonkey)** 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 first tuple $\geq v$ and scan relation sequentially from there
 - For $\sigma_{A \leq v}(r)$, just **scan relation** sequentially till first tuple $> v$ without using index
- **A6 (secondary index, comparison)**.
 - For $\sigma_{A \geq v}(r)$, **use index** to find first index entry $\geq v$ and scan index sequentially from there, to find pointers to records.
 - For $\sigma_{A \leq v}(r)$, just **scan leaf pages of index** finding pointers to records, till first entry $> v$
 - In either case, retrieving 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 θ_i and algorithms A1 through A7 that results in the least cost for $\sigma_{\theta_i}(r)$.
 - Test other conditions on tuple after fetching it into memory buffer.
- **A8 (conjunctive selection using composite index).**
 - Use appropriate composite (multiple-key) index if available.
- **A9 (conjunctive selection by intersection of identifiers).**
 - Use corresponding index for each condition, and take intersection of all the obtained sets of record pointers. Then fetch records from file.
 - If some conditions do not have appropriate indices, apply 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).**
 - If *all* conditions have available indices, use corresponding index for each condition and take union of all the obtained sets of record pointers. Then fetch records from file.
 - Otherwise, use linear scan.
- **Negation:** $\sigma_{\neg\theta}(r)$
 - Use linear scan on file
 - If very few records satisfy $\neg\theta$, and an index is applicable to θ
 - Find satisfying records using index and fetch from file

‘Bitmap Index Scan’ of PostgreSQL

- A **hybrid** approach of ‘linear scan’ + ‘secondary index scan’
(*This is NOT a scan using the conventional bitmap index that was in Ch. 14.*)
 - Bridges gap between secondary index scan and linear file scan when number of matching records is not known before execution (☞ A6)
 - **Bitmap with 1 bit per page in relation**
 - Steps:
 - Index scan used to find record ids, and set bit of corresponding page in bitmap
 - Linear file scan fetching only pages with bit set to 1
 - Performance
 - Similar to index scan when only a few bits are set
 - Similar to linear file scan when most bits are set
 - Never behaves very badly compared to best alternative

Sorting

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

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).....*

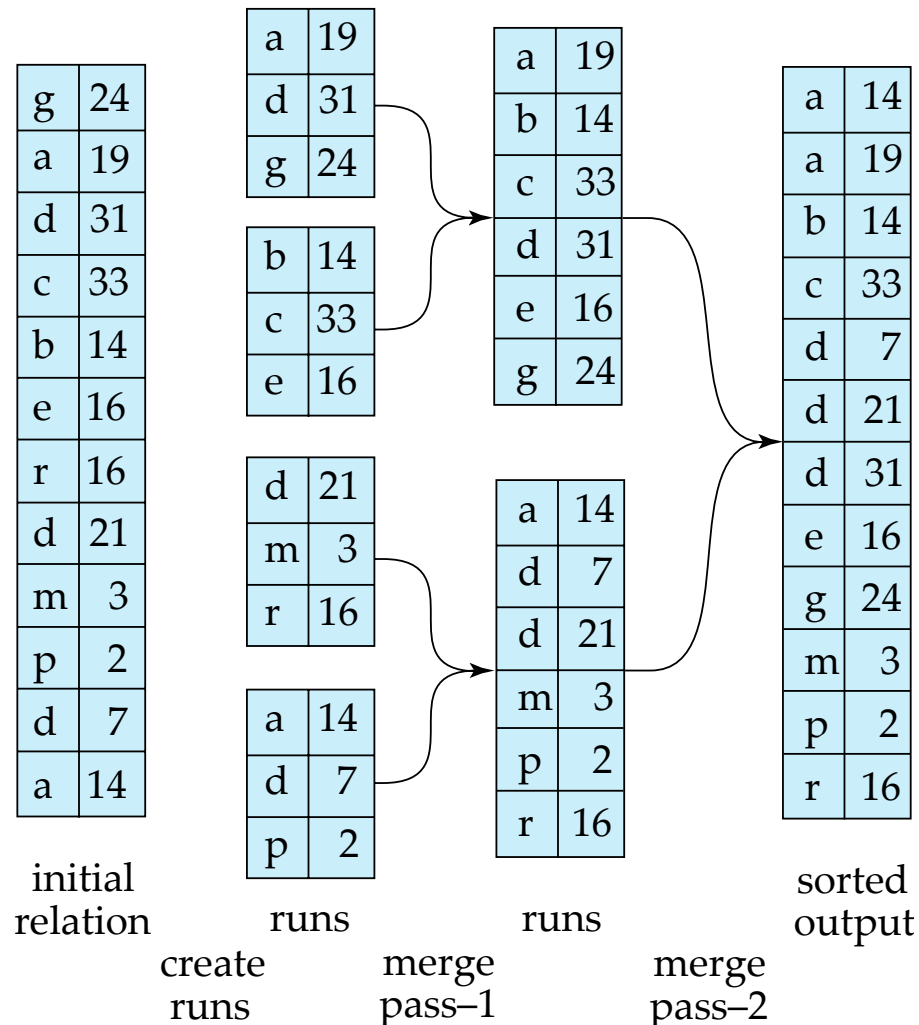
External Sort-Merge (Cont.)

2. **Merge the runs (N-way merge).** (Here, we assume $N = M-1$ for simplicity.)
 1. Use N blocks of memory to buffer input runs, and 1 block to buffer 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 next 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 10 times the size of the initial runs
 - Repeated passes are performed till all runs have been merged into one.

Example: External Sorting Using Sort-Merge



M=3, N=2
b_r=12, b_b=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 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 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
$$n_r * b_s + b_r \text{ block transfers, plus } n_r + b_r \text{ seeks}$$
- 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 worst case memory availability 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-loop join 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
```

Block Nested-Loop Join (Cont.)

- Worst case estimate: $b_r * b_s + b_r$ block transfers + $2 * b_r$ seeks
 - Each block in the inner relation s is read once for each $block$ in the outer relation
- Best case estimate: $b_r + b_s$ block transfers + 2 seeks
 - Improvements to nested loop and block nested loop algorithms:
 - In block nested-loop, use $M - 2$ disk blocks as blocking unit for outer relations, where M = memory size in blocks; use remaining two blocks to buffer inner relation and output
 - Cost = $\lceil b_r / (M-2) \rceil * b_s + b_r$ block transfers + $2 \lceil b_r / (M-2) \rceil$ seeks
- If equi-join attribute forms a key of the inner relation, stop the inner loop on the first match
- Scan the inner loop forward and backward alternately, to make use of the blocks remaining in buffer (with LRU replacement)
- Use index on the inner relation if available (next slide)

Indexed Nested-Loop Join

- Index lookups can replace file scans, if
 - join is an equi-join 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 .
- Cost of the indexed nested-loop join:
 - $(b_r + n_r * c)$ block transfers and seeks
 - Worst case: buffer has space for only one page of r , and, for each tuple in r , we perform an index lookup on s .
 - Where c is the cost (seek + transfer) of traversing index and fetching all matching s tuples for one tuple of r
 - c can be estimated as 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.

Example of Nested-Loop Join Costs

- Compute $student \bowtie takes$, with $student$ as the outer relation.
- Let $takes$ have a primary B⁺-tree index on the attribute ID , which contains 20 entries in each index node.
- Since $takes$ has 10,000 tuples, the height of the tree is 4, and one more access is needed to find the actual data
- $student$ has 5000 tuples
- Cost of block nested-loop join
 - $100 * 400 + 100 = 40,100$ block transfers and $2 * 100 = 200$ seeks
 - assuming worst case memory
 - may be significantly less with more memory
 - $(100 * 200 + 100)$ block transfers if considering ID is ordered
- Cost of indexed nested-loop join
 - $100 + 5000 * 5 = 25,100$ block transfers and seeks.
 - CPU cost likely to be less than that for block nested-loop join

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 handling of duplicate values in join attribute — every pair with same value on join attribute must be matched
 3. Detailed algorithm in book

	<i>a1</i>	<i>a2</i>
<i>pr</i> →	a	3
	b	1
	d	8
	d	13
	f	7
	m	5
	q	6
	<i>r</i>	

	<i>a1</i>	<i>a3</i>
<i>ps</i> →	a	A
	b	G
	c	L
	d	N
	m	B
	<i>s</i>	

Merge Join (Cont.)

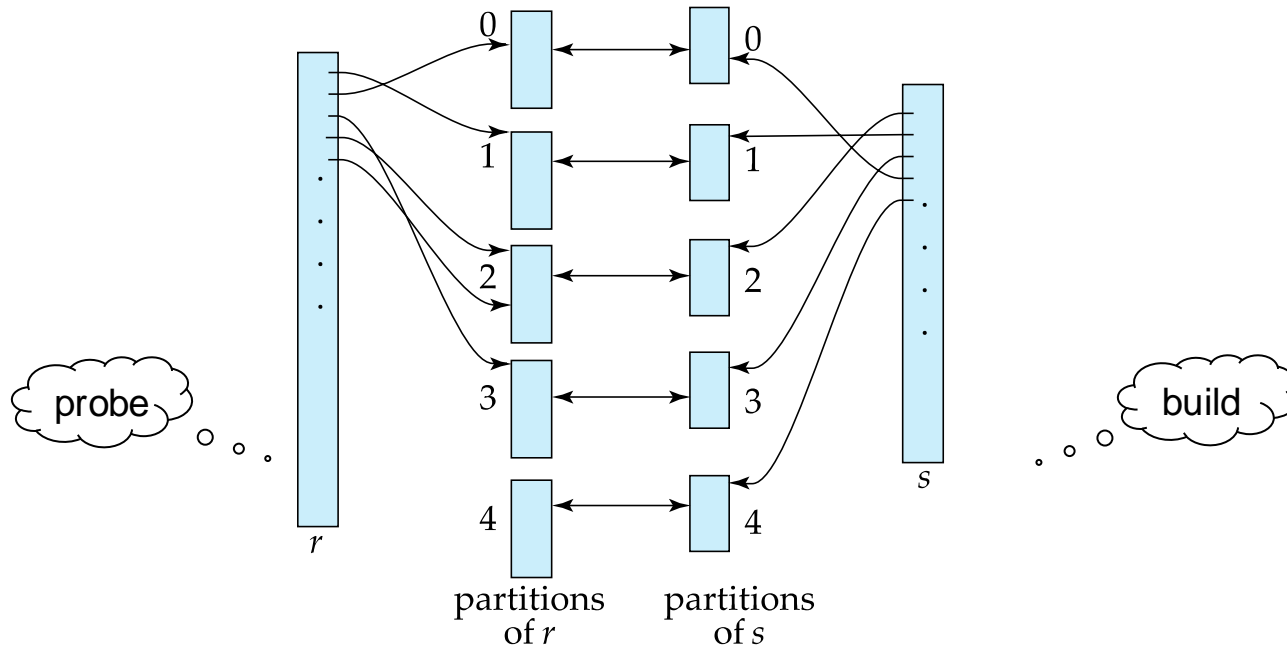
- Can be used only for equi-joins and natural joins
- Each block needs to be read only once (assuming all tuples for any given value of the join attributes fit in memory)
- Thus the cost of merge join is:
$$b_r + b_s \text{ block transfers} + \lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil \text{ seeks}$$

(+ the cost of sorting if relations are unsorted)
- **hybrid merge-join:** If one relation is sorted, and the other has a secondary B⁺-tree index on the join attribute
 - Merge the sorted relation with the leaf entries of the B⁺-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 equi-joins 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 15.10 r_i is denoted as H_{ri} , s_i is denoted as H_{si} and n is denoted as n_h .

Hash Join (Cont.)



- r tuples in r_i need only to be compared with s tuples in s_i . Need not be compared with s tuples in any other partition, since:
 - An r tuple and an s tuple that satisfy the join condition will have the same value for the join attributes.
 - If that value is hashed to some value i , the r tuple has to be in r_i and the s tuple in s_i .

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 $r \bowtie_{\theta_i} s$:

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

Joins - Etc.

- Joins over **spatial data**
 - No simple sort order for spatial joins
 - Indexed nested loops join with spatial indices
 - R-trees, quad-trees, k-d-B-trees
- **Outer join** can be computed via either of
 1. A join followed by addition of null-padded non-participating tuples
 2. Modifying the (block-nested-loop) join algorithm (except for full-outer joins)

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.
 - Performing projection on each tuple followed by duplicate elimination.
- **Set operations** (\cup , \cap and $-$): can either use variant of merge-join after sorting, or variant of hash-join.

Other Operations : 'group by'

- **'group by'** 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 on 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 count, add up the partial aggregates
 - For *avg*, keep sum and count, and divide sum by count at the end

distributive

algebraic

Answering Keyword Queries

- Indices mapping keywords to documents
 - For each keyword, store 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 intersection of lists corresponding to those keywords
- To support ranking, inverted lists store extra information
 - **TF (term frequency), IDF (inverse document frequency)**
 - **PageRank** of the document/web page

$$TF(t, d) = \frac{(\text{Number of occurrences of term } t \text{ in document } d)}{(\text{Total number of terms in the document } d)}$$

$$IDF(t, D) = \log_e \frac{(\text{Total number of documents in the corpus})}{(\text{Number of documents with term } t \text{ in them})}$$

$$TF-IDF(t, d, D) = TF(t, d) \times IDF(t, D)$$

Evaluation of Expressions

- So far: we have seen algorithms for individual operations
- Alternatives for evaluating an entire expression tree
 - **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 above alternatives in more detail.

Notice:

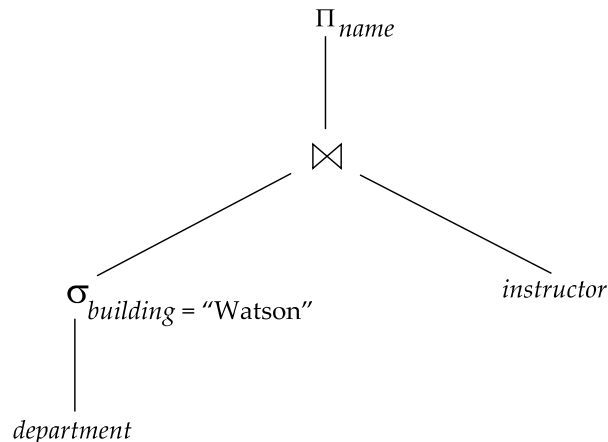
- The term '*materialize*' in the PostgreSQL's query plan denotes '*storing data in the main memory*' e.g. storing the inner node of a nested loop join.
- The term '*pipeline*' in the PostgreSQL is a feature specifically designed to enhance the efficiency of data transfer between the database server and client applications. When pipeline mode is activated, the database server sends *multiple result rows to the client application without waiting for a request for each row*. This parallel processing approach significantly reduces the overall latency and network overhead involved in data transfer.

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 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 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 previous expression tree, don't store result of

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

- instead, pass tuples directly to the join.. Similarly, don't store 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 next tuple from top level operation
 - Each operation requests next tuple from children operations as required, in order to output its next tuple
 - In between calls, 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, child puts tuples in buffer, parent removes tuples from buffer
 - if buffer is full, child waits till there is space in the buffer, and then generates more tuples
 - System schedules operations that have space in output buffer and can process more input tuples
- Alternative name: **pull** and **push** models of pipelining

Continuous-Stream Data

- **Data streams**

- Data entering database in a continuous manner
- E.g., Sensor networks, user clicks, ...

- **Continuous queries**

- Results get updated as streaming data enters the database
- Aggregation on windows is often used
 - E.g., **tumbling windows** divide time into units, e.g., hours, minutes (↔ **sliding window**)

End of Chapter 15