# Chapter 12: Query Processing
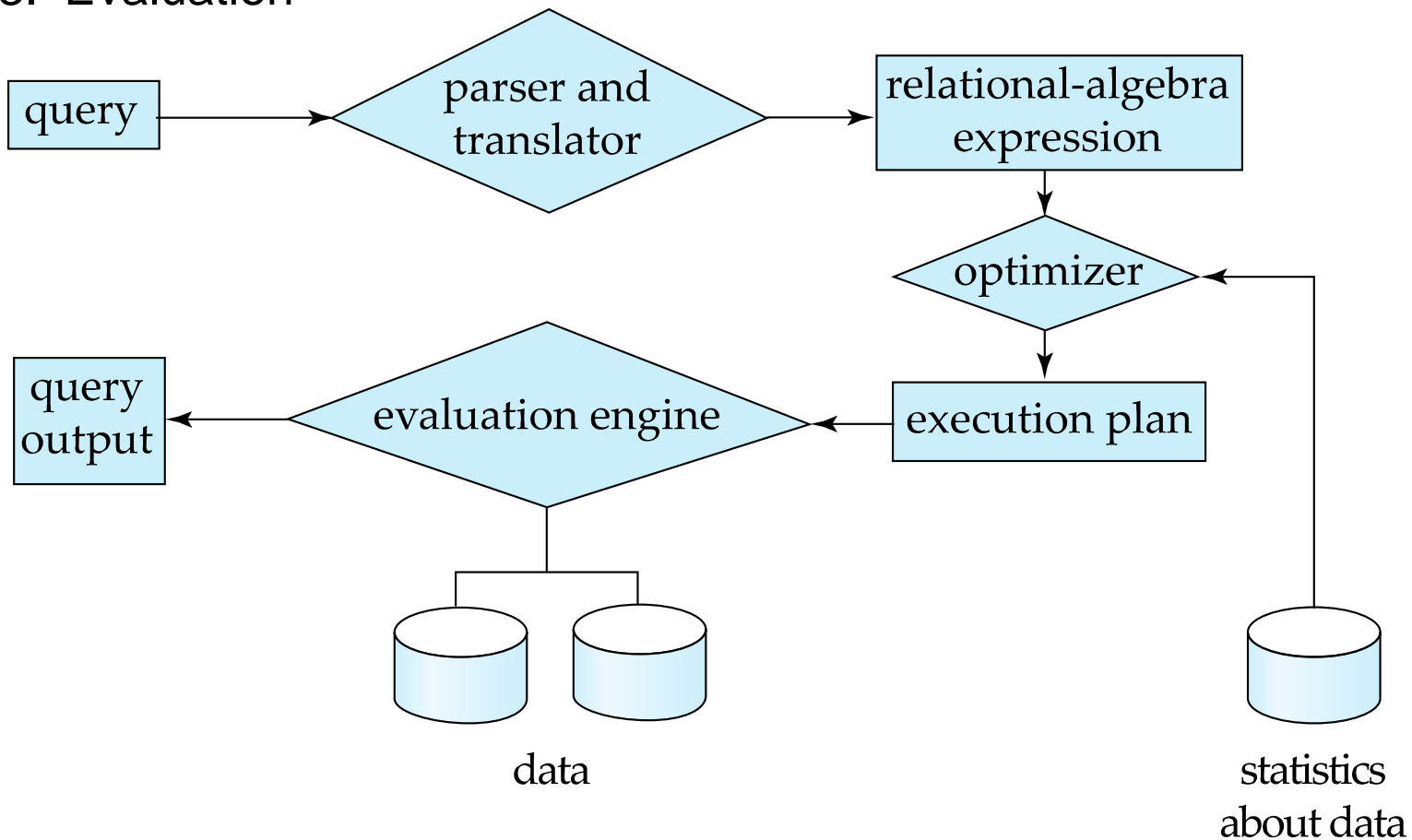
**Database System Concepts, 6th Ed.**

# Chapter 12:  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

query → parser and translator → relational-algebra expression → optimizer → execution plan → evaluation engine → query output

evaluation engine → data

optimizer ← statistics about data

# 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

- Annotated expression specifying detailed evaluation strategy is called an **evaluation-plan**.

  - E.g., can use an index on *salary* to find instructors with salary < 75000,

  - or can perform 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 lowest cost.

- Cost is estimated using statistical information from the database catalog

  ▸ e.g. number of tuples in each relation, size of tuples, etc.

# Measures of Query Cost

- Cost is generally measured as total elapsed time for answering query
  - Many factors contribute to time cost
    - *disk accesses, CPU*, or even network *communication*
- Typically disk access is the predominant cost, and is also relatively easy to estimate.   Measured by taking into account
  - Number of seeks            * average-seek-cost
  - Number of blocks read     * average-block-read-cost
  - Number of blocks written * average-block-write-cost
    - Cost to write a block is greater than cost to read a block
      - data is read back after being written to ensure that the write was successful

# Measures of Query Cost (Cont.)

- 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

  - $t_S$ – time for one seek

  - Cost for b block transfers plus S seeks
    $$b * t_T + S * t_S$$

- We ignore CPU costs for simplicity

  - Real systems do take CPU cost into account

- We do not include cost to writing output to disk in our cost formulae

- Required data may be buffer resident already, avoiding disk I/O

  - But hard to take into account for cost estimation, so we assume worst case

# Selection Operation

- **File scan**

- Algorithm **A1** (**linear search**).  Scan each file block and test all records to see whether they satisfy the selection condition.

  - 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** (**primary index, equality on key**).  Retrieve a single record that satisfies the corresponding equality condition
  - $Cost = (h_i + 1) * (t_T + t_S)$
- **A3** (**primary 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 nonkey**)*.*

  - 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!

# Implementation of Complex Selections

- **Conjunction:** $\sigma_{\theta 1 \wedge \theta 2 \wedge \ldots \theta n}(r)$

- **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 tuple after fetching it into memory buffer.

- **A8** (**conjunctive selection using composite index**).

  - Use appropriate composite (multiple-key) index if available.

# Quiz Time

**Quiz Q1:**

Given a choice between using a secondary index and a file scan to answer a query select * from r where P, where P is a predicate,

(1) It is always better to use the secondary index

(2) It is always better to use file scan

(3) It depends on the number of records fetched by the query
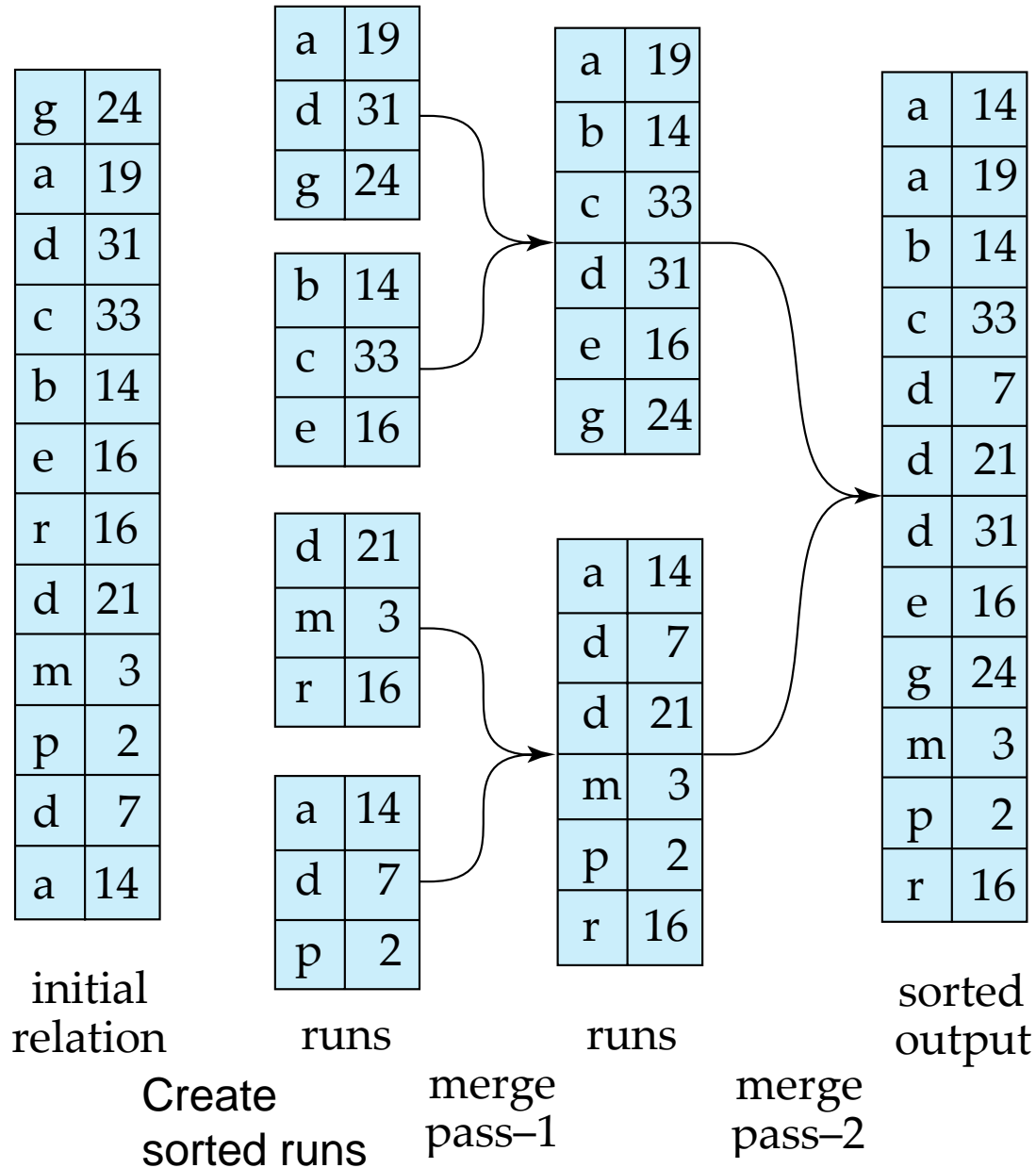
(4) None of the above

# 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, techniques like quicksort can be used. For relations that don't fit in memory, **external sort-merge** is a good choice.

| | |
|---|---|
| g | 24 |
| a | 19 |
| d | 31 |
| c | 33 |
| b | 14 |
| e | 16 |
| r | 16 |
| d | 21 |
| m | 3 |
| p | 2 |
| d | 7 |
| a | 14 |

initial
relation

| | |
|---|---|
| a | 19 |
| d | 31 |
| g | 24 |

| | |
|---|---|
| b | 14 |
| c | 33 |
| e | 16 |

| | |
|---|---|
| d | 21 |
| m | 3 |
| r | 16 |

| | |
|---|---|
| a | 14 |
| d | 7 |
| p | 2 |

runs

| | |
|---|---|
| a | 19 |
| b | 14 |
| c | 33 |
| d | 31 |
| e | 16 |
| g | 24 |

| | |
|---|---|
| a | 14 |
| d | 7 |
| d | 21 |
| m | 3 |
| p | 2 |
| r | 16 |

runs

| | |
|---|---|
| a | 14 |
| a | 19 |
| b | 14 |
| c | 33 |
| d | 7 |
| d | 21 |
| d | 31 |
| e | 16 |
| g | 24 |
| m | 3 |
| p | 2 |
| r | 16 |

sorted
output

Create
sorted runs

merge
pass–1

merge
pass–2

# 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 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.
  - Total number of merge passes required: $\lceil \log_{M-1}(b_r/M) \rceil$ where $b_r$ is the number of blocks in relation $r$
- Cost analysis: each pass reads and writes all the records, times number of passes

# 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 $\quad r \bowtie_\theta s$
  **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.

# 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 $t_r$.
  - Output matching pairs of tuples
- Worst case:  buffer has space for only one page of $r$.
- Cost of the join:  $b_r (t_T + t_S) + n_r * c$
  - Where $c$ is the cost of traversing index and fetching all matching $s$ tuples for one tuple or $r$
  - $c$ can be estimated as cost of a single selection on $s$ using the join condition..

# 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

$pr$

| $a1$ | $a2$ |
|------|------|
| a    | 3    |
| b    | 1    |
| d    | 8    |
| d    | 13   |
| f    | 7    |
| m    | 5    |
| q    | 6    |

$r$

$ps$

| $a1$ | $a3$ |
|------|------|
| a    | A    |
| b    | G    |
| c    | L    |
| d    | N    |
| m    | B    |

$s$

# 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:
  $$b_r + b_s \text{ block transfers } + \lceil b_r / b_b \rceil + \lceil b_s / b_b \rceil \text{ seeks}$$
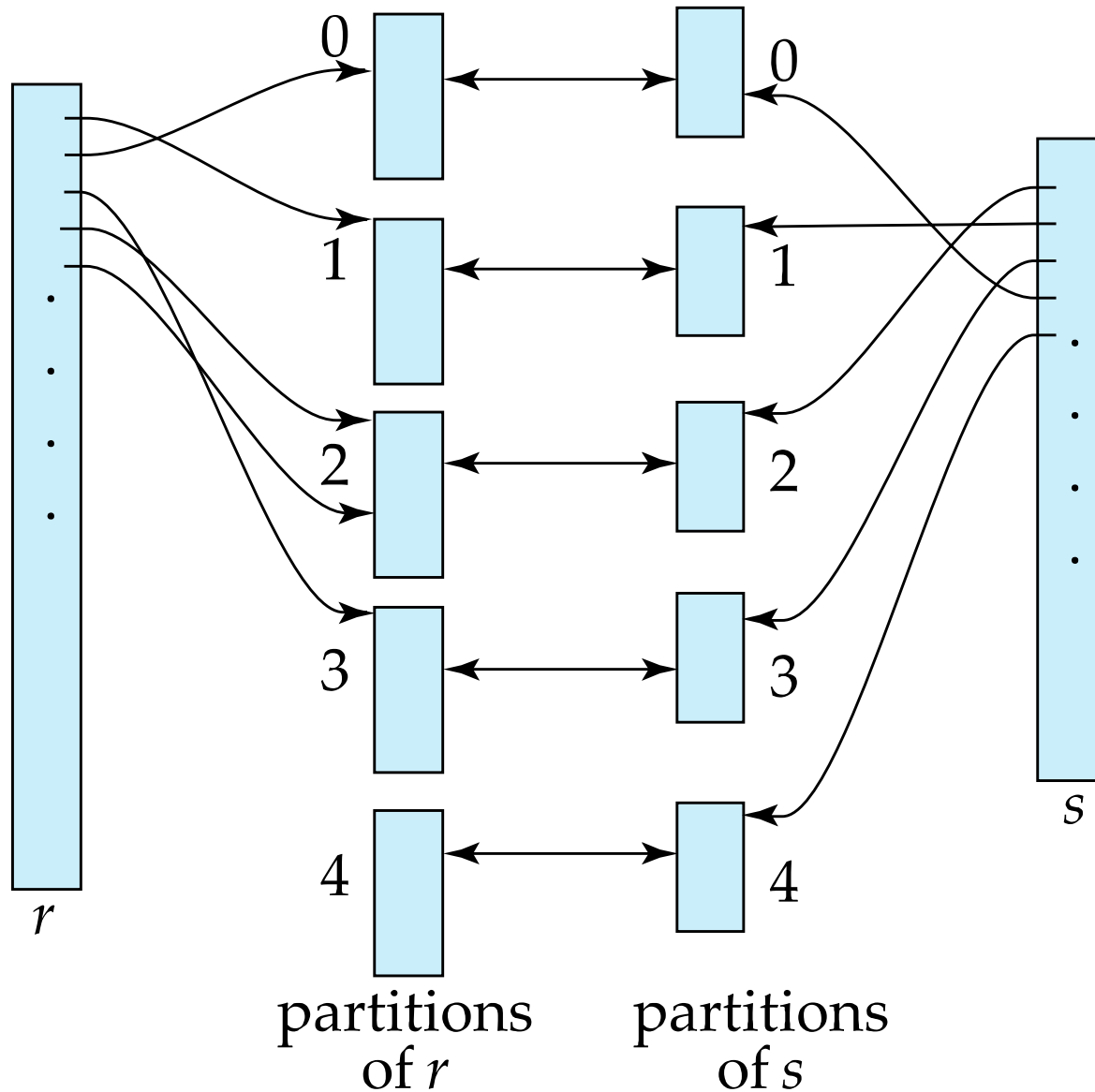  - + the cost of sorting if relations are unsorted.

# 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 $\{0, 1, ..., n\}$, where *JoinAttrs* denotes the common attributes of $r$ and $s$ used in the natural join.

  - $r_0, r_1, . . ., r_n$ denote partitions of $r$ tuples
    - Each tuple $t_r \in r$ is put in partition $r_i$ where $i = h(t_r[JoinAttrs])$.

  - $r_0, r_1. . ., r_n$ denotes partitions of $s$ tuples
    - Each tuple $t_s \in s$ is put in partition $s_i$, where $i = h(t_s[JoinAttrs])$.

# Hash-Join (Cont.)



partitions of $r$

partitions of $s$

$r$

$s$

# 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 $s_i$.

# 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**.

# Hash-Join algorithm (Cont.)

- The value *n* and the hash function *h* is chosen such that each $s_i$ should fit in memory.

  - Typically n is chosen as $\lceil b_S/M \rceil * f$ where f is a "**fudge factor**", typically around 1.2

  - The probe relation partitions $s_i$ need not fit in memory

# Complex Joins

- Join with a conjunctive condition:

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

  - 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

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

# Quiz Time

**Quiz Q2:**

If two relations r(A,B) and s(A,C) are given, sorted on attribute A, then the natural join of r and s can be computed fastest using
(1) nested loops join
(2) indexed nested loops join
(3) merge join
(4) hash join

**Quiz Q3:**

If data is stored on a solid state (flash) disk instead of a hard disk, which of the following join methods will benefit the most:
(1) nested loops join
(2) indexed nested loops join
(3) merge join
(4) hash join
  Hint: which one has the maximum number of random IO operations?

# 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 on each group.

  - *Optimization:* combine tuples in the same group during run generation and intermediate merges, by computing partial aggregate values

# 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 : Outer Join

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

# 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

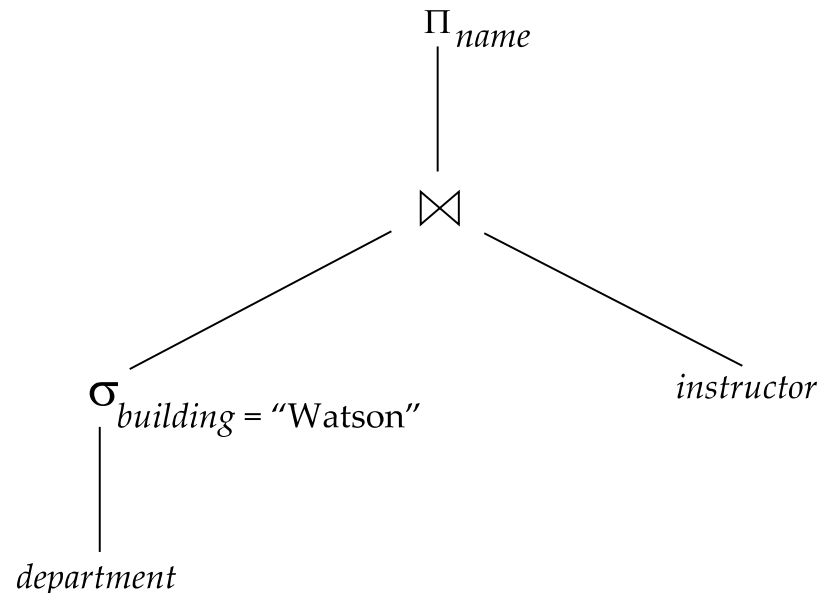- We study above alternatives in more detail

# 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.*

# 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

# 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 beginning of file
            - – E.g.merge join: sort relations;
                - » state: pointers to 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 merge from earlier state till
              next output tuple is found.  Save pointers as iterator state.
        - ▸ close()

# End of Chapter

# External Merge Sort (Cont.)

- **Cost analysis:**

  - Total number of merge passes required: $\lceil \log_{M-1}(b_r/M) \rceil$.

  - Block transfers for initial run creation as well as in each pass is $2b_r$

    - for final pass, we don't count write cost

      - we ignore 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 total number of block transfers for external sorting:
      $$b_r \, ( \, 2 \lceil \log_{M-1}(b_r / M) \rceil + 1)$$

  - 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

    ▸ Buffer size: $b_b$ (read/write $b_b$ blocks at a time)

    ▸ 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:
      $2 \lceil b_r / M \rceil + \lceil b_r / b_b \rceil (2 \lceil \log_{M-1}(b_r / M) \rceil - 1)$

# 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 end of $s_i$ add remaining tuples in the hash index to the result.

# 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 build relation, when probing keep track of which $r$ tuples matched $s$ tuples. At end of $s_i$ output non-matched $r$ tuples padded with nulls

# Measures of Query Cost (Cont.)

■ 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

    ▸ We often use worst case estimates, assuming only the minimum amount of memory needed for the operation is available