# Lecture 9:
# Query Processing & Optimization

Dr. Kyong-Ha Lee
(kyongha@kisti.re.kr)

# Contents

## Brief overview of this lecture

- Basic theories and principles about Index mechanisms used in database systems
- Not much discussion on implementation or tools, but will be happy to discuss them if there are any questions

| | |
|---|---|
| Contents 1 | Measures of Query Cost |
| Contents 2 | Selection and Sorting |
| Contents 3 | Join Operations |
| Contents 4 | Cost Estimation |
| Contents 5 | Cost-based Optimization |
| Contents 6 | Dynamic Programming for Choosing Plan |
| Contents 5 | Materialized View |

*Disclaimer: these slides are based on the slides created by the authors of Database System Concepts 7th ed. and modified by K.H. Lee.

# Basic Steps in Query Processing

1. Parsing and translation
2. Optimization
3. Evaluation (or exection)

Check syntax, verifies relations

$\sigma_{salary<75000}(\prod_{salary}(instructor))$

query

parser and translator

relational-algebra expression

Query written in SQL
e.g., select salary from instructor
where salary < 75000

Choose a plan with the lowest cost amongst all equivalent plan

$\sigma_{salary<75000}(\prod_{salary}(instructor))$ ?
$\prod_{salary}(\sigma_{salary<75000}(instructor))$ ?

Evaluation(or exection)

optimizer

query output

evaluation engine

execution plan

Annotated expression specifying detailed evaluation strategy

data

statistics about data

# Measures of Query Cost

- *Many factors contribute to time cost*
  - Disk accesses, 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

- We ignore CPU costs for simplicity

- We describe how estimate the cost of each operation
  - Not consider the cost to writing output to disks

# 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 # **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 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

# Selection Operation $\sigma$

- 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 $r$
- **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 = $h_i * (t_T + t_S) + t_S + t_T * b$

# Selection Operation σ

- **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 $> v$; do not use index
- **A6 (clustering 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, retrieve records that are pointed to
  - requires an I/O per record; Linear file scan may be cheaper!

# Implementation of Complex Selections

- **Conjunction:** $\sigma_{\theta_1 \wedge \theta_2 \wedge \ldots \theta_n}(r)$
- **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.
- **A9** (**conjunctive selection by intersection of identifiers**).
  - Requires indices with record pointers.
  - 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.
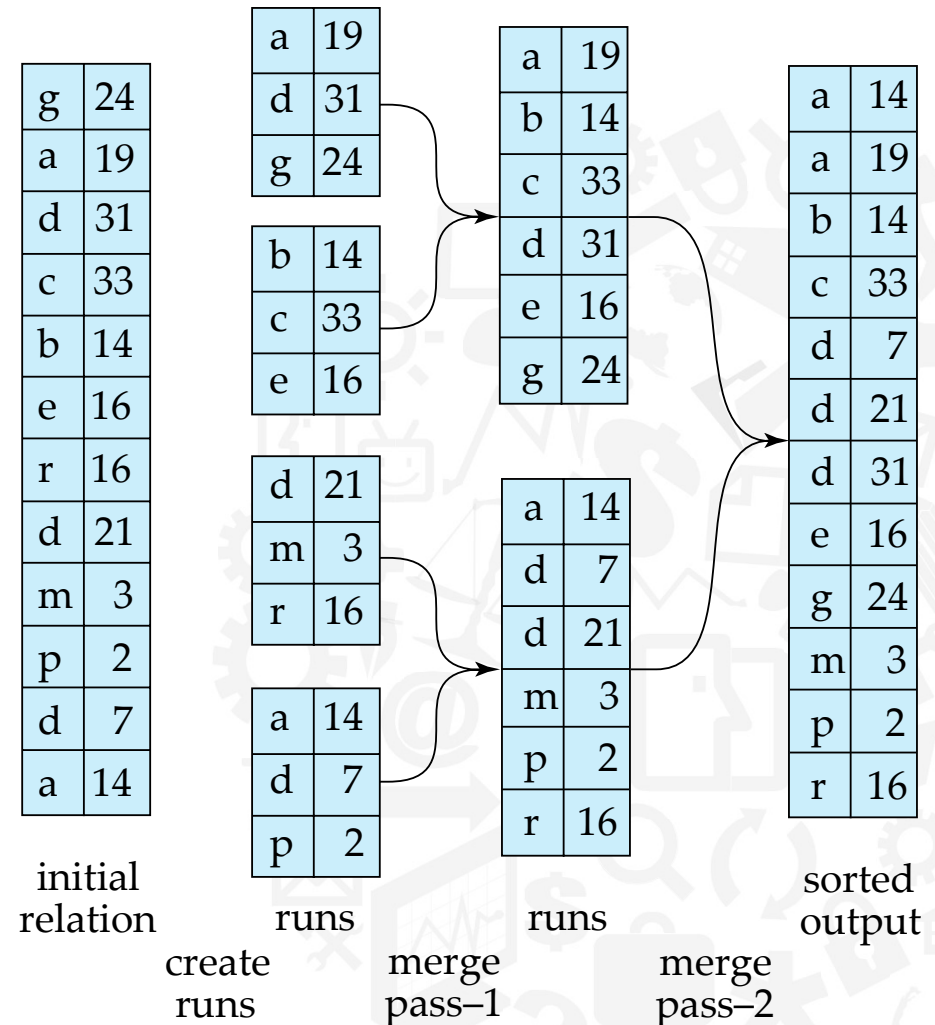
# Algorithms for Complex Selections

- **Disjunction:** $\sigma_{\theta_1 \vee \theta_2 \vee \ldots \theta_n}(r)$.

- **A10** (**disjunctive selection by union of identifiers**).
  - Applicable if *all* conditions have available indices.
    - Otherwise use linear scan.
  - Use corresponding index for each condition, and take union of all the obtained sets of record pointers.
  - Then fetch records from 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. 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

create 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

merge pass–1

| | |
|---|---|
| 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

merge pass–2

11

# 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 (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:

# Cost

- Cost analysis:
  - 1 block per run leads to too many seeks during 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: $\lceil \log_{\lfloor M/bb \rfloor - 1}(b_r/M) \rceil$.
  - 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:
      $$b_r \left( 2 \lceil \log_{\lfloor M/bb \rfloor - 1}(b_r/M) \rceil + 1 \right)$$
    - with substantial seek cost

# Join Operations

- Several different algorithms to implement joins
  - Nested-loop join
  - Block nested-loop join
  - Indexed nested-loop join
  - Merge-join
  - Hash-join
- Choice a join algorithm based on cost estimation
- 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 $t_r \bullet t_s$ 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

  $n_r * b_s + b_r$ block transfers, plus $n_r + b_r$ seeks

  - $n_r$ : # of tuples in relation r   $b_s$ : # of blocks that include records in relation s

- 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-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 $t_r \bullet t_s$ 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: $b_r + b_s$ 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 = $\lceil b_r / (M\text{-}2) \rceil * b_s + b_r$ block transfers + $2 \lceil b_r / (M\text{-}2) \rceil$ seeks
  - If equi-join attribute forms a key or inner relation, stop inner loop on first match
  - Scan inner loop forward and backward alternately, to make use of the blocks remaining in buffer (with LRU replacement)
  - Use index on 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 $t_r$.
- 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 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.
- 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* ⋈ *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 loops join
  - 400*100 + 100 =  40,100 block transfers + 2 * 100 = 200 seeks
    - assuming worst case memory
    - may be significantly less with more memory
- Cost of indexed nested loops join
  - 100 + 5000 * 5 = 25,100  block transfers and seeks.
  - CPU cost likely to be less than that for block nested loops 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

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

$pr$

$r$

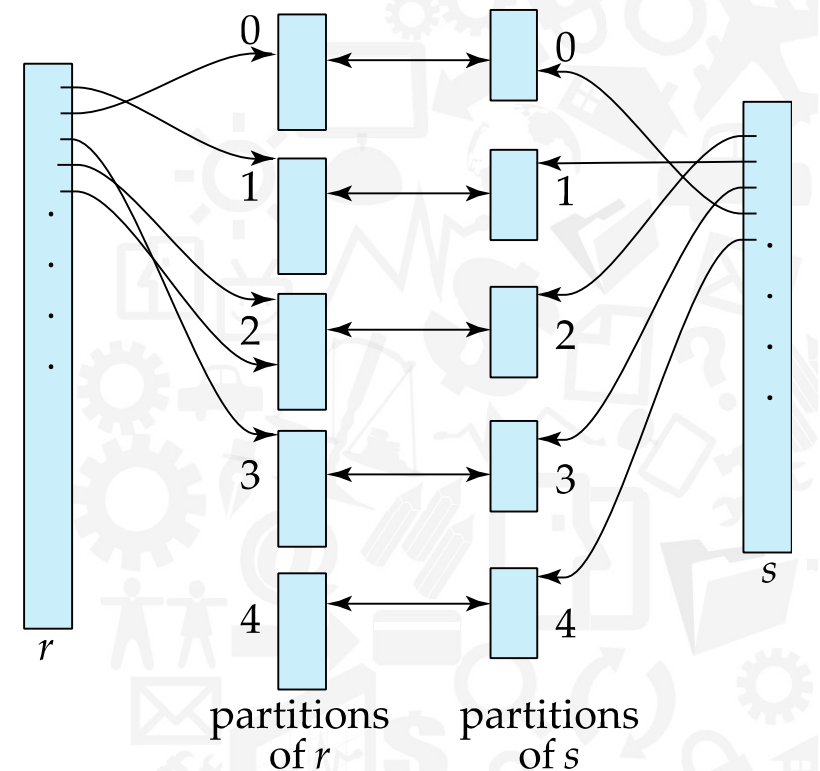| a1 | a3 |
|----|----|
| a  | A  |
| b  | G  |
| c  | L  |
| d  | N  |
| m  | B  |

$ps$

$s$

# 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.

- **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 $\{0, 1, ..., n\}$, 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 $i = h(t_s[JoinAttrs])$



partitions of $r$   partitions of $s$

# 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
- **Recursive partitioning** required if number of partitions *n* is greater than number of pages *M* of memory.
  - instead of partitioning *n* ways, use  *M* – 1 partitions for s
  - Further partition the *M* – 1 partitions using a different hash function
  - Use same partitioning method on *r*
  - Rarely required: e.g., with block size of 4 KB, recursive partitioning not needed for relations of < 1GB with memory size of 2MB, or relations of < 36 GB with memory of 12 MB

# Example of Cost of Hash-Join

- Assume that memory size is 20 blocks

- $b_{instructor}$ = 100 and $b_{teaches}$ = 400.

- *instructor* is to be used as build input.  Partition it into five partitions, each of size 20 blocks.  This partitioning can be done in one pass.

- Similarly, partition *teaches* into five partitions, each of size 80.  This is also done in one pass.

- Therefore total cost, ignoring cost of writing partially filled blocks:

  - $3(100 + 400) = 1500$ block transfers + $2(\lceil 100/3 \rceil + \lceil 400/3 \rceil) = 336$ seeks

- Useful when memory sized are relatively large, and the build input is bigger than memory.
- **Main feature of hybrid hash join:**

  **Keep the first partition of the build relation in memory.**
- E.g. With memory size of 25 blocks, *instructor* can be partitioned into five partitions, each of size 20 blocks.
  - Division of memory:
    - The first partition occupies 20 blocks of memory
    - 1 block is used for input, and 1 block each for buffering the other 4 partitions.
- *teaches* is similarly partitioned into five partitions each of size 80
  - the first is used right away for probing, instead of being written out
- Cost of 3(80 + 320) + 20 +80 = 1300 block transfers for hybrid hash join, instead of 1500 with plain hash-join.
- Hybrid hash-join most useful if $M >> \sqrt{b_s}$

# 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

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

- Join with a disjunctive condition

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

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

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

# 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 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

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

# 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 a **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
  - "**Term frequency**" of the keyword in the document
  - "**Inverse document frequency**" of the keyword
  - **Page rank** of the document/web page

# 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
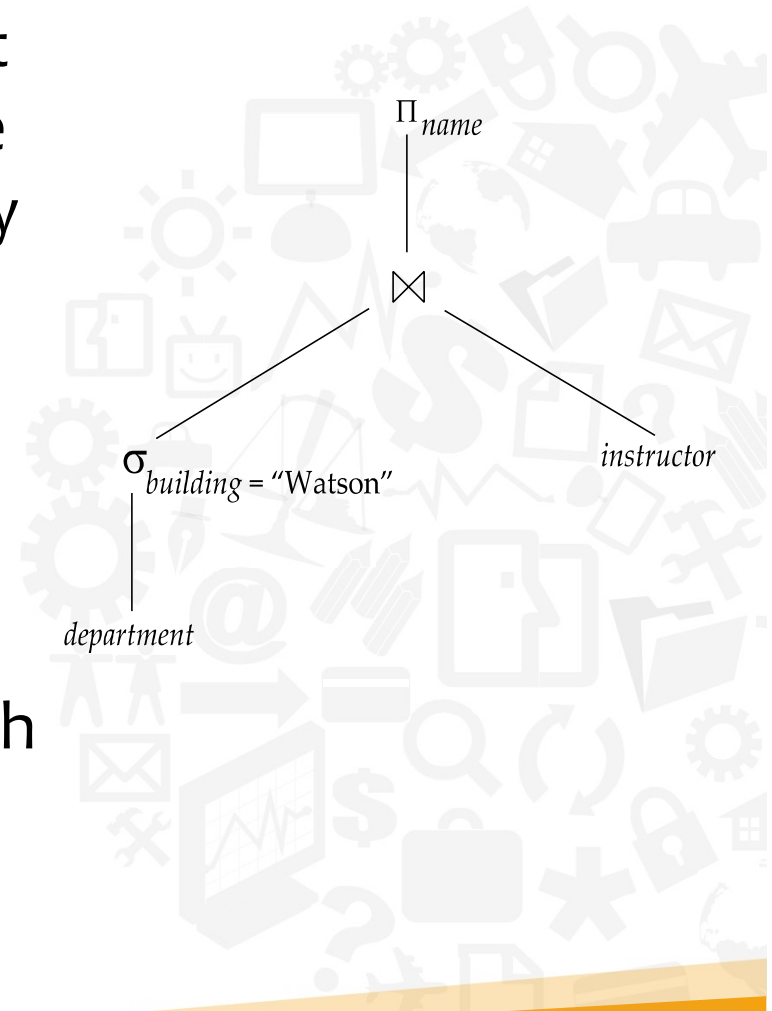
- We study above alternatives in more detail

# 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*

$\Pi_{name}$

$\bowtie$

$\sigma_{building\ =\ "Watson"}$

*instructor*

*department*

# Pipelining

- **Pipelined evaluation**:  evaluate several operations simultaneously, passing the results of one operation on to the next.
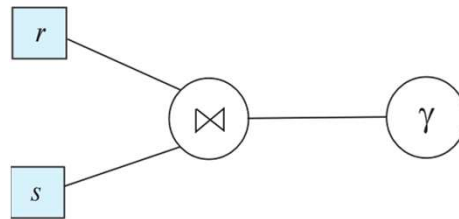- E.g., in previous expression tree, don't store result of

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

  - 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**
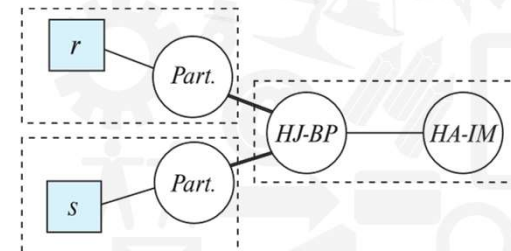
# Blocking Operations

- **Blocking operations:** cannot generate any output until all input is consumed
  - E.g., sorting, aggregation, …
- But can often consume inputs from a pipeline, or produce outputs to a pipeline
- Key idea: blocking operations often have two suboperations
  - E.g., for sort: run generation and merge
  - For hash join: partitioning and build-probe
- Treat them as separate operations



(a) Logical Query

(b) Pipelined Plan

- **Pipeline stages:**
  - All operations in a stage run concurrently
  - A stage can start only after preceding stages have completed execution
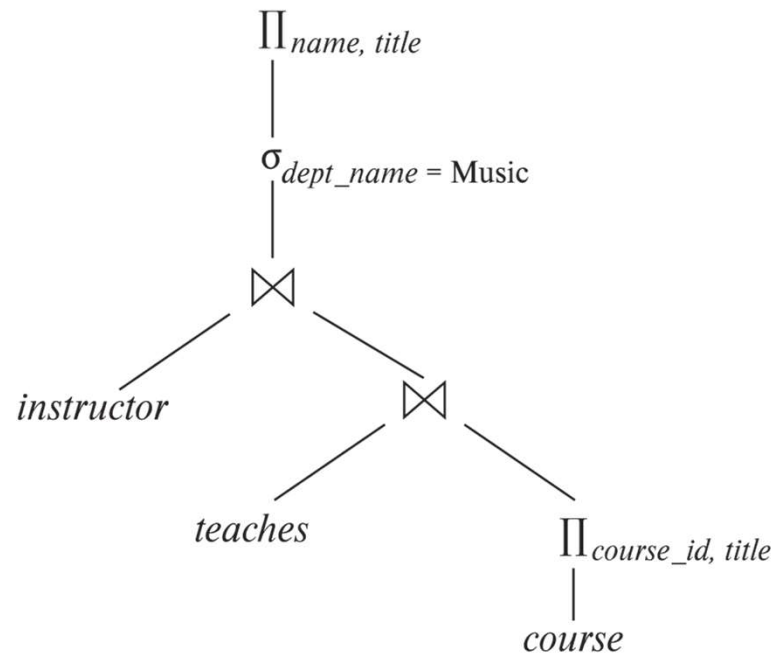
# Pipelining for Continuous-Stream Data

- **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

- Need to use pipelined processing algorithms
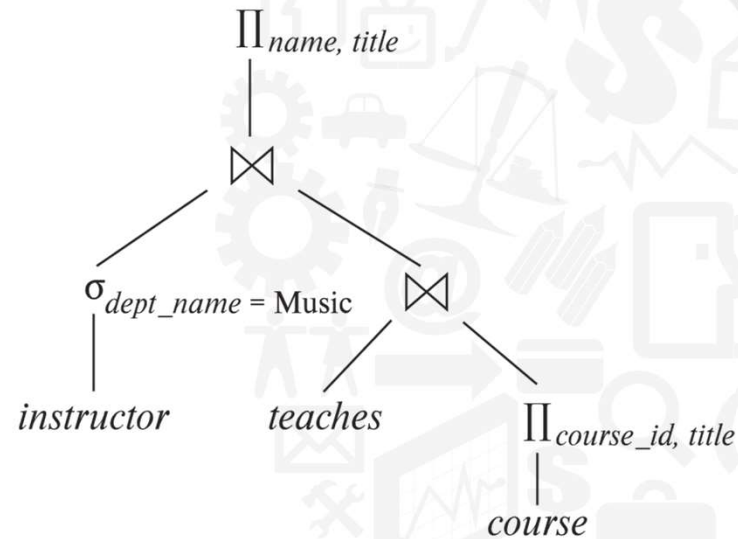  - **Punctuations** used to infer when all data for a window has been received

# Query Evaluation

- Alternative ways of evaluating a given query
  - Equivalent expressions
  - Different algorithms for each operation



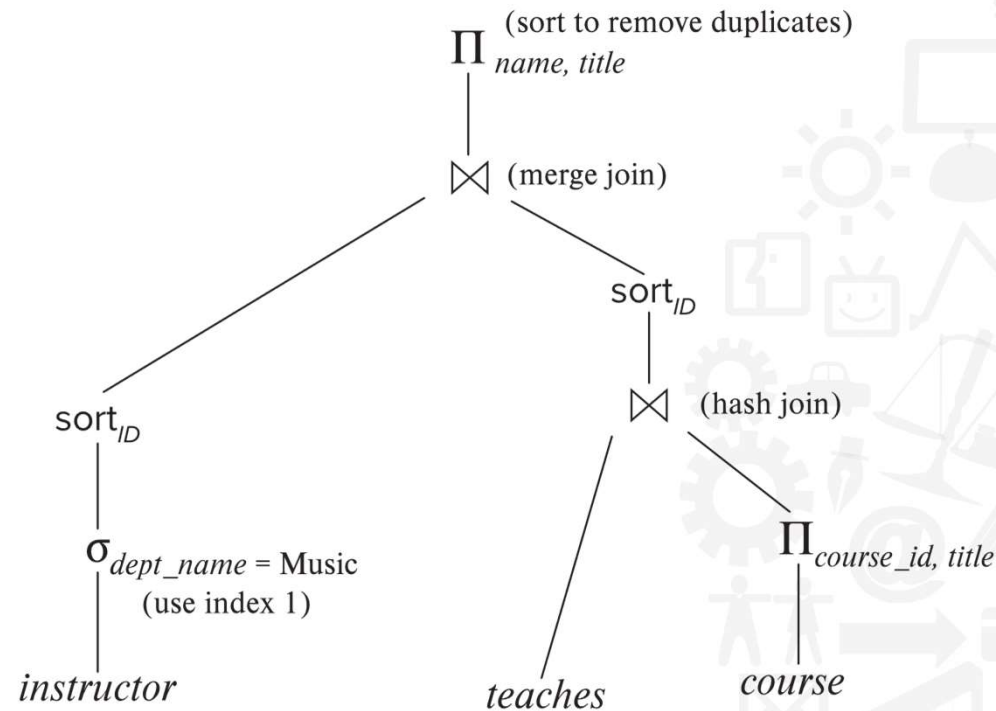(a) Initial expression tree          (b) Transformed expression tree

# Query Evaluation Plan

- An **evaluation plan** defines exactly what algorithm is used for each operation, and how the execution of the operations is coordinated.

$\Pi_{name,\ title}$ (sort to remove duplicates)

$\bowtie$ (merge join)

$sort_{ID}$

$\sigma_{dept\_name}$ = Music
(use index 1)

*instructor*

*teaches*

$sort_{ID}$

$\bowtie$ (hash join)

$\Pi_{course\_id,\ title}$

*course*

- Find out how to view query execution plans on your favorite database

# Viewing Query Evaluation Plans

- Most database support **explain** <query>
  - Displays plan chosen by query optimizer, along with cost estimates
  - Some syntax variations between databases
    - Oracle: **explain plan for** <query> followed by **select * from** table (*dbms_xplan.display*)
    - Maria DB: **explain** {select, update or delete statement}
- Some databases (e.g. MariaDB) support **analyze** <query>
  - Shows actual runtime statistics found by running the query, in addition to showing the plan

# Generating Equivalent Expressions

- Two relational algebra expressions are said to be **equivalent** if the two expressions generate the same set of tuples on every *legal* database instance
  - Note: order of tuples is irrelevant
  - we don't care if they generate different results on databases that violate integrity constraints
- In SQL, inputs and outputs are multisets of tuples
  - Two expressions in the multiset version of the relational algebra are said to be equivalent if the two expressions generate the same multiset of tuples on every legal database instance.
- An **equivalence rule** says that expressions of two forms are equivalent
  - Can replace expression of first form by second, or vice versa

# Equivalence Rules

1. Conjunctive selection operations can be decomposed into a sequence of individual selections.

$$\sigma_{\theta_1 \wedge \theta_2}(E) \equiv \sigma_{\theta_1}(\sigma_{\theta_2}(E))$$

2. Selection operations are commutative.

$$\sigma_{\theta_1}(\sigma_{\theta_2}(E)) \equiv \sigma_{\theta_2}(\sigma_{\theta_1}(E))$$

3. Only the last in a sequence of projection operations is needed, the others can be omitted.

$$\prod_{L_1}(\prod_{L_2}(\dots(\prod_{L_n}(E))\dots)) \equiv \prod_{L_1}(E)$$

where $L_1 \subseteq L_2 \dots \subseteq L_n$

4. Selections can be combined with Cartesian products and theta joins.

   a. $\sigma_{\theta}(E_1 \times E_2) \equiv E_1 \bowtie_{\theta} E_2$

   b. $\sigma_{\theta_1}(E_1 \bowtie_{\theta_2} E_2) \equiv E_1 \bowtie_{\theta_1 \wedge \theta_2} E_2$

# Equivalence Rules(Cont.)

5. Theta-join operations (and natural joins) are commutative.

$$E_1 \bowtie E_2 \equiv E_2 \bowtie E_1$$

6. (a) Natural join operations are associative:

$$(E_1 \bowtie E_2) \bowtie E_3 \equiv E_1 \bowtie (E_2 \bowtie E_3)$$

(b) Theta joins are associative in the following manner:

$$(E_1 \bowtie_{\theta_1} E_2) \bowtie_{\theta_2 \wedge \theta_3} E_3 \equiv E_1 \bowtie_{\theta_1 \wedge \theta_3} (E_2 \bowtie_{\theta_2} E_3)$$

where $\theta_2$ involves attributes from only $E_2$ and $E_3$.

7. The selection operation distributes over the theta join operation under the following two conditions:
(a) When all the attributes in $\theta_0$ involve only the attributes of one of the expressions ($E_1$) being joined.
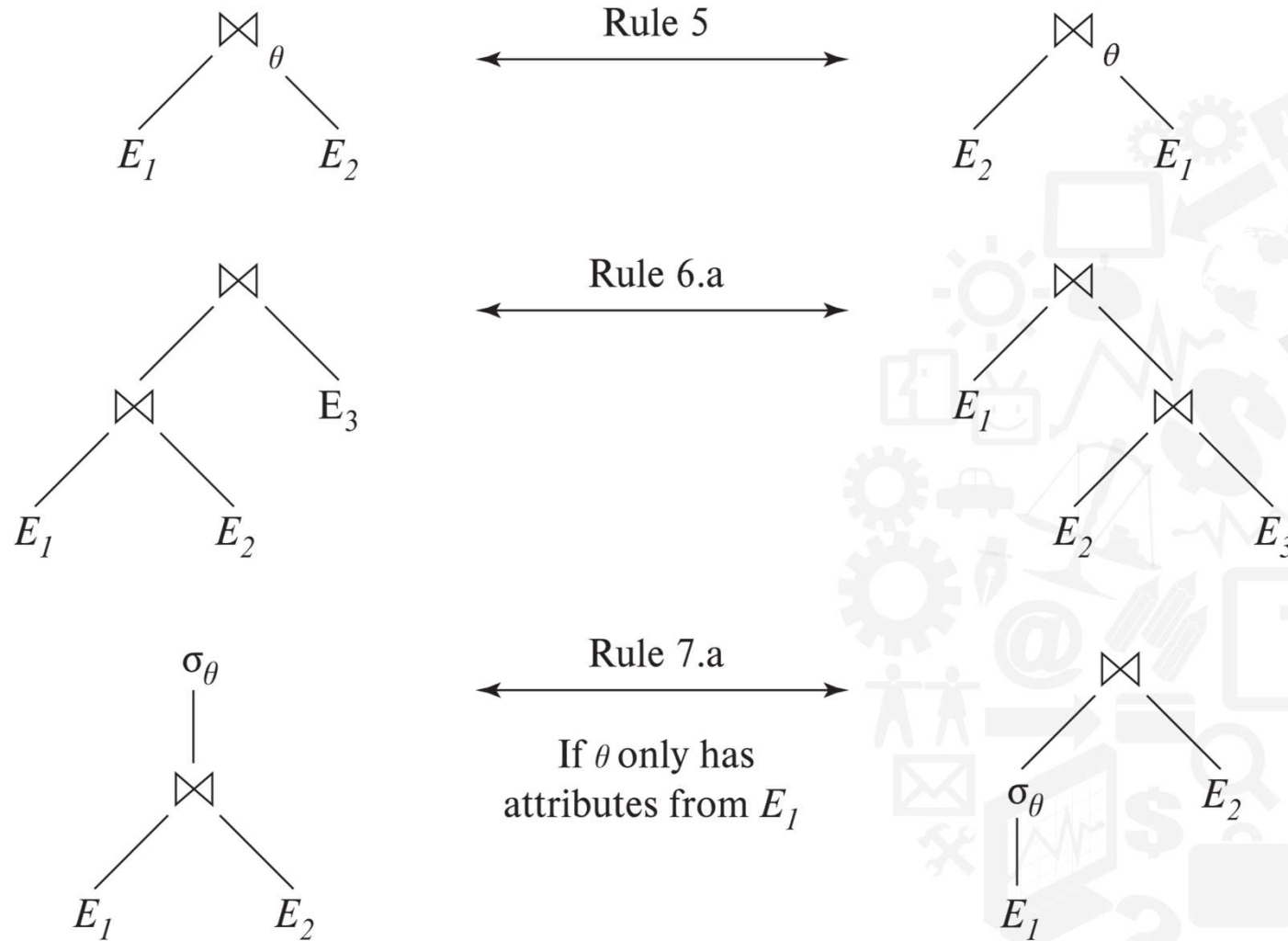
$$\sigma_{\theta_0} (E_1 \bowtie_{\theta} E_2) \equiv (\sigma_{\theta_0}(E_1)) \bowtie_{\theta} E_2$$

(b) When $\theta_1$ involves only the attributes of $E_1$ and $\theta_2$ involves only the attributes of $E_2$.

$$\sigma_{\theta_1 \wedge \theta_2} (E_1 \bowtie_{\theta} E_2) \equiv (\sigma_{\theta_1}(E_1)) \bowtie_{\theta} (\sigma_{\theta_2}(E_2))$$

# Pictorial Depiction of Equivalence Rules



$E_1 \bowtie_\theta E_2$     Rule 5     $E_2 \bowtie_\theta E_1$

$(E_1 \bowtie E_2) \bowtie E_3$     Rule 6.a     $E_1 \bowtie (E_2 \bowtie E_3)$

$\sigma_\theta(E_1 \bowtie E_2)$     Rule 7.a     $(\sigma_\theta(E_1)) \bowtie E_2$

If $\theta$ only has attributes from $E_1$

# Equivalence Rules(Cont.)

8. The projection operation distributes over the theta join operation as follows:

(a) if $\theta$ involves only attributes from $L_1 \cup L_2$:

$$\prod_{L_1 \cup L_2}(E_1 \bowtie_\theta E_2) \equiv \prod_{L_1}(E_1) \bowtie_\theta \prod_{L_2}(E_2)$$

(b) In general, consider a join $E_1 \bowtie_\theta E_2$.

- Let $L_1$ and $L_2$ be sets of attributes from $E_1$ and $E_2$, respectively.

- Let $L_3$ be attributes of $E_1$ that are involved in join condition $\theta$, but are not in $L_1 \cup L_2$, and

- let $L_4$ be attributes of $E_2$ that are involved in join condition $\theta$, but are not in $L_1 \cup L_2$.

$$\prod_{L_1 \cup L_2}(E_1 \bowtie_\theta E_2) \equiv \prod_{L_1 \cup L_2}(\prod_{L_1 \cup L_3}(E_1) \bowtie_\theta \prod_{L_2 \cup L_4}(E_2))$$

Similar equivalences hold for outerjoin operations: $⋈$, $⋉$, and $⋊$

# Equivalence Rules(Cont.)

9. The set operations union and intersection are commutative

$$E_1 \cup E_2 \equiv E_2 \cup E_1$$
$$E_1 \cap E_2 \equiv E_2 \cap E_1$$

(set difference is not commutative).

10. Set union and intersection are associative.

$$(E_1 \cup E_2) \cup E_3 \equiv E_1 \cup (E_2 \cup E_3)$$
$$(E_1 \cap E_2) \cap E_3 \equiv E_1 \cap (E_2 \cap E_3)$$

11. The selection operation distributes over $\cup$, $\cap$ and $-$.

    a. $\sigma_\theta (E_1 \cup E_2) \equiv \sigma_\theta (E_1) \cup \sigma_\theta(E_2)$
    b. $\sigma_\theta (E_1 \cap E_2) \equiv \sigma_\theta (E_1) \cap \sigma_\theta(E_2)$
    c. $\sigma_\theta (E_1 - E_2) \equiv \sigma_\theta (E_1) - \sigma_\theta(E_2)$
    d. $\sigma_\theta (E_1 \cap E_2) \equiv \sigma_\theta(E_1) \cap E_2$
    e. $\sigma_\theta (E_1 - E_2) \equiv \sigma_\theta(E_1) - E_2$

preceding equivalence does not hold for $\cup$

12. The projection operation distributes over union

$$\Pi_L(E_1 \cup E_2) \equiv (\Pi_L(E_1)) \cup (\Pi_L(E_2))$$

13. 13. Selection distributes over aggregation as below

$$\sigma_\theta(_G\gamma_A(E)) \equiv {_G}\gamma_A(\sigma_\theta(E))$$

provided $\theta$ only involves attributes in G

# Pushing Down Selections

- Query: Find the names of all instructors in the Music department, along with the titles of the courses that they teach

  - $\Pi_{name,\ title}(\sigma_{dept\_name=\ 'Music'}$
    $(instructor \bowtie (teaches \bowtie \Pi_{course\_id,\ title}\ (course))))$

- Transformation using rule 7a.

  - $\Pi_{name,\ title}((\sigma_{dept\_name=\ 'Music'}(instructor)) \bowtie$
    $(teaches \bowtie \Pi_{course\_id,\ title}\ (course)))$

- Performing the selection as early as possible reduces the size of the relation to be joined.

# Multiple Transformations

- Query: Find the names of all instructors in the Music department who have taught a course in 2017, along with the titles of the courses that they taught
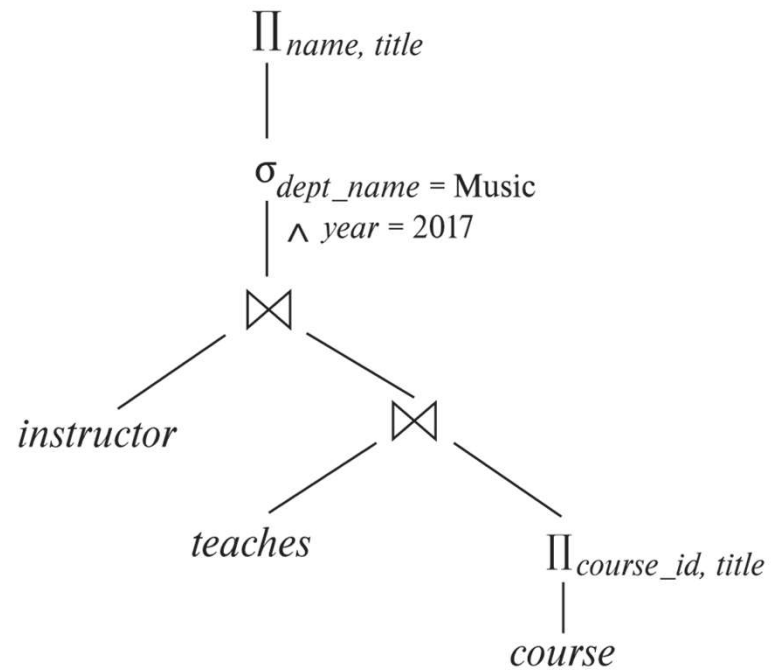
  - $\Pi_{name,\ title}(\sigma_{dept\_name=\ "Music"\wedge year\ =\ 2017}$
  $(instructor \bowtie (teaches \bowtie \Pi_{course\_id,\ title}\ (course))))$

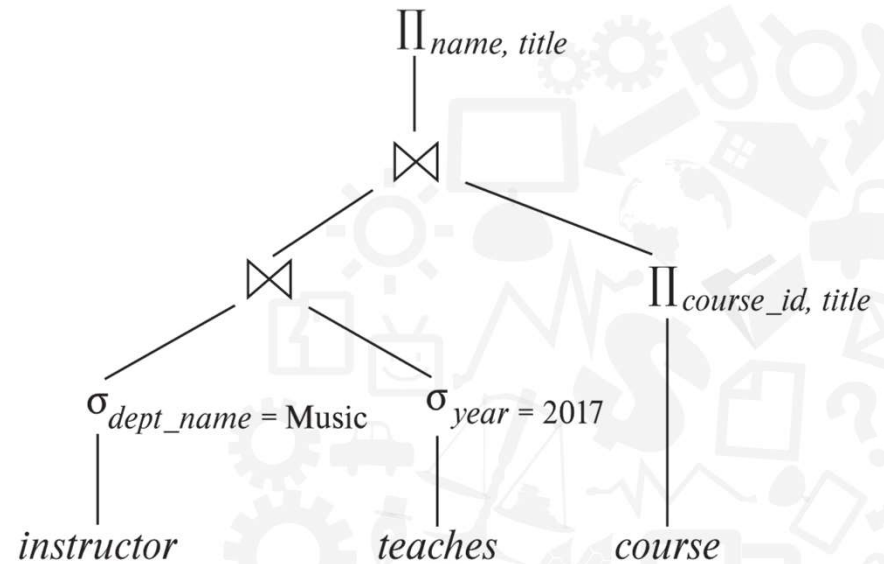- Transformation using join associatively (Rule 6a):

  - $\Pi_{name,\ title}(\sigma_{dept\_name=\ "Music"\wedge year\ =\ 2017}$
  $((instructor \bowtie teaches) \bowtie \Pi_{course\_id,\ title}\ (course)))$

- Second form provides an opportunity to apply the "perform selections early" rule, resulting in the subexpression

  $\sigma_{dept\_name\ =\ "Music"}\ (instructor) \bowtie \sigma_{year\ =\ 2017}\ (teaches)$

$\prod_{name,\ title}$

$\sigma_{dept\_name\ =\ Music\ \wedge\ year\ =\ 2017}$

instructor

teaches

$\prod_{course\_id,\ title}$

course

(a) Initial expression tree

$\prod_{name,\ title}$

$\sigma_{dept\_name\ =\ Music}$

$\sigma_{year\ =\ 2017}$

instructor

teaches

$\prod_{course\_id,\ title}$

course

(b) Tree after multiple transformations

# Pushing Down Projections

- Consider: $\Pi_{name, title}(\sigma_{dept\_name=\ \text{"Music"}}\ (instructor) \bowtie teaches)$
  $\bowtie \Pi_{course\_id, title}\ (course))))$

- When we compute

  $(\sigma_{dept\_name =\ \text{"Music"}}\ (instructor \bowtie teaches)$

  we obtain a relation whose schema is:
  (*ID, name, dept_name, salary, course_id, sec_id, semester, year*)

- Push projections using equivalence rules 8a and 8b; eliminate unneeded attributes from intermediate results to get:

  $\Pi_{name, title}(\Pi_{name, course\_id}\ ($
  $\sigma_{dept\_name=\ \text{"Music"}}\ (instructor) \bowtie teaches))$
  $\bowtie\ \Pi_{course\_id, title}\ (course))))$

- Performing the projection as early as possible reduces the size of the relation to be joined.

# Join Ordering

- For all relations $r_1$, $r_2$, and $r_3$,
$$(r_1 \bowtie r_2) \bowtie r_3 = r_1 \bowtie (r_2 \bowtie r_3)$$

- If $r_2 \bowtie r_3$ is quite large and $r_1 \bowtie r_2$ is small, we choose

$$(r_1 \bowtie r_2) \bowtie r_3$$

  so that we compute and store a smaller temporary relation.

- Consider the expression

$$\Pi_{name,\ title}(\sigma_{dept\_name=\ \text{"Music"}}(instructor) \bowtie teaches)$$
$$\bowtie \Pi_{course\_id,\ title}(course))))$$

- Could compute $teaches \bowtie \Pi_{course\_id,\ title}(course)$ first, and join result with

$$\sigma_{dept\_name=\ \text{"Music"}}(instructor)$$

  but the result of the first join is likely to be a large relation.

- Only a small fraction of the university's instructors are likely to be from the Music department
  - it is better to compute

$$\sigma_{dept\_name=\ \text{"Music"}}(instructor) \bowtie teaches \qquad \text{first.}$$

# Cost-based Optimization

- Query optimizers use equivalence rules to **systematically** generate expressions equivalent to the given expression
  - But, generating all equivalent expressions are expensive
  - Use of Dynamic Programming
- Choice of Evaluation Plans
  - Choosing the cheapest plan is imperative
- Cost of Evaluation Plan
  - Need statistics of input relations
  - Need to estimate statistics of expression results

# Cost-based Optimization(Cont.)

- Consider finding the best join-order for $r_1 \bowtie r_2 \bowtie \ldots \bowtie r_n$.

- There are $(2(n-1))!/(n-1)!$ different join orders for above expression. With $n = 7$, the number is 665280, with $n = 10$, the number is greater than 176 billion!

- No need to generate all the join orders. Using dynamic programming, the least-cost join order for any subset of $\{r_1, r_2, \ldots r_n\}$ is computed only once and stored for future use.

# Dynamic Programming

- To find best join tree for a set of *n* relations:
  - To find best plan for a set *S* of *n* relations, consider all possible plans of the form: $S_1 \bowtie (S - S_1)$ where $S_1$ is any non-empty subset of S.
  - Recursively compute costs for joining subsets of *S* to find the cost of each plan. Choose the cheapest of the $2^n - 2$ alternatives.
  - Base case for recursion: single relation access plan
    - Apply all selections on $R_i$ using best choice of indices on $R_i$
  - When plan for any subset is computed, store it and reuse it when it is required again, instead of recomputing it
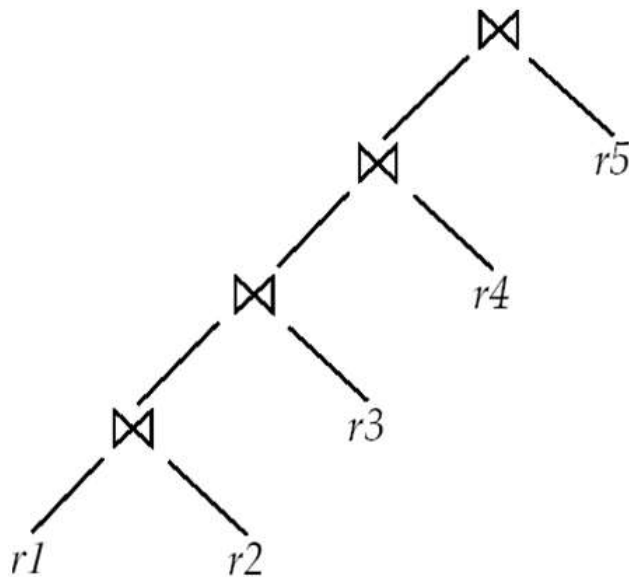    - Dynamic programming

# Join Order Optimization Algorithm

- **procedure** findbestplan(*S*)
  **if** (*bestplan*[*S*].*cost* ≠ ∞)
        **return** *bestplan*[*S*]
  // else *bestplan*[*S*] has not been computed earlier, compute it now
  **if** (*S* contains only 1 relation)
     set *bestplan*[*S*].*plan* and *bestplan*[*S*].*cost* based on the best way
     of accessing *S*  using selections on S and indices (if any) on S **else**
  **for each** non-empty subset *S1* of *S* such that *S1* ≠ *S*
       P1= findbestplan(*S1*)
       P2= findbestplan(*S - S1*)
       **for each** algorithm A for joining results of *P1* and *P2*
         ***...  compute plan and cost of using A (see next page) ..***
          **if** *cost* < *bestplan*[*S*].*cost*
            *bestplan*[*S*].*cost* = cost
            *bestplan*[*S*].*plan* = *plan*;
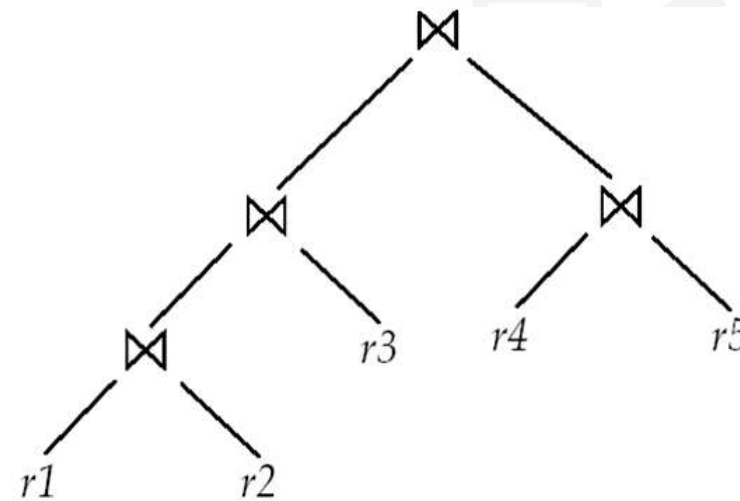  **return** *bestplan*[*S*]

# Left Deep Join Trees

- In **left-deep join trees,** the right-hand-side input for each join is a relation, not the result of an intermediate join.



(a) Left-deep join tree          (b) Non-left-deep join tree

# Heuristic Optimization

- Cost-based optimization is expensive, even with dynamic programming.

- Systems may use *heuristics* to reduce the number of choices that must be made in a cost-based fashion.

- Heuristic optimization transforms the query-tree by using a set of rules that typically (but not in all cases) improve execution performance:

  – Perform selection early (reduces the number of tuples)

  – Perform projection early (reduces the number of attributes)

  – Perform most restrictive selection and join operations (i.e., with smallest result size) before other similar operations.

  – Some systems use only heuristics, others combine heuristics with partial cost-based optimization.

# Query Optimizers

- Many optimizers considers only left-deep join orders.
  - Plus heuristics to push selections and projections down the query tree
  - Reduces optimization complexity and generates plans amenable to pipelined evaluation.

- Heuristic optimization used in some versions of Oracle:
  - Repeatedly pick "best" relation to join next
    - Starting from each of n starting points.  Pick best among these

- Even with the use of heuristics, cost-based query optimization imposes a substantial overhead.
  - But is worth it for expensive queries
  - Optimizers often use simple heuristics for very cheap queries, and perform exhaustive enumeration for more expensive queries

# Statistical Information for Cost Estimation

- $n_r$: number of tuples in a relation $r$.

- $b_r$: number of blocks containing tuples of $r$.

- $l_r$: size of a tuple of $r$.

- $f_r$: blocking factor of $r$ — i.e., the number of tuples of $r$ that fit into one block.

- $V(A, r)$: number of distinct values that appear in $r$ for attribute $A$; same as the size of $\prod_A(r)$.

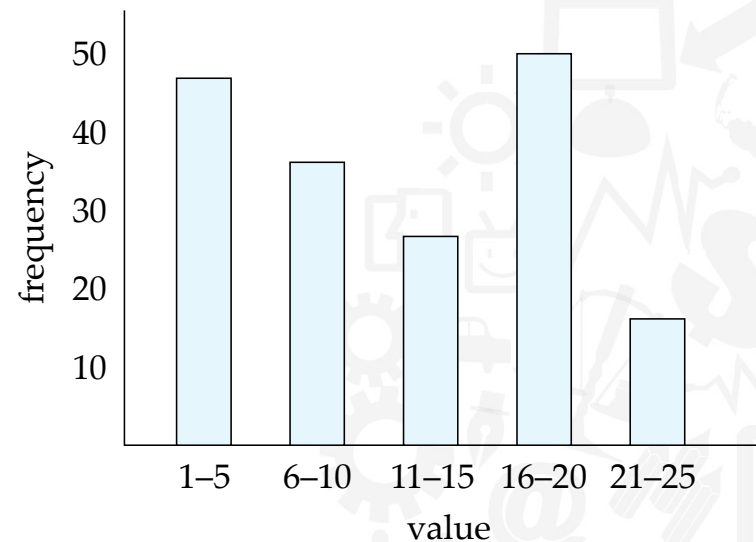- If tuples of $r$ are stored together physically in a file, then:

$$b_r = \left\lceil \frac{n_r}{f_r} \right\rceil$$

# Histogram

- Histograms and other statistics usually computed based on a **random sample**
- Histogram on attribute *age* of relation *person*



- **Equi-width** histograms
- **Equi-depth** histograms break up range such that each range has (approximately) the same number of tuples
  - E.g. (4, 8, 14, 19)
- Many databases also store *n* **most-frequent values** and their counts
  - Histogram is built on remaining values only

# Selection Size Estimation

- $\sigma_{A=v}(r)$
  - $n_r / V(A, r)$ : number of records that will satisfy the selection
  - Equality condition on a key attribute: *size estimate* = 1
- $\sigma_{A \leq V}(r)$ (case of $\sigma_{A \geq V}(r)$ is symmetric)
  - Let c denote the estimated number of tuples satisfying the condition.
  - If min(A,r) and max(A,r) are available in catalog
    - c = 0 if v < min(A,r)
    - c = $n_r \cdot \dfrac{v - \min(A, r)}{\max(A, r) - \min(A, r)}$
  - If histograms available, can refine above estimate
  - In absence of statistical information $c$ is assumed to be $n_r / 2$.

# Join Size Estimation

Running example:
$$student \bowtie takes$$

Catalog information for join examples:

- $n_{student} = 5{,}000$.
- $f_{student} = 50$, which implies that $b_{student} = 5000/50 = 100$.
- $n_{takes} = 10000$.
- $f_{takes} = 25$, which implies that $b_{takes} = 10000/25 = 400$.
- $V(ID, takes) = 2500$, which implies that on average, each student who has taken a course has taken 4 courses.
  - Attribute *ID* in *takes* is a foreign key referencing *student*.
  - $V(ID, student) = 5000$ (*primary key!*)

# Join Size Estimation

- The Cartesian product $r \times s$ contains $n_r . n_s$ tuples; each tuple occupies $s_r + s_s$ bytes.

- If $R \cap S = \varnothing$, then $r \bowtie s$ is the same as $r \times s$.

- If $R \cap S$ is a key for $R$, then a tuple of $s$ will join with at most one tuple from $r$
  - therefore, the number of tuples in $r \bowtie s$ is no greater than the number of tuples in $s$.

- If $R \cap S$ *in* $S$ is a foreign key in $S$ referencing $R$, then the number of tuples in $r \bowtie s$ is exactly the same as the number of tuples in $s$.
  - The case for $R \cap S$ being a foreign key referencing $S$ is symmetric.

- In the example query *student* $\bowtie$ *takes, ID* in *takes* is a foreign key referencing *student*
  - hence, the result has exactly $n_{takes}$ tuples, which is 10000

# Materialized Views

- A **materialized view** is a view whose contents are computed and stored.
- Consider the view
  **create view** *department_total_salary(dept_name, total_salary)* **as**
  **select** *dept_name*, **sum**(*salary*)
  **from** *instructor*
  **group by** *dept_name*
- Materializing the above view would be very useful if the total salary by department is required frequently
  – Saves the effort of finding multiple tuples and adding up their amounts
- The task of keeping a materialized view up-to-date with the underlying data is known as **materialized view maintenance**
- Materialized views can be maintained by recomputation on every update
- A better option is to use **incremental view maintenance**
  – **Changes to database relations are used to compute changes to the materialized view, which is then updated**

# Join Minimization

> **select** r.A, r.B
> **from** r, s
> **where** r.B = s.B

- Check if join with s is redundant, drop it
  - E.g., join condition is on foreign key from r to s, r.B is declared as not null, and no selection on s
  - Other sufficient conditions possible
    > **select** r.A, s2.B
    > **from** r, s **as** s1, s **as** s2
    > **where** r.B=s1.B **and** r.B = s2.B **and** s1.A < 20 **and** s2.A < 10
    - join with s1 is redundant and can be dropped (along with selection on s1)
  - Lots of research in this area since 70s/80s!

# Multi-Query Optimization

- Example

    Q1: **select * from** (r **natural join** t) **natural join** s

    Q2: **select * from** (r **natural join** u) **natural join** s

    – Both queries share common subexpression (r natural join s)

    – May be useful to compute (r natural join s) once and use it in both queries

    - But this may be more expensive in some situations

        – e.g. (r natural join s) may be expensive, plans as shown in queries may be cheaper

- **Multiquery optimization**: find best overall plan for a set of queries, expoiting sharing of common subexpressions between queries where it is useful

# Multiquery Optimization (Cont.)

- Simple heuristic used in some database systems:
  - optimize each query separately
  - detect and exploiting common subexpressions in the individual optimal query plans
    - May not always give best plan, but is cheap to implement
  - **Shared scans**: widely used special case of multiquery optimization
- Set of materialized views may share common subexpressions
  - As a result, view maintenance plans may share subexpressions
  - Multiquery optimization can be useful in such situations

# Question?

–source: https://www.fox.com/the-simpsons