

Estimation of the Size of Joins

- The Cartesian product $r \times s$ contains $n_r \cdot n_s$ tuples; each tuple occupies $s_r + s_s$ bytes.
- If $R \cap S = \emptyset$, 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

Estimation of the Size of Joins (Cont.)



- If $R \cap S = \{A\}$ is not a key for R or S .
If we assume that every tuple t_r in R produces tuples in $R \bowtie S$, the number of tuples in $R \bowtie S$ is estimated to be:
$$\frac{n_r \cdot n_s}{V(A, S)}$$

If the reverse is true, the estimate obtained will be:
$$\frac{n_r \cdot n_s}{V(A, R)}$$

The lower of these two estimates is probably the more accurate one.

- Can improve on above if histograms are available
 - Use formula similar to above, for each cell of histograms on the two relations

Estimation of the Size of Joins (Cont.)



- Compute the size estimates for *depositor* *customer* without using information about foreign keys:
 - $V(ID, takes) = 2500$, and
 $V(ID, student) = 5000$
 - The two estimates are $5000 * 10000/2500 = 20,000$ and $5000 * 10000/5000 = 10000$
 - We choose the lower estimate, which in this case, is the same as our earlier computation using foreign keys.

Size Estimation for Other Operations

- Projection: estimated size of $\Pi_A(r) = V(A,r)$
- Aggregation : estimated size of $g_F(r) = V(A,r)$
- Set operations
 - For unions/intersections of selections on the same relation: rewrite and use size estimate for selections
 - E.g. $\sigma_{\theta_1}(r) \cup \sigma_{\theta_2}(r)$ can be rewritten as $\sigma_{\theta_1 \vee \theta_2}(r)$
 - For operations on different relations:
 - estimated size of $r \cup s = \text{size of } r + \text{size of } s.$
 - estimated size of $r \cap s = \text{minimum size of } r \text{ and size of } s.$
 - estimated size of $r - s = r.$
 - All the three estimates may be quite inaccurate, but provide upper bounds on the sizes.

Size Estimation (Cont.)



- Outer join:

- Estimated size of $r \bowtie s = \text{size of } r + \text{size of } r$

- Case of right outer join is symmetric

- Estimated size of $r \bowtie s = \text{size of } r + \text{size of } r + \text{size of } s$

Estimation of Number of Distinct Values

Selections: $\sigma_{\theta}(r)$

- If θ forces A to take a specified value: $V(A, \sigma_{\theta}(r)) = 1$.
 - e.g., $A = 3$
- If θ forces A to take on one of a specified set of values:
 $V(A, \sigma_{\theta}(r)) = \text{number of specified values.}$
 - (e.g., $(A = 1 \vee A = 3 \vee A = 4)$),
- If the selection condition θ is of the form $A \text{ op } r$
estimated $V(A, \sigma_{\theta}(r)) = V(A.r) * s$
 - where s is the selectivity of the selection.
- In all the other cases: use approximate estimate of
 $\min(V(A,r), n_{\sigma_{\theta}(r)})$
 - More accurate estimate can be got using probability theory, but this one works fine generally

Estimation of Distinct Values (Cont.)

Joins: $r \bowtie s$

- If all attributes in A are from r

$$\text{estimated } V(A, r \bowtie s) = \min(V(A, r), n_{r \bowtie s})$$

- If A contains attributes $A1$ from r and $A2$ from s , then estimated $V(A, r \bowtie s) =$

$$\min(V(A1, r) * V(A2 - A1, s), V(A1 - A2, r) * V(A2, s), n_{r \bowtie s})$$

- More accurate estimate can be got using probability theory, but this one works fine generally

Estimation of Distinct Values (Cont.)

- Estimation of distinct values are straightforward for projections.
 - They are the same in $\Pi_{A(r)}$ as in r .
- The same holds for grouping attributes of aggregation.
- For aggregated values
 - For $\min(A)$ and $\max(A)$, the number of distinct values can be estimated as $\min(V(A,r), V(G,r))$ where G denotes grouping attributes
 - For other aggregates, assume all values are distinct, and use $V(G,r)$

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

Materialized View Maintenance

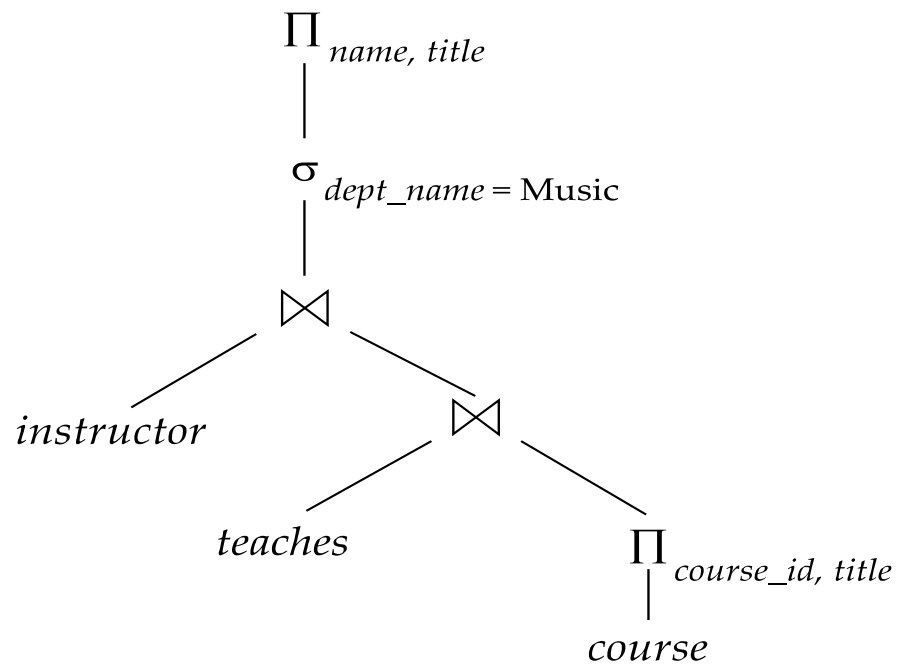
- The task of keeping a materialized view up-to-date with the underlying data is known as **materialized view maintenance**
- Materialized views can be maintained by recomputation on every update
- A better option is to use **incremental view maintenance**
 - **Changes to database relations are used to compute changes to the materialized view, which is then updated**
- View maintenance can be done by
 - Manually defining triggers on insert, delete, and update of each relation in the view definition
 - Manually written code to update the view whenever database relations are updated
 - Periodic recomputation (e.g. nightly)
 - Above methods are directly supported by many database systems
 - Avoids manual effort/correctness issues

Incremental View Maintenance

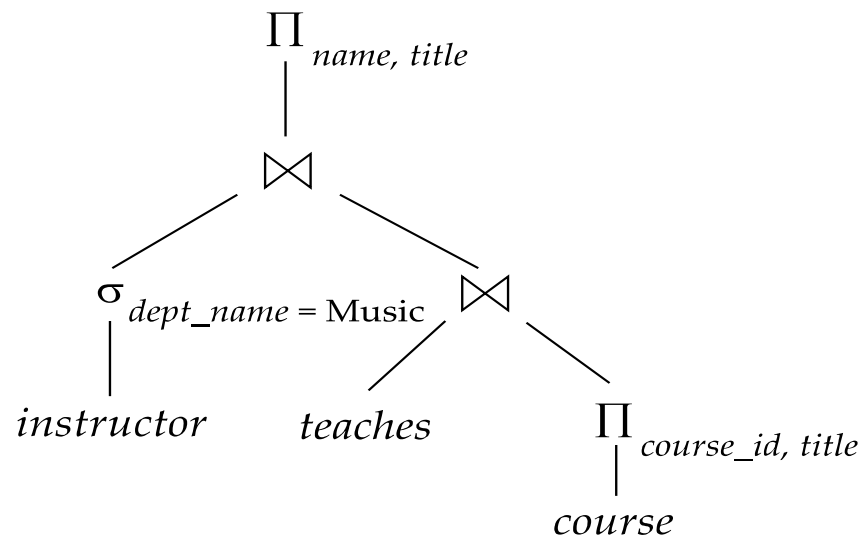
- The changes (inserts and deletes) to a relation or expressions are referred to as its **differential**
 - Set of tuples inserted to and deleted from r are denoted \mathbf{i}_r and \mathbf{d}_r
- To simplify our description, we only consider inserts and deletes
 - We replace updates to a tuple by deletion of the tuple followed by insertion of the update tuple
- We describe how to compute the change to the result of each relational operation, given changes to its inputs
- We then outline how to handle relational algebra expressions

Materialized View Selection

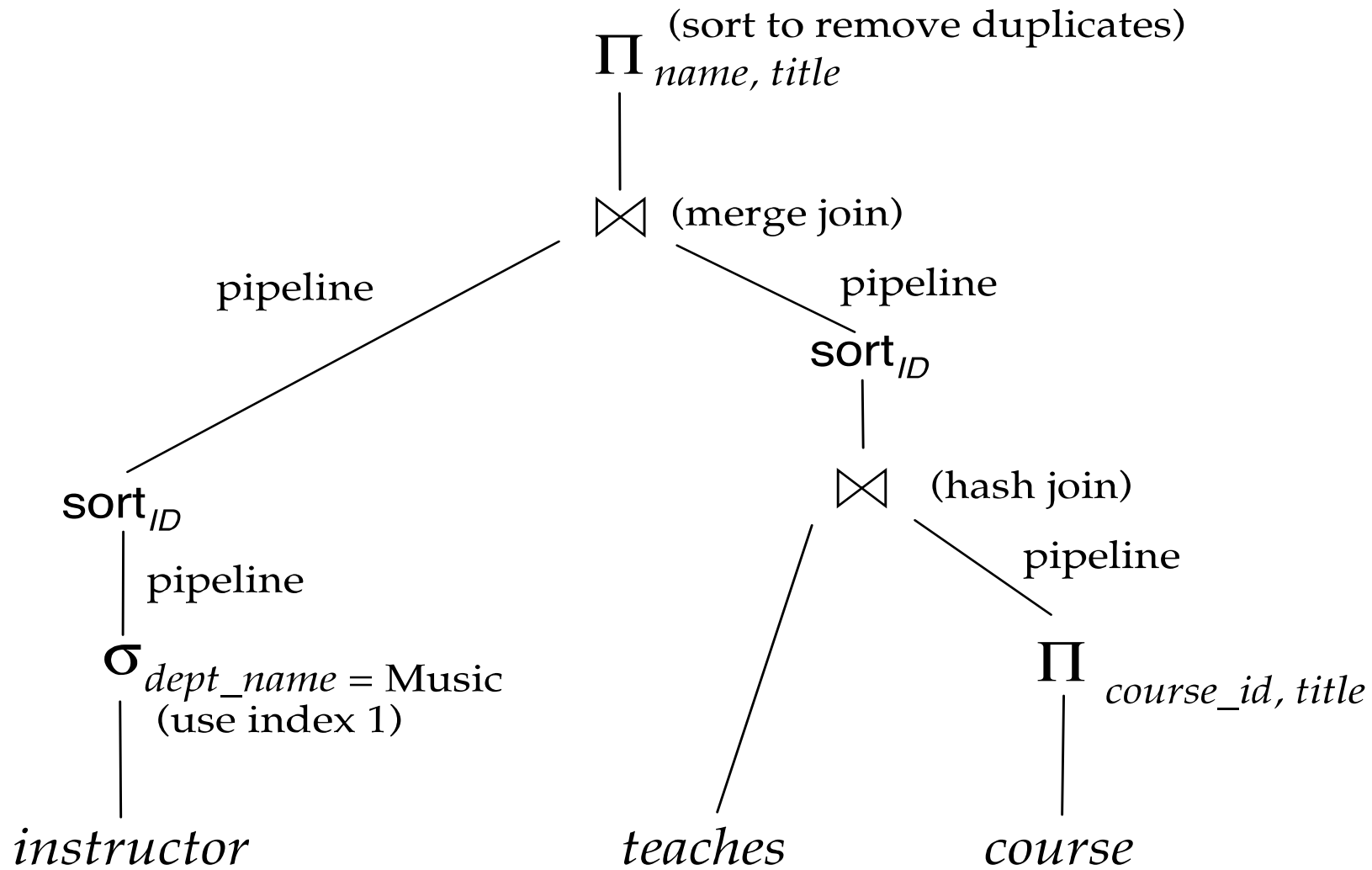
- **Materialized view selection**: “What is the best set of views to materialize?”.
- **Index selection**: “what is the best set of indices to create”
 - closely related, to materialized view selection
 - but simpler
- Materialized view selection and index selection based on typical system **workload** (queries and updates)
 - Typical goal: minimize time to execute workload , subject to constraints on space and time taken for some critical queries/updates
 - One of the steps in database tuning
 - more on tuning in later chapters
- Commercial database systems provide tools (called “tuning assistants” or “wizards”) to help the database administrator choose what indices and materialized views to create

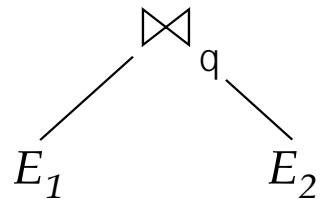


(a) Initial expression tree

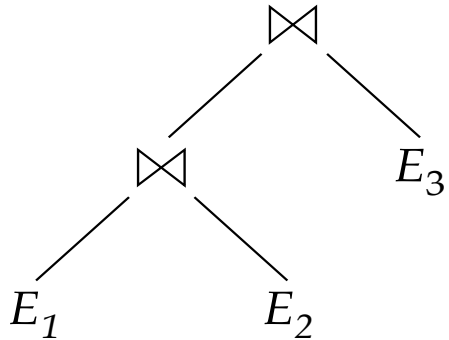
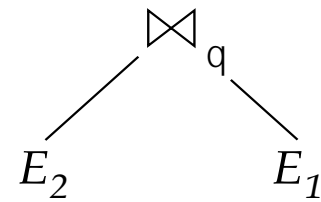


(b) Transformed expression tree

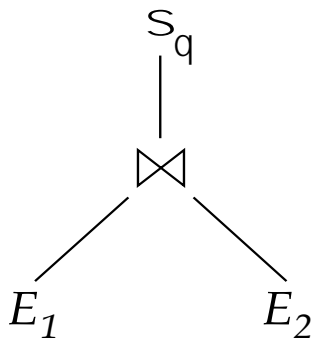
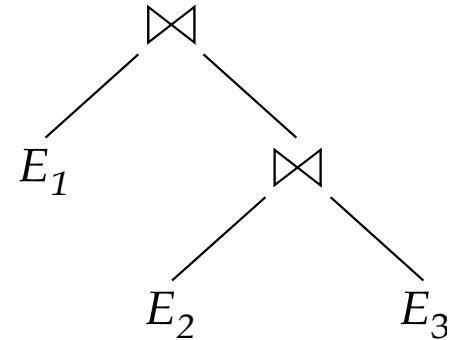




Rule 5



Rule 6.a



Rule 7.a

If q only has attributes from E_1

