

Principles of Database Systems (CS307)

Lecture 9: Relational Algebra

Ran Cheng

Department of Computer Science and Engineering
Southern University of Science and Technology

- Most contents are from slides made by Stéphane Faroult and the authors of Database System Concepts (7th Edition).
- Their original slides have been modified to adapt to the schedule of CS307 at SJUSTech.

Relational Algebra

- A procedural language consisting of a set of operations that take one or two relations as input and produce a new relation as their result.
- 6 Basic Operators:
 - select: σ
 - project: Π
 - union: \cup
 - set difference: $-$
 - Cartesian product: \times
 - rename: ρ

Cartesian product

- The term "Cartesian Product" is derived from the mathematician René Descartes.
- His work in coordinate geometry showed how to describe the position of points using ordered pairs of numbers.
- In a similar fashion, the Cartesian Product pairs elements from two sets.



"Cogito, ergo sum" ("I think, therefore I am")

Select Operation

- The select operation selects tuples that satisfy a given predicate
 - Notation: $\sigma_p(r)$
 - p is called the selection predicate
 - SQL: **select ***
- Example
 - Select those tuples of the *instructor* relation where the *instructor* is in the “Physics” department

$$\sigma_{dept_name = "Physics"}(instructor)$$

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
22222	Einstein	Physics	95000
33456	Gold	Physics	87000

```
sql
SELECT *
FROM instructor
WHERE dept_name = 'Physics';
```

Select Operation

- We allow comparisons using $=, \neq, >, \geq, <, \leq$ in the selection predicate
- We can combine several predicates into a larger predicate by using the connectives:
 \wedge (and), \vee (or), \neg (not)
- Example: Find the instructors in Physics with a salary greater than \$90,000, we write:

$$\sigma_{dept_name = "Physics" \wedge salary > 90,000} (instructor)$$

- The select predicate may include comparisons between two attributes.
 - Example, find all departments whose name is the same as their building name:

$$\sigma_{dept_name = building} (department)$$

Project Operation

- A unary (一元) operation that returns its argument relation, with certain attributes left out. It is applied to **Columns**.

Notation: $\Pi_{A_1, A_2, A_3 \dots A_k}(r)$

where A_1, A_2, \dots, A_k are attribute names and r is a relation name.

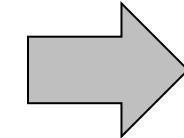
- The result is defined as **the relation of k columns** obtained by **erasing the columns** that are not listed
- Duplicate rows removed from result, since relations are sets

Project Operation

- Example: eliminate the *dept_name* attribute of instructor
 - Query:

$$\Pi_{ID, name, salary} (instructor)$$

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000



<i>ID</i>	<i>name</i>	<i>salary</i>
10101	Srinivasan	65000
12121	Wu	90000
15151	Mozart	40000
22222	Einstein	95000
32343	El Said	60000
33456	Gold	87000
45565	Katz	75000
58583	Califieri	62000
76543	Singh	80000
76766	Crick	72000
83821	Brandt	92000
98345	Kim	80000

Composition of Relational Operations

- The result of a relational-algebra operation is relation
 - ... and therefore, relational-algebra operations can be composed together into a relational-algebra expression
- Consider the query: Find the names of all instructors in the Physics department

$$\Pi_{name}(\sigma_{dept_name = "Physics"}(instructor))$$

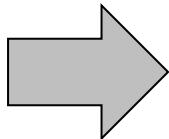
- Instead of giving the name of a relation as the argument of the projection operation, we give an expression that evaluates to a relation

Cartesian-Product Operation

- The Cartesian-product operation (denoted by \times) allows us to combine information from any two relations.
 - Example: the Cartesian product of the relations `instructor` and `teaches` is written as:
$$\text{instructor} \times \text{teaches}$$
- We construct a tuple of the result out of each possible pair of tuples
 - ... one from the `instructor` relation and one from the `teaches` relation (see next slide)
 - Since the `instructor` ID appears in both relations, we **distinguish** between these attributes by attaching to the attribute the name of the relation from which the attribute originally came – using ‘`xxx.xxx`’ .
 - `instructor.ID`
 - `teaches.ID`

The “instructor teaches” table

<i>ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>
10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000



<i>course_id</i>	<i>title</i>	<i>dept_name</i>	<i>credits</i>
BIO-101	Intro. to Biology	Biology	4
BIO-301	Genetics	Biology	4
BIO-399	Computational Biology	Biology	3
CS-101	Intro. to Computer Science	Comp. Sci.	4
CS-190	Game Design	Comp. Sci.	4
CS-315	Robotics	Comp. Sci.	3
CS-319	Image Processing	Comp. Sci.	3
CS-347	Database System Concepts	Comp. Sci.	3
EE-181	Intro. to Digital Systems	Elec. Eng.	3
FIN-201	Investment Banking	Finance	3
HIS-351	World History	History	3
MU-199	Music Video Production	Music	3
PHY-101	Physical Principles	Physics	4

Join Operation

- **Problem:** The Cartesian-Product “instructor \times teaches” associates **every tuple of instructor** with **every tuple of teaches**
 - Most of the resulting rows have information about instructors who did NOT teach a particular course

<i>instructor.ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>	<i>teaches.ID</i>	<i>course_id</i>	<i>sec_id</i>	<i>semester</i>	<i>year</i>
10101	Srinivasan	Comp. Sci.	65000	10101	CS-101	1	Fall	2017
10101	Srinivasan	Comp. Sci.	65000	10101	CS-315	1	Spring	2018
10101	Srinivasan	Comp. Sci.	65000	10101	CS-347	1	Fall	2017
10101	Srinivasan	Comp. Sci.	65000	12121	FIN-201	1	Spring	2018
10101	Srinivasan	Comp. Sci.	65000	15151	MU-199	1	Spring	2018
10101	Srinivasan	Comp. Sci.	65000	22222	PHY-101	1	Fall	2017

Join Operation

- To get only those tuples of “instructor \times teaches” that pertain to instructors and the courses that they taught, we write:

$$\sigma_{instructor.id = teaches.id} (instructor \times teaches)$$

- We get only those tuples of “instructor \times teaches” that pertain to instructors and the courses that they taught
 - i.e., those tuples where $instructor.id = teaches.id$

Join Operation

- The table corresponding to $\sigma_{instructor.id = teaches.id} (instructor \times teaches)$:

<i>instructor.ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>	<i>teaches.ID</i>	<i>course_id</i>	<i>sec_id</i>	<i>semester</i>	<i>year</i>
10101	Srinivasan	Comp. Sci.	65000	10101	CS-101	1	Fall	2017
10101	Srinivasan	Comp. Sci.	65000	10101	CS-315	1	Spring	2018
10101	Srinivasan	Comp. Sci.	65000	10101	CS-347	1	Fall	2017
12121	Wu	Finance	90000	12121	FIN-201	1	Spring	2018
15151	Mozart	Music	40000	15151	MU-199	1	Spring	2018
22222	Einstein	Physics	95000	22222	PHY-101	1	Fall	2017
32343	El Said	History	60000	32343	HIS-351	1	Spring	2018
45565	Katz	Comp. Sci.	75000	45565	CS-101	1	Spring	2018
45565	Katz	Comp. Sci.	75000	45565	CS-319	1	Spring	2018
76766	Crick	Biology	72000	76766	BIO-101	1	Summer	2017
76766	Crick	Biology	72000	76766	BIO-301	1	Summer	2018
83821	Brandt	Comp. Sci.	92000	83821	CS-190	1	Spring	2017
83821	Brandt	Comp. Sci.	92000	83821	CS-190	2	Spring	2017
83821	Brandt	Comp. Sci.	92000	83821	CS-319	2	Spring	2018
98345	Kim	Elec. Eng.	80000	98345	EE-181	1	Spring	2017

Join Operation

- The table corresponding to $\sigma_{instructor.id = teaches.id} (instructor \times teaches)$:

<i>instructor.ID</i>	<i>name</i>	<i>dept_name</i>	<i>salary</i>	<i>teaches.ID</i>	<i>course_id</i>	<i>sec_id</i>	<i>semester</i>	<i>year</i>
10101	Srinivasan	Comp. Sci.	65000	10101	CS-101	1	Fall	2017
10101	Srinivasan	Comp. Sci.	65000	10101	CS-315	1	Spring	2018
10101	Srinivasan	Comp. Sci.	65000	10101	CS-347	1	Fall	2017
12121	Wu	Finance	90000	12121	FIN-201	1	Spring	2018
15151	Mozart	Music	40000	15151	MU-199	1	Spring	2018
22222	Einstein	Physics	95000	22222	PHY-101	1	Fall	2017
32343	El Said	History	60000	32343	HIS-351	1	Spring	2018
45565	Katz	Comp. Sci.	75000	45565	CS-101	1	Spring	2018
45565	Katz	Comp. Sci.	75000	45565	CS-319	1	Spring	2018
76766	Crick	Biology	72000	76766	BIO-101	1	Summer	2017
76766	Crick	Biology	72000	76766	BIO-301	1	Summer	2018
83821	Brandt	Comp. Sci.	92000	83821	CS-190	1	Spring	2017
83821	Brandt	Comp. Sci.	92000	83821	CS-190	2	Spring	2017
83821	Brandt	Comp. Sci.	92000	10101	Srinivasan	Comp. Sci.	65000	12121
98345	Kim	Eng. E	...	10101	Srinivasan	Comp. Sci.	65000	15151
				10101	Srinivasan	Comp. Sci.	65000	22222
				10101	Srinivasan	Comp. Sci.	65000	PHY-101

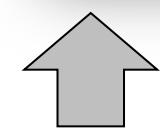
... will NOT include such tuples (rows) with different IDs

Recall: The Old Way of Writing Joins

- Use commas to separate the tables
 - Example: The solution for the same question in the previous slide
- A little bit history:
 - join was introduced in SQL-1999 (later than this original way)
- Problem:
 - If you forget a comma, it will still work sometimes (interpreted as “renaming”)



```
select m.title, c.credited_as,  
       p.first_name, p.surname  
  from movies m,  
       credits c,  
       people p  
 where c.movieid = m.movieid  
   and p.peopleid = c.peopleid  
   and m.country = 'cn'
```



The SQL syntax was derived from the form of cartesian products in relational algebra

$$\sigma_{movies.id = credits.id \wedge people.peopleid = credits.peopleid \wedge movies.country = "cn"} (Movies \times Credits \times People)$$

Recall: The Old Way of Writing Joins

- Use commas to separate the tables
 - Example: The solution for the same question in the previous slide
- A little bit history:
 - join was introduced in SQL-1999 (later than this original way)
- Problem:
 - If you forget a comma, it will still work sometimes (interpreted as “renaming”)

The “select operation” is written as the “where” clause here

$$\sigma_{movies.id = credits.id \wedge people.peopleid = credits.peopleid \wedge movies.country = "cn"}$$

```
select m.title, c.credited_as,  
       p.first_name, p.surname  
  from movies m,  
       credits c,  
       people p  
 where c.movieid = m.movieid  
   and p.peopleid = c.peopleid  
   and m.country = 'cn'
```



The SQL syntax was derived from the form of cartesian products in relational algebra

(Movies \times Credits \times People)

Join Operation

- The join operation allows us to combine a select operation and a Cartesian-Product operation into a single operation
 - Consider relations $r(R)$ and $s(S)$:
 - Let “**theta (θ)**” be a predicate (谓词) on attributes in the schema R “union” S. The join operation $r \bowtie_{\theta} s$ is defined as follows:

$$r \bowtie_{\theta} s = \sigma_{\theta} (r \times s)$$

- Thus, $\sigma_{instructor.id = teaches.id} (instructor \times teaches)$ can equivalently be written as:
 $instructor \bowtie_{Instructor.id = teaches.id} teaches$
- **Predicate** -- a condition or expression that evaluates to either true or false

Union Operation

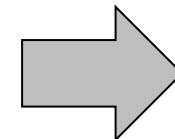
- The union operation allows us to combine two relations
 - Notation: $r \cup s$
- For $r \cup s$ to be valid:
 - r, s must have the **same arity** (same number of attributes)
 - The attribute domains must be **compatible**
 - Example: 2nd column of r deals with the same type of values as does the 2nd column of s

Union Operation

- Example: To find all courses taught in the Fall 2017 semester, or in the Spring 2018 semester, or in both

$$\Pi_{course_id} (\sigma_{semester="Fall" \wedge year=2017}(section)) \cup \Pi_{course_id} (\sigma_{semester="Spring" \wedge year=2018}(section))$$

<i>course_id</i>	<i>sec_id</i>	<i>semester</i>	<i>year</i>	<i>building</i>	<i>room_number</i>	<i>time_slot_id</i>
BIO-101	1	Summer	2017	Painter	514	B
BIO-301	1	Summer	2018	Painter	514	A
CS-101	1	Fall	2017	Packard	101	H
CS-101	1	Spring	2018	Packard	101	F
CS-190	1	Spring	2017	Taylor	3128	E
CS-190	2	Spring	2017	Taylor	3128	A
CS-315	1	Spring	2018	Watson	120	D
CS-319	1	Spring	2018	Watson	100	B
CS-319	2	Spring	2018	Taylor	3128	C
CS-347	1	Fall	2017	Taylor	3128	A
EE-181	1	Spring	2017	Taylor	3128	C
FIN-201	1	Spring	2018	Packard	101	B
HIS-351	1	Spring	2018	Painter	514	C
MU-199	1	Spring	2018	Packard	101	D
PHY-101	1	Fall	2017	Watson	100	A



<i>course_id</i>
CS-101
CS-315
CS-319
CS-347
FIN-201
HIS-351
MU-199
PHY-101

Set-Intersection Operation

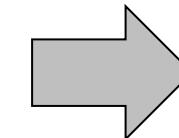
- The set-intersection operation allows us to find tuples that are in both the input relations
 - Notation: $r \cap s$
- Assume (same as Union):
 - r, s have the same arity
 - Attributes of r and s are compatible

Set-Intersection Operation

- Example: Find the set of all courses taught in both the Fall 2017 and the Spring 2018 semesters

$$\Pi_{course_id} (\sigma_{semester = "Fall"} \wedge year = 2017 (section)) \cap \Pi_{course_id} (\sigma_{semester = "Spring"} \wedge year = 2018 (section))$$

course_id	sec_id	semester	year	building	room_number	time_slot_id
BIO-101	1	Summer	2017	Painter	514	B
BIO-301	1	Summer	2018	Painter	514	A
CS-101	1	Fall	2017	Packard	101	H
CS-101	1	Spring	2018	Packard	101	F
CS-190	1	Spring	2017	Taylor	3128	E
CS-190	2	Spring	2017	Taylor	3128	A
CS-315	1	Spring	2018	Watson	120	D
CS-319	1	Spring	2018	Watson	100	B
CS-319	2	Spring	2018	Taylor	3128	C
CS-347	1	Fall	2017	Taylor	3128	A
EE-181	1	Spring	2017	Taylor	3128	C
FIN-201	1	Spring	2018	Packard	101	B
HIS-351	1	Spring	2018	Painter	514	C
MU-199	1	Spring	2018	Packard	101	D
PHY-101	1	Fall	2017	Watson	100	A



course_id
CS-101

Set Difference Operation

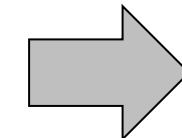
- The set-difference operation allows us to find tuples that are in one relation but are not in another
 - Notation: $r - s$
- Assume (same as Union and Set Intersection):
 - r, s have the same arity
 - Attributes of r and s are compatible

Set Difference Operation

- Example: Find all courses taught in the Fall 2017 semester, but not in the Spring 2018 semester

$$\Pi_{course_id} (\sigma_{semester="Fall" \wedge year=2017}(section)) - \Pi_{course_id} (\sigma_{semester="Spring" \wedge year=2018}(section))$$

course_id	sec_id	semester	year	building	room_number	time_slot_id
BIO-101	1	Summer	2017	Painter	514	B
BIO-301	1	Summer	2018	Painter	514	A
CS-101	1	Fall	2017	Packard	101	H
CS-101	1	Spring	2018	Packard	101	F
CS-190	1	Spring	2017	Taylor	3128	E
CS-190	2	Spring	2017	Taylor	3128	A
CS-315	1	Spring	2018	Watson	120	D
CS-319	1	Spring	2018	Watson	100	B
CS-319	2	Spring	2018	Taylor	3128	C
CS-347	1	Fall	2017	Taylor	3128	A
EE-181	1	Spring	2017	Taylor	3128	C
FIN-201	1	Spring	2018	Packard	101	B
HIS-351	1	Spring	2018	Painter	514	C
MU-199	1	Spring	2018	Packard	101	D
PHY-101	1	Fall	2017	Watson	100	A



course_id
CS-347
PHY-101

The Assignment Operation

- It is convenient at times to write a relational-algebra expression by **assigning parts of it to temporary relation variables**
 - The assignment operation is denoted by \leftarrow and works like **assignment** in a programming language
- Example: Find all instructor in the “Physics” and Music department

$$Physics \leftarrow \sigma_{dept_name = "Physics"}(instructor)$$
$$Music \leftarrow \sigma_{dept_name = "Music"}(instructor)$$
$$Physics \cup Music$$

- With the assignment operation, a query can be **written as a sequential program consisting of a series of assignments** followed by **an expression** whose value is displayed as the result of the query.

The Rename Operation

- The results of relational-algebra expressions do not have a name that we can use to refer to them
 - The rename operator, ρ , is provided for that purpose
 - The expression $\rho_x(E)$ returns the result of expression E under the name x
 - Another form of the rename operation which also renames the columns:
 - $\rho_{x(A1, A2, \dots An)}(E)$

Equivalent Queries

- There is **more than one way** to write a query in relational algebra - **why?**
 - Example: Find information about courses taught by instructors in the Physics department with salary greater than 90,000
 - Query 1
$$\sigma_{dept_name = "Physics" \wedge salary > 90,000} (instructor)$$
 - Query 2
$$\sigma_{dept_name = "Physics"} (\sigma_{salary > 90.000} (instructor))$$
 - The two queries are not identical
 - they are, however, **equivalent** -- they give the same result on any database

Equivalent Queries

- Example: Find information about courses taught by instructors in the Physics department
 - Query 1

$$\sigma_{dept_name = "Physics"}(instructor \bowtie_{instructor.ID = teaches.ID} teaches)$$

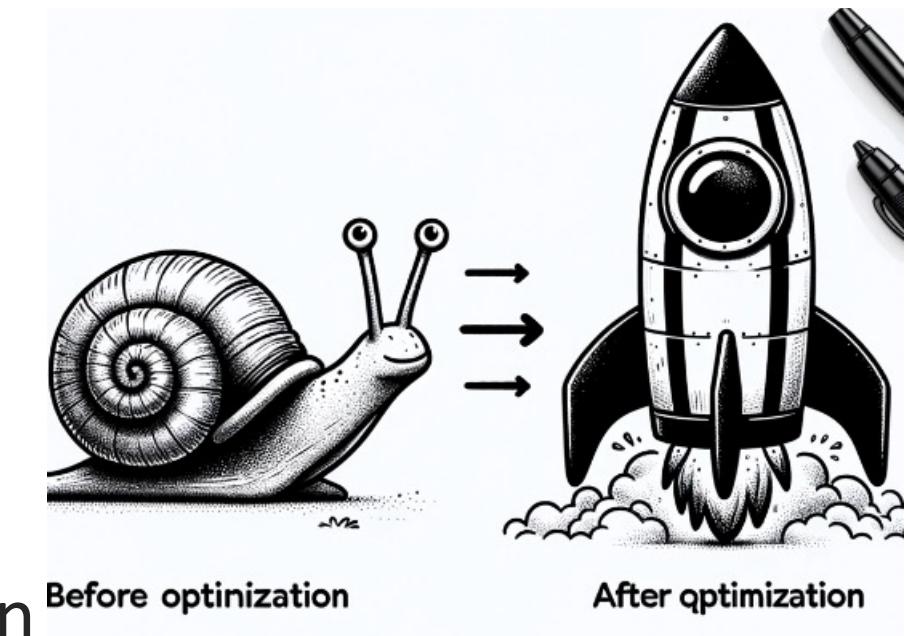
- (Join first, then select)
 - Query 2
 - (Select first, then join)

Equivalent Queries

- Application of Relational Algebra: **Query Optimization**
 - Transform queries into equivalent ones with less computational cost

Query Optimization

- **Definition:** Query optimization is the process of improving the efficiency of a query to minimize execution time and resource usage.
- **Purpose:** Optimization aims to provide rapid query results and reduce the load on the database system.
- **Process:** Involves selecting the most efficient way to **execute** a given query by considering different query plans.
- **Why It Matters:** The speed of data retrieval is critical in high-volume database environments, impacting user experience and system performance.



Query Optimization

The Role of the Database Management System (DBMS)

- **DBMS Functions:** The DBMS translates and executes SQL queries, ensuring data integrity and performance.
- **Query Optimizer:** A core component of the DBMS that determines the most efficient way to execute a query.
- **Execution Plan:** The blueprint created by the optimizer that the database engine follows to retrieve data.
- **Query Executor:** The DBMS component that carries out the execution plan and returns the query results.

Query Optimization

Factors Affecting Query Performance

- **Database Schema Design:** Proper normalization and table design can significantly impact performance.
- **Index Usage:** Correct indexing strategies can reduce data retrieval times.
- **Statistics:** Up-to-date data distribution statistics help the optimizer create efficient query plans.
- **Query Complexity:** Simpler queries typically run faster; complex queries may need optimization techniques to improve performance.

Query Optimization

Writing Efficient SQL Queries

- **Selective Data Retrieval:** Request only the data you need using specific column names instead of '*'.
- **Join Types:** Understand different join operations and use the most appropriate one for your query.
- **Subquery Usage:** Subqueries can be powerful but may lead to inefficiency; consider alternative methods like joins.
- **Indexing:** Use indexes wisely to speed up searches on large tables.

Query Optimization

Tools and Techniques for Query Optimization

- **Analyzing Tools:** Use tools like EXPLAIN PLAN to understand how a query will be executed.
- **Identifying Bottlenecks:** Look for slow parts of your query that could be causing delays.
- **Optimization Techniques:** Apply strategies such as indexing, query refactoring, and materialized views to improve performance.
- **Regular Maintenance:** Keep statistics up-to-date and periodically review query performance.