# DTS207TC Database Development and Design

## Lecture 2
## Mid-level SQL

Di Zhang, Autumn 2025

**Page titles with * will not be assessed**

# *SQL Interview Questions Website

- China Big Company

  - https://www.nowcoder.com/ta/sql

  - https://www.mianshiya.com/tag/SQL

- International Company

  - https://leetcode.cn/problemset/database/

  - https://www.stratascratch.com

-- You will be asked to be tested without AI and Internet in real interview…

# *Tips on CW

- Correctness

- Efficiency

- Completeness

- Document quality

- **Join operations** take two relations and return as a result another relation.

- A join operation is a Cartesian product which requires that tuples in the two relations match (under some condition).  It also specifies the attributes that are present in the result of the join

- The join operations are typically used as subquery expressions in the **from** clause

- Three types of joins:

  - Natural join

  - Inner join

  - Outer join

# Natural Join in SQL

- Natural join matches tuples with the same values for all common attributes, and retains only one copy of each common column.

- List the names of instructors along with the course ID of the courses that they taught

  - **select** *name*, *course_id*
    **from** *students, takes*
    **where** *student.ID = takes.ID*;

- Same query in SQL with "natural join" construct

  - **select** *name*, *course_id*
    **from** *student* **natural join** *takes*;

- The **from** clause can have multiple relations combined using natural join:

  **select** $A_1, A_2, \ldots A_n$
  **from** $r_1$ **natural join** $r_2$ **natural join** $..$ **natural join** $r_n$
  **where** $P$ ;

| ID | name | dept_name | tot_cred |
|-------|----------|------------|----------|
| 00128 | Zhang | Comp. Sci. | 102 |
| 12345 | Shankar | Comp. Sci. | 32 |
| 19991 | Brandt | History | 80 |
| 23121 | Chavez | Finance | 110 |
| 44553 | Peltier | Physics | 56 |
| 45678 | Levy | Physics | 46 |
| 54321 | Williams | Comp. Sci. | 54 |
| 55739 | Sanchez | Music | 38 |
| 70557 | Snow | Physics | 0 |
| 76543 | Brown | Comp. Sci. | 58 |
| 76653 | Aoi | Elec. Eng. | 60 |
| 98765 | Bourikas | Elec. Eng. | 98 |
| 98988 | Tanaka | Biology | 120 |

| ID | course_id | sec_id | semester | year | grade |
|-------|----------|--------|----------|------|-------|
| 00128 | CS-101 | 1 | Fall | 2017 | A |
| 00128 | CS-347 | 1 | Fall | 2017 | A- |
| 12345 | CS-101 | 1 | Fall | 2017 | C |
| 12345 | CS-190 | 2 | Spring | 2017 | A |
| 12345 | CS-315 | 1 | Spring | 2018 | A |
| 12345 | CS-347 | 1 | Fall | 2017 | A |
| 19991 | HIS-351 | 1 | Spring | 2018 | B |
| 23121 | FIN-201 | 1 | Spring | 2018 | C+ |
| 44553 | PHY-101 | 1 | Fall | 2017 | B- |
| 45678 | CS-101 | 1 | Fall | 2017 | F |
| 45678 | CS-101 | 1 | Spring | 2018 | B+ |
| 45678 | CS-319 | 1 | Spring | 2018 | B |
| 54321 | CS-101 | 1 | Fall | 2017 | A- |
| 54321 | CS-190 | 2 | Spring | 2017 | B+ |
| 55739 | MU-199 | 1 | Spring | 2018 | A- |
| 76543 | CS-101 | 1 | Fall | 2017 | A |
| 76543 | CS-319 | 2 | Spring | 2018 | A |
| 76653 | EE-181 | 1 | Spring | 2017 | C |
| 98765 | CS-101 | 1 | Fall | 2017 | C- |
| 98765 | CS-315 | 1 | Spring | 2018 | B |
| 98988 | BIO-101 | 1 | Summer | 2017 | A |
| 98988 | BIO-301 | 1 | Summer | 2018 | null |

| ID | name | dept_name | tot_cred | course_id | sec_id | semester | year | grade |
|---|---|---|---|---|---|---|---|---|
| 00128 | Zhang | Comp. Sci. | 102 | CS-101 | 1 | Fall | 2017 | A |
| 00128 | Zhang | Comp. Sci. | 102 | CS-347 | 1 | Fall | 2017 | A- |
| 12345 | Shankar | Comp. Sci. | 32 | CS-101 | 1 | Fall | 2017 | C |
| 12345 | Shankar | Comp. Sci. | 32 | CS-190 | 2 | Spring | 2017 | A |
| 12345 | Shankar | Comp. Sci. | 32 | CS-315 | 1 | Spring | 2018 | A |
| 12345 | Shankar | Comp. Sci. | 32 | CS-347 | 1 | Fall | 2017 | A |
| 19991 | Brandt | History | 80 | HIS-351 | 1 | Spring | 2018 | B |
| 23121 | Chavez | Finance | 110 | FIN-201 | 1 | Spring | 2018 | C+ |
| 44553 | Peltier | Physics | 56 | PHY-101 | 1 | Fall | 2017 | B- |
| 45678 | Levy | Physics | 46 | CS-101 | 1 | Fall | 2017 | F |
| 45678 | Levy | Physics | 46 | CS-101 | 1 | Spring | 2018 | B+ |
| 45678 | Levy | Physics | 46 | CS-319 | 1 | Spring | 2018 | B |
| 54321 | Williams | Comp. Sci. | 54 | CS-101 | 1 | Fall | 2017 | A- |
| 54321 | Williams | Comp. Sci. | 54 | CS-190 | 2 | Spring | 2017 | B+ |
| 55739 | Sanchez | Music | 38 | MU-199 | 1 | Spring | 2018 | A- |
| 76543 | Brown | Comp. Sci. | 58 | CS-101 | 1 | Fall | 2017 | A |
| 76543 | Brown | Comp. Sci. | 58 | CS-319 | 2 | Spring | 2018 | A |
| 76653 | Aoi | Elec. Eng. | 60 | EE-181 | 1 | Spring | 2017 | C |
| 98765 | Bourikas | Elec. Eng. | 98 | CS-101 | 1 | Fall | 2017 | C- |
| 98765 | Bourikas | Elec. Eng. | 98 | CS-315 | 1 | Spring | 2018 | B |
| 98988 | Tanaka | Biology | 120 | BIO-101 | 1 | Summer | 2017 | A |
| 98988 | Tanaka | Biology | 120 | BIO-301 | 1 | Summer | 2018 | null |

# Dangerous in Natural Join

- Beware of unrelated attributes with same name which get equated incorrectly

- Example -- List the names of students instructors along with the titles of courses that they have taken

  - Correct version

    **select** *name*, *title*
     **from** *student* **natural join** *takes*, *course*
     **where** *takes*.*course_id* = *course*.*course_id*;

  - Incorrect version

    **select** *name*, *title*
     **from** *student* **natural join** *takes* **natural join** *course*;

  - This query omits all (student name, course title) pairs where the student takes a course in a department other than the student's own department.

  - The  correct  version (above), correctly outputs such pairs.

# Natural Join with Using Clause

- To avoid the danger of equating attributes erroneously, we can use the "**using**" construct that allows us to specify exactly which columns should be equated.

- Query example

  **select** *name*, *title*
  **from** (*student* **natural join** *takes*) **join** *course* **using** (*course_id*)

- **Join operations** take two relations and return as a result another relation.

- These additional operations are typically used as subquery expressions in the **from** clause

- **Join condition** – defines which tuples in the two relations match.

- **Join type** – defines how tuples in each relation that do not match any tuple in the other relation (based on the join condition) are treated.

| Join types |
|---|
| inner join |
| left outer join |
| right outer join |
| full outer join |

| Join conditions |
|---|
| natural |
| on $<$ predicate $>$ |
| using $(A_1, A_2, \ldots, A_n)$ |

# Join Condition

- The **on** condition allows a general predicate over the relations being joined

- This predicate is written like a **where** clause predicate except for the use of the keyword **on**

- Query example

  **select** *
   **from** *student* **join** *takes* **on** *student_ID* = *takes_ID*

  - The **on** condition above specifies that a tuple from *student* matches a tuple from *takes* if their *ID* values are equal.

- Equivalent to:

  **select** *
   **from** *student , takes*
   **where** *student_ID* = *takes_ID*

- The **on** condition allows a general predicate over the relations being joined.

- This predicate is written like a **where** clause predicate except for the use of the keyword **on**.

- Query example

  **select** *
   **from** *student* **join** *takes* **on** *student_ID* = *takes_ID*

  - The **on** condition above specifies that a tuple from *student* matches a tuple from *takes* if their *ID* values are equal.

- Equivalent to:

  **select** *
   **from** *student , takes*
   **where** *student_ID* = *takes_ID*

# Outer Join

An extension of the join operation that avoids loss of information.

Computes the join and then adds tuples form one relation that does not match tuples in the other relation to the result of the join.

- Uses *null* values.

- Three forms of outer join:

  - left outer join

  - right outer join

  - full outer join

- Relation *course*

| course_id | title | dept_name | credits |
|-----------|-------|-----------|---------|
| BIO-301 | Genetics | Biology | 4 |
| CS-190 | Game Design | Comp. Sci. | 4 |
| CS-315 | Robotics | Comp. Sci. | 3 |

- Relation *prereq*

| course_id | prereq_id |
|-----------|-----------|
| BIO-301 | BIO-101 |
| CS-190 | CS-101 |
| CS-347 | CS-101 |

- Observe that

   *course* information is missing CS-347

- *course* **natural left outer join** *prereq*

| course_id | title | dept_name | credits | prereq_id |
|-----------|-------|-----------|---------|-----------|
| BIO-301 | Genetics | Biology | 4 | BIO-101 |
| CS-190 | Game Design | Comp. Sci. | 4 | CS-101 |
| CS-315 | Robotics | Comp. Sci. | 3 | *null* |

- *course* **natural right outer join** *prereq*

| course_id | title | dept_name | credits | prereq_id |
|-----------|-------------|------------|---------|-----------|
| BIO-301 | Genetics | Biology | 4 | BIO-101 |
| CS-190 | Game Design | Comp. Sci. | 4 | CS-101 |
| CS-347 | null | null | null | CS-101 |

- *course* **natural full outer join** *prereq*

| course_id | title | dept_name | credits | prereq_id |
|-----------|-------------|-------------|---------|-----------|
| BIO-301 | Genetics | Biology | 4 | BIO-101 |
| CS-190 | Game Design | Comp. Sci. | 4 | CS-101 |
| CS-315 | Robotics | Comp. Sci. | 3 | *null* |
| CS-347 | *null* | *null* | *null* | CS-101 |

- The following minimalist but real-world example demonstrates how different approaches can lead to explosive intermediate results and vastly different computational effort, yet the final answer remains the same.
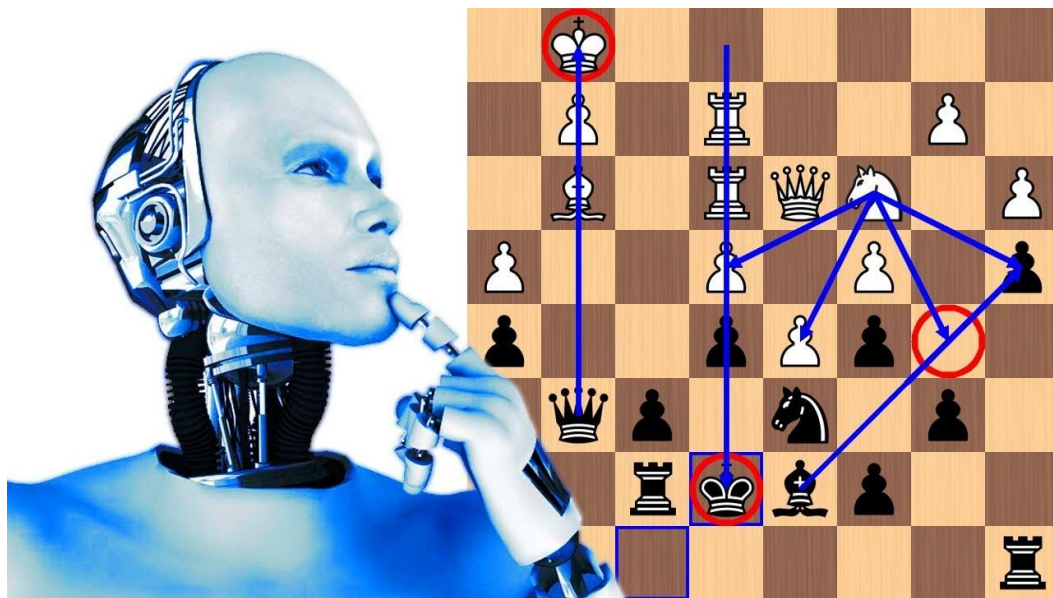
  - See '3_compare_join.txt'

- In most cases, PostgreSQL will automatically optimize. However, if you write the JOIN conditions as complex expressions or use views/CTEs, the optimizer may not be able to see the filter conditions, resulting in performance degradation.

- Similar approach to playing chesses with AI

  - Define legal moves

  - Define an evaluation function

  - Search for the best solution（can be solved by various methods, such as Paper）

- In some cases, it is not desirable for all users to see the entire logical model (that is, all the actual relations stored in the database.)

- Consider a person who needs to know an instructors name and department, but not the salary.  This person should see a relation described, in SQL, by

      **select** *ID*, *name*, *dept_name*
      **from** *instructor*

- A **view** provides a mechanism to hide certain data from the view of certain users.

- Any relation that is not of the conceptual model but is made visible to a user as a "virtual relation" is called a **view**.

# View Definition

- A view is defined using the **create view** statement which has the form

    **create view** *v* **as** < query expression >

    where <query expression> is any legal SQL expression.  The view name is represented by *v*.

- Once a view is defined, the view name can be used to refer to the virtual relation that the view generates.

- View definition is not the same as creating a new relation by evaluating the query expression

    - Rather, a view definition causes the saving of an expression; the expression is substituted into queries using the view.

- A view of instructors without their salary

  **create view *faculty* as**
      **select** *ID*, *name*, *dept_name*
      **from** *instructor*

- Find all instructors in the Biology department

  **select** *name*
  **from *faculty***
  **where** *dept_name* = 'Biology'

- Create a view of department salary totals

  **create view *departments_total_salary(dept_name, total_salary)* as**
      **select** *dept_name*, **sum** (*salary*)
      **from** *instructor*
      **group by** *dept_name*;

- One view may be used in the expression defining another view

- A view relation $v_1$ is said to **depend directly** on a view relation $v_2$ if $v_2$ is used in the expression defining $v_1$

- A view relation $v_1$ is said to **depend on** view relation $v_2$ if either $v_1$ depends directly to $v_2$ or there is a path of dependencies from $v_1$ to $v_2$

- A view relation $v$ is said to be **recursive** if it depends on itself.

- **create view *physics_fall_2017* as**
  **select** *course.course_id*, *sec_id*, *building*, *room_number*
  **from** *course*, *section*
  **where** *course.course_id = section.course_id*
        **and** *course.dept_name* = 'Physics'
        **and** *section.semester* = 'Fall'
        **and** *section.year* = '2017';

- **create view *physics_fall_2017*_*watson* as**
  **select** *course_id*, *room_number*
  **from** *physics_fall_2017*
  **where** *building*= 'Watson';

- Expand the view :

  **create view *physics_fall_2017_watson*  as**
      **select** *course_id*, *room_number*
      **from** *physics_fall_2017*
      **where** *building*= 'Watson'

- To:

                    **create view *physics_fall_2017_watson* as**
                        **select** *course_id*, *room_number*
                        **from** (**select** *course.course_id*, *building*, *room_number*
                            **from** *course*, *section*
                            **where** *course.course_id* = *section.course_id*
                                **and** *course.dept_name* = 'Physics'
                                **and** *section.semester* = 'Fall'
                                **and** *section.year* = '2017')
                      **where** *building*= 'Watson';

- A way to define the meaning of views defined in terms of other views.

- Let view $v_1$ be defined by an expression $e_1$ that may itself contain uses of view relations.

- View expansion of an expression repeats the following replacement step:

  **repeat**
         Find any view relation $v_i$ in $e_1$
         Replace the view relation $v_i$ by the expression defining $v_i$
  **until** no more view relations are present in $e_1$

- As long as the view definitions are not recursive, this loop will terminate

# Materialized Views

- Certain database systems allow view relations to be physically stored.

  - Physical copy created when the view is defined.

  - Such views are called **Materialized view**:

- If relations used in the query are updated, the materialized view result becomes out of date

  - Need to **maintain** the view, by updating the view whenever the underlying relations are updated.

# Update of a View

- Add a new tuple to *faculty* view which we defined earlier

  **insert into** *faculty*

  **values** ('30765', 'Green', 'Music');

- This insertion must be represented by the insertion into the *instructor* relation

  - Must have a value for salary.

- Two approaches

  - Reject the insert

  - Insert the tuple

    ('30765', 'Green', 'Music', null)

  into the *instructor* relation

- **create view** *instructor_info* **as**
  **select** *ID*, *name*, *building*
   **from** *instructor*, *department*
   **where** *instructor.dept_name* = *department.dept_name*;

- **insert into** *instructor_info*

  **values** ('69987', 'White', 'Taylor');

- Issues

  - Which department, if multiple departments in Taylor?

  - What if no department is in Taylor?

- **create view** *history_instructors* **as**
  **select** *
  **from** *instructor*
  **where** *dept_name*= 'History';

- What happens if we insert

  ('25566', 'Brown', 'Biology', 100000)

  into *history_instructors?*

# View Updates in SQL

- Most SQL implementations allow updates only on simple views

  - The **from** clause has only one database relation.

  - The **select** clause contains only attribute names of the relation, and does not have any expressions, aggregates, or **distinct** specification.

  - Any attribute not listed in the **select** clause can be set to null

  - The query does not have a **group** by or **having** clause.

# Built-in Data Types in SQL

- **date:** Dates, containing a (4 digit) year, month and date

  - Example: **date** '2005-7-27'

- **time:** Time of day, in hours, minutes and seconds.

  - Example: **time** '09:00:30'     **time** '09:00:30.75'

- **timestamp:** date plus time of day

  - Example: **timestamp** '2005-7-27 09:00:30.75'

- **interval:** period of time

  - Example: interval '1' day

  - Subtracting a date/time/timestamp value from another gives an interval value

  - Interval values can be added to date/time/timestamp values

- Large objects (photos, videos, CAD files, etc.) are stored as a *large object*:

  - **blob**: binary large object -- object is a large collection of uninterpreted binary data (whose interpretation is left to an application outside of the database system)

  - **clob**: character large object -- object is a large collection of character data

- When a query returns a large object, a pointer is returned rather than the large object itself.

- **create type** construct in SQL creates user-defined type

  **create type** *Dollars* **as numeric (12,2) final**

- Example:

  **create table** *department*
  (*dept_name* **varchar** (20),
  *building* **varchar** (15),
  *budget Dollars*);

- **create domain** construct in SQL-92 creates user-defined domain types

  **create domain** *person_name* **char**(20) **not null**

- Types and domains are similar.  Domains can have constraints, such as **not null**, specified on them.

- Example:

  **create domain** *degree_level* **varchar**(10)
      **constraint** *degree_level_test*
          **check** (**value in** ('Bachelors', 'Masters', 'Doctorate'));

# *Data Cleaning

- [Extract, transform, load – Wikipedia](#)

- Benefits

  - 1. The original table remains unchanged, allowing for rollback at any time.

  - 2. The cleaning logic is entirely contained within the VIEW, allowing version control with a SQL file.

  - 3. If you want to add new rules, simply wrap it in another VIEW; the old VIEW remains available for historical tasks.

  - 4. If you want to materialize the data later, simply create a materialized view mv_clean AS select * from v2_orders_valid; to speed things up.

  - This is the typical "using VIEWs for data cleaning" in PostgreSQL: dirty data remains in place, while clean data is "projected" as needed.

# *Comparison of Tools

|  | Easy to use | Performance | Memory | Scalability |
|---|---|---|---|---|
| Linux Command Line | Not good | Good | Based on disks | No |
| Python/Pandas | Not bad | Slow | Based on main mem | No |
| SQL/DB | Good | Slower than Linux, faster than Pandas | Based on disks | Hadoop can do it best |
| GUI Tools | Great for small tasks | Depends on the underlying implementation | Depends on the underlying implementation | Depends on the underlying implementation |

- What is 'legal moves'?

  - Where a and where b ⇔ Where b and where a : **commutative law**

  - Nature Join satisfy **association law**

  - Exchange join and where: **distributive law**

  - …

- It is similar to what we have learned in math…

# **Modeling with a kind of… *Algebra*??**

- Informally speaking, algebra usually refers to something supporting 'add' and 'multiply', as well as their inversion ('subtract' and 'divide')

|  | **+** | **✖** |
|---|---|---|
| Elementary algebra (for children) | add | multiply |
| Polynomial Algebra | add | multiply |

- How many algebras do you know?

# **Relational Algebra

| | **+** | **✕** |
|---|---|---|
| Boolean algebra (in computer science) | or | and |
| Linear algebra | add | product |
| Sigma algebra (in probilitstic theory) | union | interset |
| Relational Algebra | Union | Join |

- Example (Not strictly speaking)
  - Two tables with the same schema can be united
  - Two tables with common keys can be Joined

# ***Application

- Find all equivalent queries

  - Traverse the equivalent form of an expression according to the transformation rules. Then a selection of the best form can be conducted.

- Query = Solve Boolean equations

  - Combined with other types of equations (such as linear equations), more complex queries can be implemented, such as: Q: Under what conditions can Alice receive a scholarship? A: improve GPA at least 0.3.

- Automatic normal form conversion

  - Polynomial factorization minimizes redundancy under a multiplication structure; database normal form decomposition minimizes redundancy under a connection structure. The normal form can be converted automatically.