# Part 8: SQL II

**References:**

- Elmasri/Navathe:Fundamentals of Database Systems, 3rd Edition, 1999.
  Chap. 8, "SQL — The Relational Database Standard" (Sect. 8.2, 8.3.3, part of 8.3.4.)

- Silberschatz/Korth/Sudarshan: Database System Concepts, 3rd Edition.
  McGraw-Hill, 1999: Chapter 4: "SQL".

- Kemper/Eickler: Datenbanksysteme (in German), Ch. 4, Oldenbourg, 1997.

- Lipeck: Skript zur Vorlesung Datenbanksysteme (in German), Univ. Hannover, 1996.

- Heuer/Saake: Datenbanken, Konzepte und Sprachen (in German), Thomson, 1995.

- Date/Darwen: A Guide to the SQL Standard, Fourth Edition, Addison-Wesley, 1997.

- Date: A Guide to the SQL Standard, First Edition, Addison-Wesley, 1987.

- van der Lans: SQL, Der ISO-Standard (in German). Hanser, 1990.

- Sunderraman: Oracle Programming, A Primer. Addison-Wesley, 1999.

- Oracle 8i SQL Reference, Release 2 (8.1.6), Dec. 1999, Part No. A76989-01.

- Chamberlin: A Complete Guide to DB2 Universal Database. Morgan Kaufmann, 1998.

- Microsoft SQL Server Books Online: Accessing and Changing Data.

- Microsoft Jet Database Engine Programmer's Guide, 2nd Edition (Part of MSDN Library Visual Studio 6.0).

- DuBois: MySQL. New Riders Publishing, 2000, ISBN 0-7357-0921-1, 756 pages.

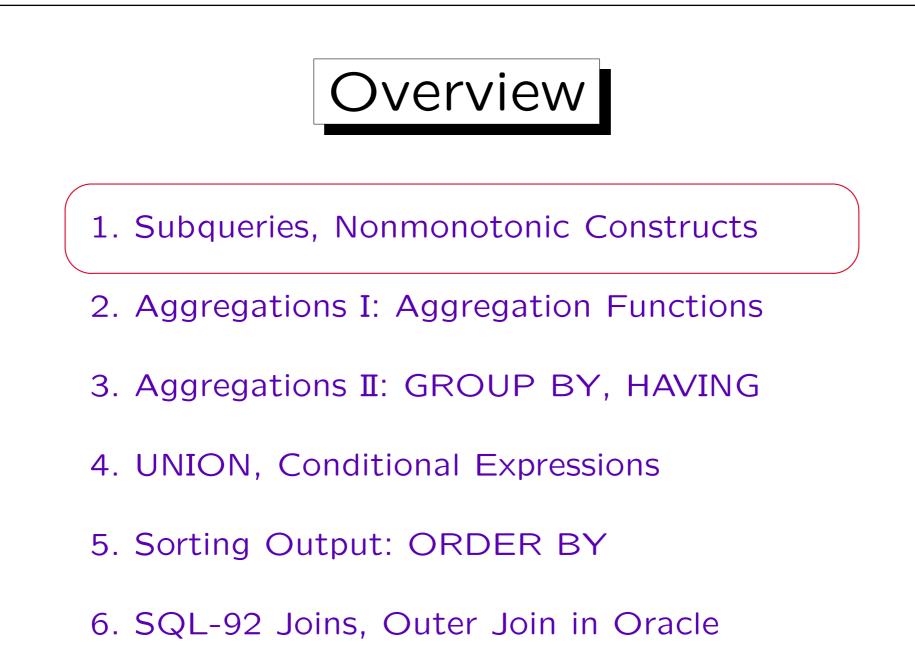# Objectives

After completing this chapter, you should be able to:

- write advanced queries in SQL including

  aggregations, subqueries, and UNION.

- enumerate and explain the clauses of an SQL query.

  SELECT, FROM, WHERE, GROUP BY, HAVING, ..., ORDER BY

- explain joins in SQL-92.

- evaluate the correctness of a given query.

- evaluate the portability of certain constructs.

# Overview

# Example Database (again)

### STUDENTS

| SID | FIRST | LAST | EMAIL |
|-----|-------|------|-------|
| 101 | Ann | Smith | ⋯ |
| 102 | Michael | Jones | (null) |
| 103 | Richard | Turner | ⋯ |
| 104 | Maria | Brown | ⋯ |

### EXERCISES

| CAT | ENO | TOPIC | MAXPT |
|-----|-----|-------|-------|
| H | 1 | Rel. Algeb. | 10 |
| H | 2 | SQL | 10 |
| M | 1 | SQL | 14 |

### RESULTS

| SID | CAT | ENO | POINTS |
|-----|-----|-----|--------|
| 101 | H | 1 | 10 |
| 101 | H | 2 | 8 |
| 101 | M | 1 | 12 |
| 102 | H | 1 | 9 |
| 102 | H | 2 | 9 |
| 102 | M | 1 | 10 |
| 103 | H | 1 | 5 |
| 103 | M | 1 | 7 |

# Nonmonotonic Behaviour (1)

- SQL queries using only the constructs introduced above compute monotonic functions on the existing tables: If further rows are inserted, one gets at least the same answers as before, and maybe more.

- However, not all queries behave monotonically in this way: E.g. print students who have not yet sub-mitted any homework.

  Currently Maria Brown would be a correct answer. But if a homework result were inserted for her, she would no longer qualify.

- Therefore, this query cannot be formulated with the SQL constructs that were introduced so far.

# Nonmonotonic Behaviour (2)

- In the natural language version of queries, formulations like "there is no", "does not exist" indicate nonmonotonic behaviour.

- Furthermore, "for all", "the minimal/maximal", also indicate nonmonotonic behaviour: In this case a violation of the "for all" condition must not exist.

    For some such queries, a formulation with aggregations (HAVING) might be natural, see below.

- When formulating queries in SQL, it is important to check whether the query requires that certain tuples do not exist.

# NOT IN (1)

- With IN (∈) and NOT IN (∉) it is possible to check whether an attribute value appears in a set that is computed by another SQL query.

- E.g. students without any homework result:

```
SELECT FIRST, LAST
FROM   STUDENTS
WHERE  SID NOT IN (SELECT SID
                   FROM   RESULTS
                   WHERE  CAT = 'H')
```

| FIRST | LAST  |
|-------|-------|
| Maria | Brown |

# NOT IN (2)

- At least conceptually, the subquery is evaluated, before the execution of the main query starts:

| STUDENTS | | | |
|---|---|---|---|
| SID | FIRST | LAST | EMAIL |
| 101 | Ann | Smith | ... |
| 102 | Michael | Jones | (null) |
| 103 | Richard | Turner | ... |
| 104 | Maria | Brown | ... |

| Result of Subquery |
|---|
| SID |
| 101 |
| 101 |
| 102 |
| 102 |
| 103 |

- Then for every STUDENTS tuple, a matching SID is searched in the subquery result. If there is none, the student name is printed.

# NOT IN (3)

- It is possible to use DISTINCT in the subquery:

```
SELECT FIRST, LAST
FROM   STUDENTS
WHERE  SID NOT IN (SELECT DISTINCT SID    ?
                   FROM    RESULTS
                   WHERE   CAT = 'H')
```

- This is logically equivalent, and the effect on the performance depends on the data and the DBMS.

  I would expect that a reasonable optimizer knows that duplicates are not important in this case and that conversely writing DISTINCT might have the effect that the optimizer does not consider certain evaluation stragegies that do not really materialize the result of the subquery.

# NOT IN (4)

- It is also possible to use `IN` (without `NOT`) for an element test.

- This is relatively seldom done, since it is equivalent to a join, which could be written without a subquery.

- But sometimes this formulation is more elegant. It might also help to avoid duplicates.

  Or to get exactly the required duplicates (see example on next page).

# NOT IN (5)

- E.g. topics ("names") of homeworks that were already solved by at least one student:

```
SELECT TOPIC
FROM   EXERCISES
WHERE  CAT='H' AND ENO IN (SELECT ENO
                           FROM   RESULTS
                           WHERE  CAT='H')
```

- Exercise: Is there a difference to this query (with or without DISTINCT)?

```
SELECT DISTINCT TOPIC
FROM   EXERCISES E, RESULTS R
WHERE  E.CAT='H' AND E.ENO=R.ENO AND R.CAT='H'
```

# NOT IN (6)

- In SQL-86, the subquery on the right-hand side of IN must have a single output column.

    So that the subquery result is really a set (or multiset), and not an arbitrary relation.

- In SQL-92, comparisons were extended to the tuple level, and therefore it is possible to write e.g.

```
WHERE (FIRST, LAST) NOT IN (SELECT FIRST, LAST
                            FROM   ...)
```

    But is not very portable. E.g. SQL Server and Access do not support it (and MySQL does not permit any subqueries, see below). An EXISTS subquery (see below) might be better if one has to compare more than one column. Oracle and DB2 do allow IN with multiple columns.

# NOT IN (7)

Atomic Formula (Form 6):



- The Subquery must result in a table with a single column (a set).

- However, in SQL-92, Oracle, and DB2 it is possible to write a tuple on the left hand side in the form ($Term_1$, ..., $Term_n$). Then the subquery must result in a table with exactly $n$ columns.

- MySQL does not support subqueries.

- The column names on the left and right hand side of IN do not have to match, but the data types must be compatible.

# NOT EXISTS (1)

- It is possible to check in the outer query whether the result of the subquery is empty (NOT EXISTS).

- In the inner query, tuple variables declared in the FROM clause of the outer query can be accessed.

  > This is actually also possible for IN subqueries, but there it is an unnecessary and unexpected complication (bad style).

- This means that the subquery has to be evaluated once for every assignment of values to the accessed tuple variables in the outer query. The subquery can be seen as parameterized.

# NOT EXISTS (2)

- Students that have not submitted any homework:

```
SELECT FIRST, LAST
FROM    STUDENTS S
WHERE   NOT EXISTS (SELECT * FROM RESULTS R
                           WHERE  R.CAT = 'H'
                           AND    R.SID = S.SID )
```

- The tuple variable S loops over the four rows in STUDENTS. Conceptually, the subquery is evaluated four times. Each time, S.SID is replaced by the SID value of the current tuple S.

    The DBMS is free to choose another, more efficient evaluation strategy if that evaluation strategy is guaranteed to give the same result.

# NOT EXISTS (3)

- First, S points to the STUDENTS tuple

| SID | FIRST | LAST | EMAIL |
|-----|-------|------|-------|
| 101 | Ann | Smith | $\cdots$ |

- S.SID in the subquery is conceptually replaced by 101 and the following query is executed:

```
SELECT * FROM RESULTS R
WHERE  R.CAT = 'H'
AND    R.SID = 101
```

| SID | CAT | ENO | POINTS |
|-----|-----|-----|--------|
| 101 | H | 1 | 10 |
| 101 | H | 2 | 8 |

- The result is not empty. Thus, the NOT EXISTS condition in the outer query is not satisfied for this S.

# NOT EXISTS (4)

- The same happens for the second row in STUDENTS.

  The subquery is executed for S.SID = 102:

```
SELECT * FROM RESULTS R
WHERE   R.CAT = 'H'
AND     R.SID = 102
```

| SID | CAT | ENO | POINTS |
|-----|-----|-----|--------|
| 102 | H   | 1   | 9      |
| 102 | H   | 2   | 9      |

- The result is not empty, therefore the NOT EXISTS condition is not satisfied.

- Also for the third row in STUDENTS, the condition is not satisfied.

# NOT EXISTS (5)

- Finally, S points to the STUDENTS tuple

| SID | FIRST | LAST | EMAIL |
|-----|-------|------|-------|
| 104 | Maria | Brown | ... |

- For S.SID=104, the result of the subquery is empty:

```
SELECT * FROM RESULTS R
WHERE   R.CAT = 'H'              no rows selected
AND     R.SID = 104
```

- Thus, the NOT EXISTS condition is satisfied for this tuple S. Maria Brown is printed as the query result.

# NOT EXISTS (6)

- While in the inner query, tuple variables from the outer query can be accessed, the converse is illegal:

```
SELECT  FIRST, LAST, R.ENO   Wrong!
FROM    STUDENTS S
WHERE   NOT EXISTS (SELECT * FROM RESULTS R
                          WHERE  R.CAT = 'H'
                          AND    R.SID = S.SID)
```

- This works like global and local variables: Variables defined in the outer query are valid for the entire query, variables defined in the subquery are valid only in the subquery ($\sim$ block structure in Pascal).

# NOT EXISTS (7)

- Subqueries that access variables from the outer query are called "correlated subqueries".

  Correlated subqueries can be understood as being parameterized with the tuples chosen in the outer query. There can be optimizations, but conceptually they are executed once for every assignment of tuples to the tuple variables in the outer query.

- Subqueries that do not access variables from the outer query are called "non-correlated subqueries".

  It suffices to evaluate a non-correlated subquery only once (since the result does not depend on the tuples chosen for the tuple variables of the outer query).

# NOT EXISTS (8)

- Non-correlated subqueries with `NOT EXISTS` are al-
  most always an error (but they are ok with `IN`):

```
SELECT FIRST, LAST    Wrong!
FROM    STUDENTS S
WHERE   NOT EXISTS (SELECT * FROM RESULTS R
                         WHERE  CAT = 'H')
```

  Here the join-condition in the subquery was forgotten, and it became
  a non-correlated subquery.

- If there is at least one homework entry in `RESULTS`,
  no matter for what student, the `NOT EXISTS` will be
  false, and the query result empty.

# NOT EXISTS (9)

- Until now, for attribute references without tuple variable ("unqualified attribute name"), there had to be a unique tuple varible to which it can refer.

- For subqueries, SQL only requires that there is a unique nearest tuple variable which has this attribute, e.g. this is legal (but bad style):

```
SELECT FIRST, LAST
FROM   STUDENTS S
WHERE  NOT EXISTS (SELECT * FROM RESULTS R
                          WHERE  CAT = 'H'
                          AND    SID  = S.SID)
```

# NOT EXISTS (10)

- In general, for attribute reference without tuple variables, the SQL parser searches the `FROM`-clauses beginning from the current subquery towards outer queries (there can be several nesting levels).

- The first `FROM`-clause that declares a tuple variable with this attribute must have exactly one such variable. Then the attribute refers to this variable.

- This rule helps that non-correlated subqueries can be developed independently and inserted into another query without any change (so it makes sense).

# NOT EXISTS (11)

- It is also legal to declare tuple variables in the subquery that have the same name as tuple variables in the outer query.

  ```
  SELECT FIRST, LAST
  FROM    STUDENTS X
  WHERE   NOT EXISTS (SELECT * FROM RESULTS X
                                  WHERE ???)
  ```

- References to X in the subquery mean RESULTS X. The variable declared in the outer query becomes "shadowed": It cannot be accessed in the subquery.

# NOT EXISTS (12)

- It is legal to specify a `SELECT`-list in the subquery, but since for `NOT EXISTS` the returned columns do not matter, "`SELECT *`" should be used.

- Some authors say that in some systems `SELECT null` or `SELECT 1` is actually faster than `SELECT *`.

  > "`SELECT null`" is used by Oracle's programmers (in "`catalog.sql`"). But this does not work in DB2 (null cannot be used as a term here). Today, resonably good optimizers should know that the column values are not really needed, and the `SELECT`-list should not matter, not even for performance.

# NOT EXISTS (13)

Atomic Formula (Form 7):

$$\longrightarrow \boxed{\text{EXISTS}} \rightarrow \boxed{(} \rightarrow \boxed{\text{Subquery}} \rightarrow \boxed{)} \longrightarrow$$

- A subquery is an expression of the form SELECT ...
  FROM ... [WHERE ...] [GROUP BY ...] [HAVING ...].

  [...] means that these parts are optional. SQL-92 also allows UNION (see below) in subqueries (as do Oracle, DB2, and SQL Server), SQL-86 does not (and Access really does not support it).

- ORDER BY is not allowed in subqueries.

  It would make no sense there, it is only for the final output.

- Subqueries must be enclosed in parentheses (...).

# NOT EXISTS (14)

- It is possible to use EXISTS without negation.

- Who has submitted at least one homework?

```
SELECT  SID, FIRST, LAST
FROM    STUDENTS S
WHERE   EXISTS (SELECT * FROM RESULTS R
                      WHERE   R.SID = S.SID
                      AND     R.CAT = 'H')
```

- But the same query can be done with a usual join:

```
SELECT DISTINCT S.SID, FIRST, LAST
FROM    STUDENTS S, RESULTS R
WHERE   S.SID = R.SID AND R.CAT = 'H'
```

# "For all" (1)

- Who got the best result for Homework 1?

```
SELECT FIRST, LAST, POINTS
FROM   STUDENTS S, RESULTS X
WHERE  S.SID = X.SID
AND    X.CAT = 'H' AND X.ENO = 1
AND    NOT EXISTS
       (SELECT * FROM RESULTS Y
         WHERE  Y.CAT = 'H' AND Y.ENO = 1
         AND    Y.POINTS > X.POINTS)
```

- I.e. a result X for Homework 1 is selected if there is no result Y for this exercise with more points than X.

# "For all" (2)

- In mathematical logic, there are two quantifiers:

  ◇ $\exists s X\colon F$: There is an $X$ that satisfies $F$.
  (existential quantifier)

  ◇ $\forall s X\colon F)$: For all $X$, $F$ is true.
  (universal quantifier)

- In tuple relational calculus, the maximal number of points for Homework 1 is expressed e.g. as follows:

$$\{\text{X.POINTS [RESULTS X]} \mid \text{X.CAT} = \text{'H'} \land \text{X.ENO} = 1 \land$$
$$\forall \text{RESULTS Y: Y.CAT} = \text{'H'} \land \text{Y.ENO} = 1$$
$$\rightarrow \text{Y.POINTS} \leq \text{X.POINTS})\}$$

# "For all" (3)

- The pattern $\forall s\, X\colon (F_1 \rightarrow F_2)$ is very typical:
  $F_2$ must be true for all $X$ that satisfy $F_1$.

- SQL has only an existential quantifier ("EXISTS"),
  but not a universal quantifier.

    However, see ">= ALL" below.

- This is no problem, because $\forall s\, X\colon F$ is equivalent
  to $\neg\exists s\, X\colon \neg F$. One type of quantifier suffices.

    "$F$ is true for all $X$" is the same as "$F$ is false for no $X$".

- The above pattern is equivalent to $\neg\exists s\, X\colon F_1 \wedge \neg F_2$.

# "For all" (4)

- The above example is logically equivalent to:

$$\{\text{X.POINTS } [\text{RESULTS X}] \mid \text{X.CAT} = \text{'H'} \wedge \text{X.ENO} = 1 \wedge$$
$$\neg\exists \text{RESULTS Y: Y.CAT} = \text{'H'} \wedge \text{Y.ENO} = 1$$
$$\wedge \text{Y.POINTS} > \text{X.POINTS})\}$$

- In SQL, this is written as:

```
SELECT  X.POINTS
FROM    RESULTS X
WHERE   X.CAT = 'H' AND X.ENO = 1
AND     NOT EXISTS
        (SELECT * FROM RESULTS Y
          WHERE  Y.CAT = 'H' AND Y.ENO = 1
          AND    Y.POINTS > X.POINTS)
```

# Nested Subqueries

- Subqueries can be nested to any reasonable depth.

- List the students who solved all homeworks:

```
SELECT FIRST, LAST
FROM   STUDENTS S
WHERE  NOT EXISTS
          (SELECT * FROM EXERCISES E
            WHERE  CAT = 'H'
            AND    NOT EXISTS
                     (SELECT * FROM RESULTS R
                       WHERE  R.SID = S.SID
                       AND    R.ENO = E.ENO
                       AND    R.CAT = 'H'))
```

# Common Errors (1)

Exercises:

- Would this query find students without homeworks in the database? If not, what does it compute?

```
SELECT DISTINCT S.SID, FIRST, LAST
FROM    STUDENTS S, RESULTS R
WHERE   S.SID <> R.SID AND R.CAT = 'H'
```

- Would this query find exercises that were not yet solved?

```
SELECT DISTINCT E.CAT, E.ENO
FROM    EXERCISES E, RESULTS R
WHERE   E.CAT = R.CAT AND E.ENO = R.ENO
AND     R.SID IS NULL
```

# Common Errors (2)

- It is important to understand that the absence/non-existence of a row is very different than the existence of a row with a different value.

    If the requested query behaves in a non-monotonic fashion (i.e. insertion of a row could invalidate an answer), then `NOT EXISTS`, `NOT IN`, `<> ALL` etc. are required.

- There is no way to write it without a subquery.

    Except possibly using an outer join. Aggregations also change when tuples are inserted, but without subquery, they cannot express "for all" or "not exists".

# Common Errors (3)

- Does this query compute the student with the best result for Homework 1?

```
SELECT DISTINCT S.FIRST, S.LAST, X.POINTS
FROM    STUDENTS S, RESULTS X, RESULTS Y
WHERE   S.SID = X.SID
AND     X.CAT = 'H' AND X.ENO = 1
AND     Y.CAT = 'H' AND Y.ENO = 1
AND     X.POINTS > Y.POINTS
```

- If not, what does it compute?

# Common Errors (4)

- As mentioned above, using a non-correlated sub-query with `NOT EXISTS` is normally an error.

- Does this also apply in this case (there is a join condition in the subquery)?

```
SELECT FIRST, LAST
FROM   STUDENTS S
WHERE  NOT EXISTS              Wrong!
          (SELECT *
            FROM   RESULTS R, STUDENTS S
            WHERE  R.SID = S.SID
            AND    R.CAT = 'H' AND R.ENO = 1)
```

# Common Errors (5)

- What is the error in this query? It is supposed to find students that have neither submitted a home-work nor participated in an exam.

```
SELECT FIRST, LAST      Wrong!
FROM   STUDENTS S
WHERE  SID NOT IN (SELECT SID
                   FROM   EXERCISES)
```

- This query is syntactically correct SQL. Why?

- What is the output of the query?

    Under the assumption that EXERCISES is not empty.

# Common Errors (6)

- Is there any problem with this query?

  The task is to list all students who did not yet actively participated in the course, i.e. neither submitted a homework nor took the exam.

```
SELECT S.FIRST, S.LAST
FROM   STUDENTS S, RESULTS R
WHERE  S.SID = R.SID
AND    NOT EXISTS (SELECT *
                        FROM   RESULTS R
                        WHERE  S.SID = R.SID)
```

# ALL, ANY, SOME (1)

- It is possible to compare a value with all values in a set (computed by a subquery).

- One can require that the comparison returns true for all set elements (ALL) or for at least one (ANY):

```
SELECT  S.FIRST, S.LAST, X.POINTS
FROM    STUDENTS S, RESULTS X
WHERE   S.SID=X.SID AND X.CAT='H' AND X.ENO=1
AND     X.POINTS >= ALL (SELECT Y.POINTS
                         FROM   RESULTS Y
                         WHERE  Y.CAT = 'H'
                         AND    Y.ENO = 1)
```

# ALL, ANY, SOME (2)

- The following is logically equivalent to the above query:

```
SELECT  S.FIRST, S.LAST, X.POINTS
FROM    STUDENTS S, RESULTS X
WHERE   S.SID=X.SID AND X.CAT='H' AND X.ENO=1
AND     NOT X.POINTS < ANY (SELECT Y.POINTS
                            FROM    RESULTS Y
                            WHERE   Y.CAT = 'H'
                            AND     Y.ENO = 1)
```

- Again, "for all" can be replaced by "not exists not".

  Of course, also conversely "exists" is equivalent to "not for all not".

# ALL, ANY, SOME (3)

- This construct is not strictly necessary, since e.g.

$$T_1 < \text{ANY (SELECT } T_2 \text{ FROM ... WHERE ...)}$$

  is equivalent to

  EXISTS (SELECT * FROM ... WHERE ... AND $T_1 < T_2$)

    This requires that $T_1$ explicitly mentions a tuple variable which is not redeclared in the subquery (so that the meaning of $T_1$ is not changed by moving it into the subquery).

- E.g. Oracle internally does such transformations so that the query optimizer does not have to handle too many different cases (syntactic variants).

# ALL, ANY, SOME (4)

Atomic Formula (Form 8):

# ALL, ANY, SOME (5)

Syntactic Remarks:

- `ANY` and `SOME` are synonyms.

- "`x IN S`" is equivalent to "`x = ANY S`".

- The subquery must have a single result column.

  > SQL92 allows comparisons also on a tuple basis. Oracle supports this only with `<>` and `=`, DB2 supports only `=ANY` (which is equivalent to `IN`). SQL86, SQL Server, and Access do not support tuple comparisons.

- If none of the keywords `ALL`, `ANY`, `SOME` are present, the subquery must yield at most one row.

  > Since there is also only one column, this means the subquery gives a single data value. If the subqery result is empty, the null value is used.

# Single Value Subqueries (1)

- Who got full points for Homework 1?

```
SELECT  S.FIRST, S.LAST
FROM    STUDENTS S, RESULTS R
WHERE   S.SID=R.SID AND R.CAT='H' AND R.ENO=1
AND     R.POINTS = (SELECT MAXPT
                    FROM   EXERCISES
                    WHERE  CAT='H' AND ENO=1)
```

- It is only possible to leave out ANY/ALL when the subquery is guaranteed to return at most one row.

  In the example, a key of EXERCISES is specified. But in general, this may depend on the data. The query might run during testing, but later give an error. Use constraints to ensure that the necessary assumtions are really satisfied.

# Single Value Subqueries (2)

- In SQL92, DB2, SQL Server, and Access, a sub-query returning a single data element can be used as a term/expression. Thus, this is equally legal:

    `(SELECT MAXPT FROM ...) = R.POINTS`

- In Oracle8 and SQL86, the subquery must be on the right hand side.

- One can even do further computations with the result of a subquery, e.g. (not in SQL86, Oracle8):

    `R.POINTS >= (SELECT MAXPT FROM ...) * 0.9`

# Single Value Subqueries (3)

- If the subquery has an empty result, the null value is used instead.

- E.g. this is a strange way to ask for students that have not yet solved Homework 1:

```
SELECT FIRST, LAST
FROM   STUDENTS S
WHERE  (SELECT 1                    Bad Style!
           FROM   RESULTS R
           WHERE  R.SID = S.SID
           AND    R.CAT = 'H' AND R.ENO = 1) IS NULL
```

- In SQL86 and Oracle8, this is a syntax error.

# Subqueries under FROM (1)

- Since the result of an SQL-query is a table, it is natural that one can write a subquery instead of a table name in the FROM-clause.

- This was not allowed in SQL-86, and at that time SQL was often criticized as having "not orthogonal constructs", which cannot be combined arbitrarily.

  In relational algebra, wherever one can write a relation name, one can also write a subquery (relational algebra expression).

- Subqueries under FROM are really needed only seldom, and might make the query more complex.

# Subqueries under FROM (2)

- Subqueries under `FROM` are needed e.g. for nested aggregations, see below.

- In the following example, the join of `RESULTS` and `EXERCISES` is computed in a subquery (that might result from a view definition, see below):

```
SELECT X.SID, ROUND(X.POINTS*100/X.MAXPT) AS PCT
FROM    (SELECT E.CAT, E.ENO, R.SID, R.POINTS,
                 E.MAXPT
         FROM   EXERCISES E, RESULTS R
         WHERE  E.CAT=R.CAT AND E.ENO=R.ENO) X
WHERE  X.CAT = 'H' AND X.ENO = 1
```

# Subqueries under FROM (3)

- SQL92, SQL Server, and DB2 require declaring a tuple variable for the subquery; in Oracle and Access this is optional.

- SQL92, DB2, and SQL Server (but not Oracle8 and Access) permit to rename columns in this way:

```
FROM (...) X(CATEGORY, EX_NO, ...)
```

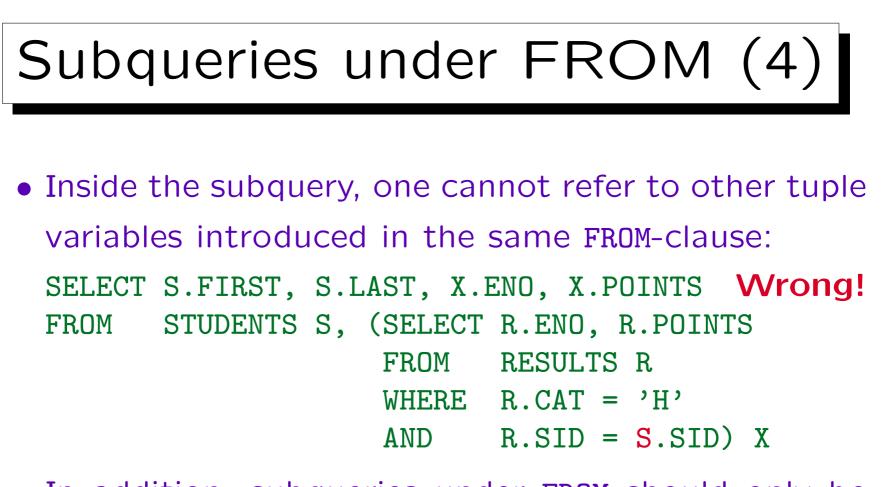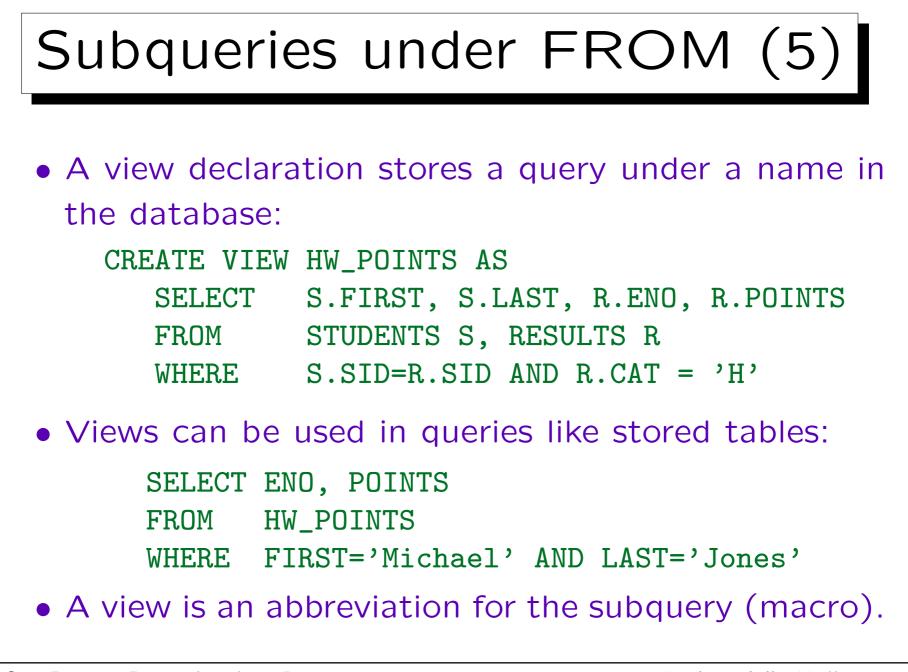- In Oracle and Access, columns can only be renamed inside the subquery.

    All systems support the specification of new column names in the SELECT-clause, so that is the more portable way.

# Subqueries under FROM (4)

- Inside the subquery, one cannot refer to other tuple variables introduced in the same `FROM`-clause:

```
SELECT S.FIRST, S.LAST, X.ENO, X.POINTS   Wrong!
FROM    STUDENTS S, (SELECT R.ENO, R.POINTS
                      FROM    RESULTS R
                      WHERE   R.CAT = 'H'
                      AND     R.SID = S.SID) X
```

- In addition, subqueries under `FROM` should only be used if needed. They can make queries much more difficult to understand.

# Subqueries under FROM (5)

- A view declaration stores a query under a name in the database:

```
CREATE VIEW HW_POINTS AS
    SELECT    S.FIRST, S.LAST, R.ENO, R.POINTS
    FROM      STUDENTS S, RESULTS R
    WHERE     S.SID=R.SID AND R.CAT = 'H'
```

- Views can be used in queries like stored tables:

```
SELECT ENO, POINTS
FROM   HW_POINTS
WHERE  FIRST='Michael' AND LAST='Jones'
```

- A view is an abbreviation for the subquery (macro).

# Subqueries under FROM (6)

- When a view used in a query, the DBMS simply replaces the view name by the query it stands for.

    Views existed already in SQL-86. However, since SQL-86 did not contain subqueries under FROM, there were complex restrictions for using views.

- By using views, one can build complex queries step by step.

    If the optimizer is not very good, it might be possible that a query built in this way runs slightly slower than a single "monolithic" query. However, there should be no difference to using subqueries under FROM. A performance improvement is only possible if one can formulate the query without such subqueries.
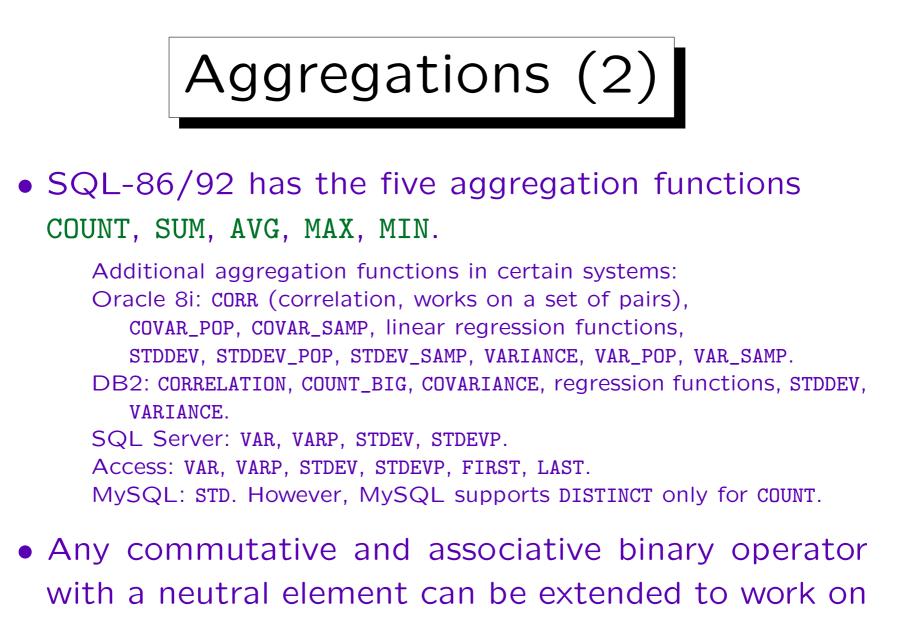
# Overview

# Aggregations (1)

- Aggregation functions are functions from a set or multiset to a single value (usually a number).

$$\text{E.g.:} \quad min\{41, 57, 19, 23, 27\} = 19$$

- Aggregation functions aggregate or summarize an entire set of values to a single value.

  Aggregation functions are also called "set functions", "group functions" or "column functions". They take not a single value as input, but an entire column (a set). The column can be constructed by means of a query (it does not have to be a column of a stored table).

- Aggregation functions are often used for statistical evaluations (e.g. average).

# Aggregations (2)

- **SQL-86/92 has the five aggregation functions**
  COUNT, SUM, AVG, MAX, MIN.

  > Additional aggregation functions in certain systems:
  >
  > Oracle 8i: CORR (correlation, works on a set of pairs),
  > COVAR_POP, COVAR_SAMP, linear regression functions,
  > STDDEV, STDDEV_POP, STDEV_SAMP, VARIANCE, VAR_POP, VAR_SAMP.
  >
  > DB2: CORRELATION, COUNT_BIG, COVARIANCE, regression functions, STDDEV,
  > VARIANCE.
  >
  > SQL Server: VAR, VARP, STDEV, STDEVP.
  >
  > Access: VAR, VARP, STDEV, STDEVP, FIRST, LAST.
  >
  > MySQL: STD. However, MySQL supports DISTINCT only for COUNT.

- **Any commutative and associative binary operator with a neutral element can be extended to work on sets. E.g. $sum$ is the set-version of $+$.**
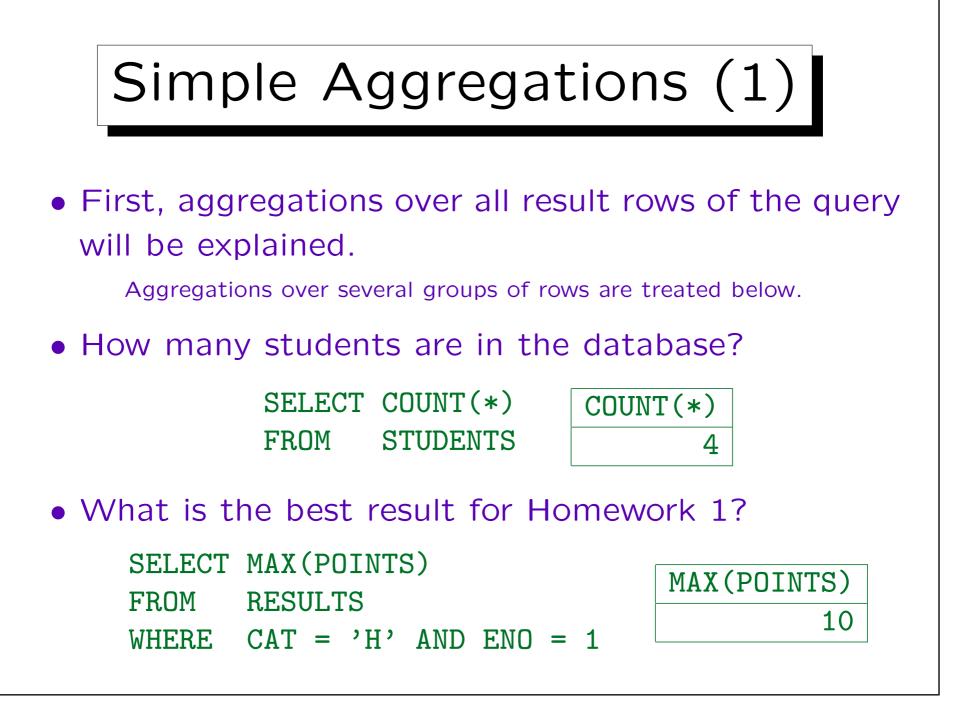
# Aggregations (3)

- Some aggregation functions are sensitive to dupli-cates (e.g. sum), others are not (e.g. minimum).
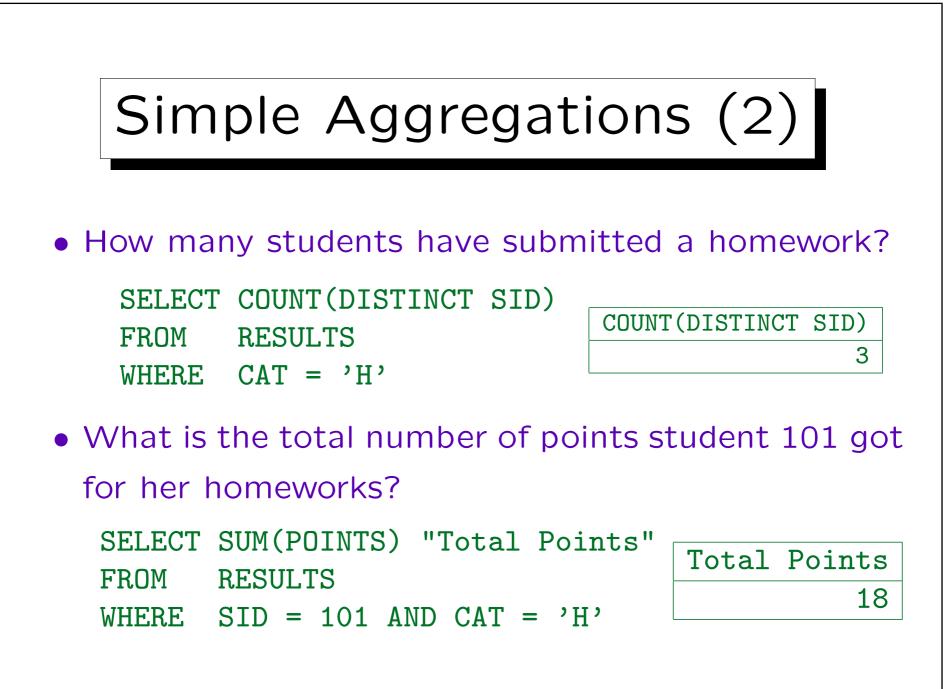
    E.g. the sum of all items of an invoice. If two items cost the same amount, nevertheless both must be added.

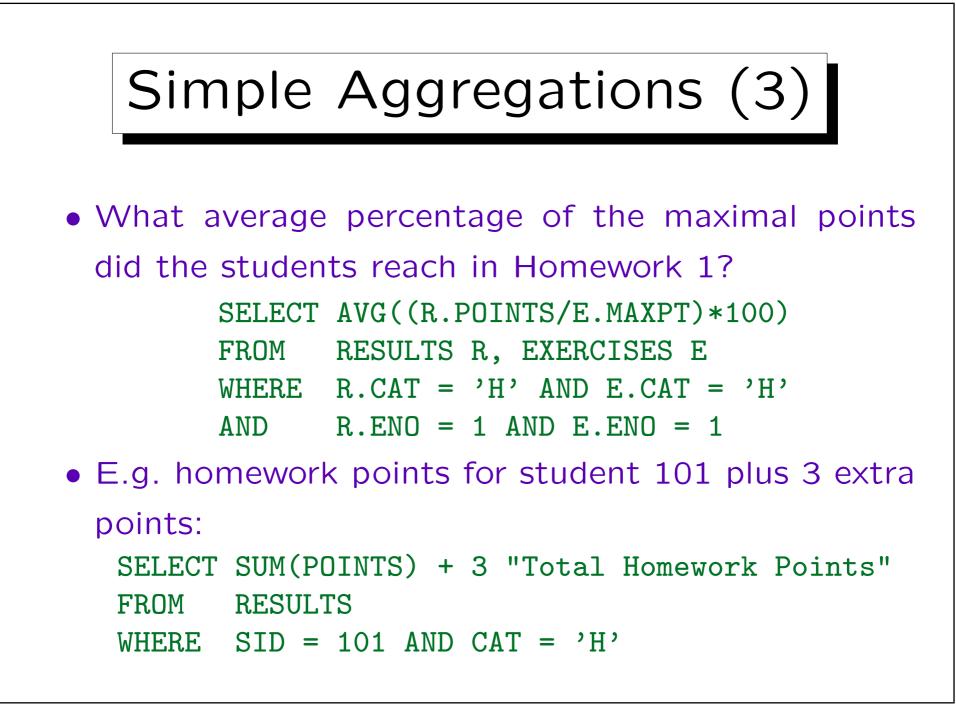- In SQL, one can request duplicate elimination (input is a set) or not (input is a multiset).

    A multiset is a set where each element has a multiplicity, e.g. an element can be contained in a multiset two times. In contrast to a list, there is still no specific order. Also the name "bag" is used.

- SUM(DISTINCT X) and AVG(DISTINCT X) are most likely an error. Some students mix up SUM and COUNT.

# Simple Aggregations (1)

- First, aggregations over all result rows of the query will be explained.

  Aggregations over several groups of rows are treated below.

- How many students are in the database?

```
SELECT COUNT(*)
FROM   STUDENTS
```

| COUNT(*) |
|---|
| 4 |

- What is the best result for Homework 1?

```
SELECT MAX(POINTS)
FROM   RESULTS
WHERE  CAT = 'H' AND ENO = 1
```

| MAX(POINTS) |
|---|
| 10 |

# Simple Aggregations (2)

- How many students have submitted a homework?

```
SELECT COUNT(DISTINCT SID)
FROM   RESULTS
WHERE  CAT = 'H'
```

| COUNT(DISTINCT SID) |
|---|
| 3 |

- What is the total number of points student 101 got for her homeworks?

```
SELECT SUM(POINTS) "Total Points"
FROM   RESULTS
WHERE  SID = 101 AND CAT = 'H'
```

| Total Points |
|---|
| 18 |

# Simple Aggregations (3)

- What average percentage of the maximal points did the students reach in Homework 1?

```
SELECT AVG((R.POINTS/E.MAXPT)*100)
FROM    RESULTS R, EXERCISES E
WHERE   R.CAT = 'H' AND E.CAT = 'H'
AND     R.ENO = 1 AND E.ENO = 1
```

- E.g. homework points for student 101 plus 3 extra points:

```
SELECT SUM(POINTS) + 3 "Total Homework Points"
FROM    RESULTS
WHERE   SID = 101 AND CAT = 'H'
```
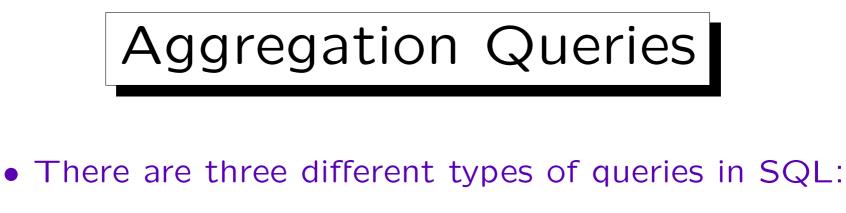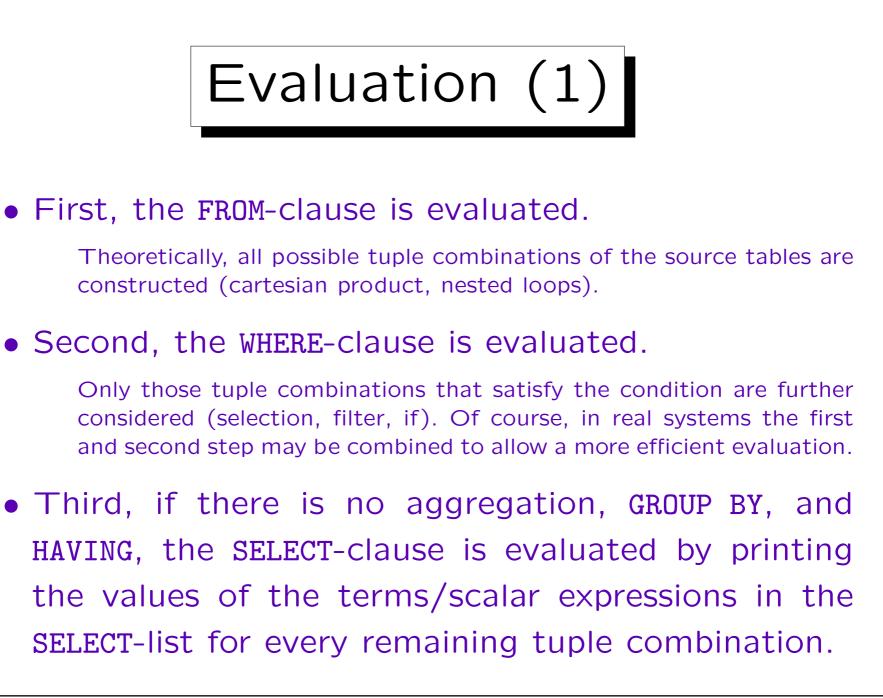
# Simple Aggregations (4)

- It is possible to compute more than one aggrega-
  tion in the `SELECT` list, e.g.: What is the minimum
  and maximum number of points for Homework 1?

```
SELECT MIN(POINTS), MAX(POINTS)
FROM   RESULTS
WHERE  CAT = 'H' AND ENO = 1
```

- The aggregations can refer to different columns:

```
SELECT COUNT(DISTINCT TOPIC), AVG(MAXPT)
FROM   EXERCISES E
```

# Aggregation Queries

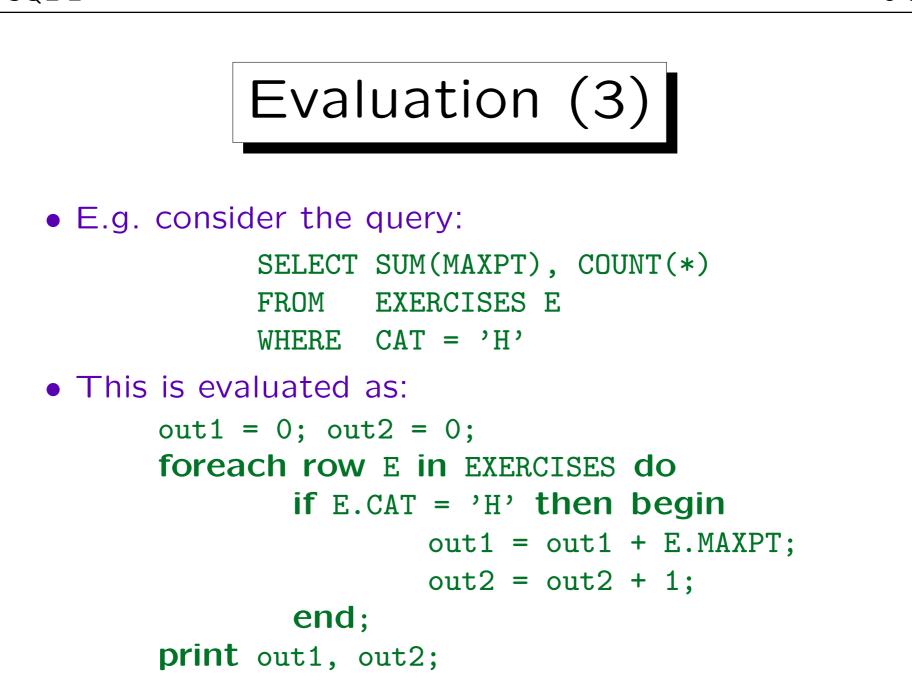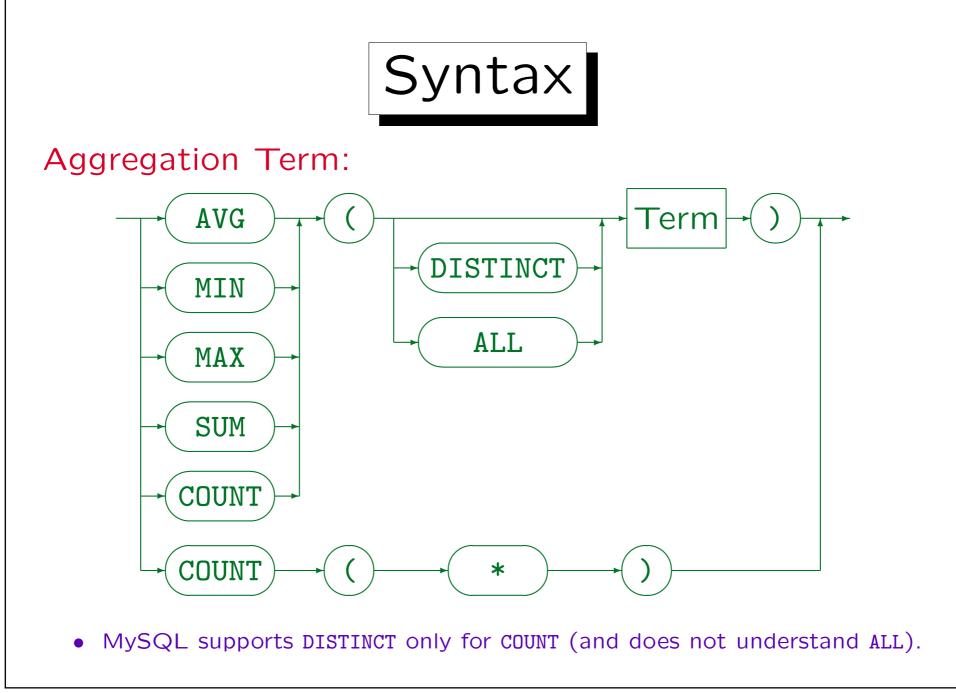- There are three different types of queries in SQL:

  ◇ Queries without aggregation functions and without `GROUP BY` and `HAVING`: See above.

  ◇ Queries with aggregation functions under `SELECT`, but no `GROUP BY` (called "simple aggregations" above): Result is always exactly one row.

  ◇ Queries with `GROUP BY`.

- Each type has different syntax restrictions and is evaluated in a different way.

# Evaluation (1)

- First, the `FROM`-clause is evaluated.

    Theoretically, all possible tuple combinations of the source tables are constructed (cartesian product, nested loops).

- Second, the `WHERE`-clause is evaluated.

    Only those tuple combinations that satisfy the condition are further considered (selection, filter, if). Of course, in real systems the first and second step may be combined to allow a more efficient evaluation.

- Third, if there is no aggregation, `GROUP BY`, and `HAVING`, the `SELECT`-clause is evaluated by printing the values of the terms/scalar expressions in the `SELECT`-list for every remaining tuple combination.
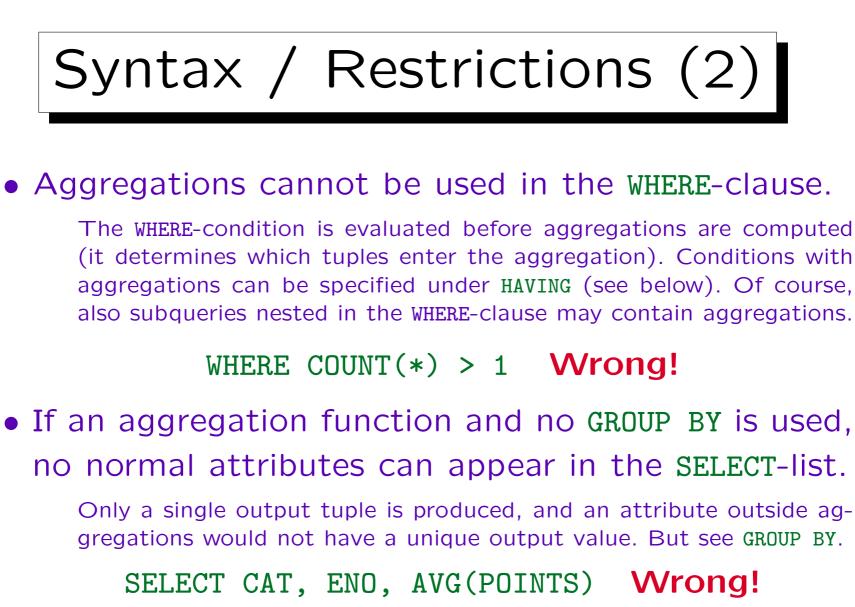
# Evaluation (2)

- When the `SELECT`-list contains an aggregation term, and there is no `GROUP BY`, only a single output row is computed by applying the aggregation operators.

- Instead of printing the values of columns as usual, the values are added to a set/multiset that is the input to the aggregation function.

  If the `SELECT`-list contains multiple aggregations, multiple such sets must be managed.

- If no `DISTINCT` is used, the aggregated values can be incrementally computed without explicitly storing a temporary set of values (see next slide).

# Evaluation (3)

- E.g. consider the query:

```
SELECT SUM(MAXPT), COUNT(*)
FROM    EXERCISES E
WHERE   CAT = 'H'
```

- This is evaluated as:

```
out1 = 0; out2 = 0;
foreach row E in EXERCISES do
        if E.CAT = 'H' then begin
                out1 = out1 + E.MAXPT;
                out2 = out2 + 1;
        end;
print out1, out2;
```

# Syntax

Aggregation Term:



- MySQL supports `DISTINCT` only for `COUNT` (and does not understand `ALL`).

# Syntax / Restrictions (1)

- The arguments of SUM and AVG must be numeric. COUNT, MIN, and MAX accept any datatype.

- Aggregations cannot be nested, e.g. the following is illegal:

  AVG(COUNT(*))     **Wrong!**

  After the COUNT only a single value remains. Thus, applying another aggregation makes no sense.
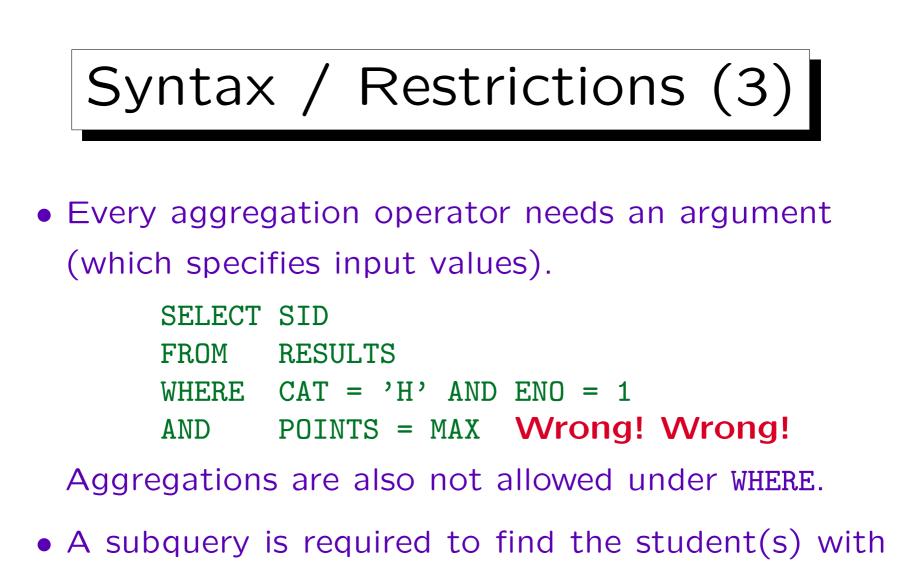
  > It is possible that aggregations are first applied to groups of rows, and then the result is input to another aggregation. E.g. what is the average over the total number of points students got for their homeworks? This is done with GROUP BY and subqueries (see below).

# Syntax / Restrictions (2)

- Aggregations cannot be used in the WHERE-clause.

    The WHERE-condition is evaluated before aggregations are computed (it determines which tuples enter the aggregation). Conditions with aggregations can be specified under HAVING (see below). Of course, also subqueries nested in the WHERE-clause may contain aggregations.

    WHERE COUNT(*) > 1    **Wrong!**

- If an aggregation function and no GROUP BY is used, no normal attributes can appear in the SELECT-list.

    Only a single output tuple is produced, and an attribute outside aggregations would not have a unique output value. But see GROUP BY.

    SELECT CAT, ENO, AVG(POINTS)    **Wrong!**
    FROM   RESULTS

# Syntax / Restrictions (3)

- Every aggregation operator needs an argument (which specifies input values).

```
SELECT  SID
FROM    RESULTS
WHERE   CAT = 'H' AND ENO = 1
AND     POINTS = MAX    Wrong! Wrong!
```

  Aggregations are also not allowed under WHERE.

- A subquery is required to find the student(s) with the best result for Homework 1 (see below).

# Null Values in Aggregations

- Usually, null values are ignored (filtered out) before the aggregation function is applied.

- Only COUNT(*) includes null values (it counts rows, not attribute values).

- The difference between COUNT(EMAIL) and COUNT(*) is that the first counts only those rows where EMAIL is not null, whereas the second counts all rows.

  Otherwise, the actual attribute value is not important for COUNT, and one probably should use COUNT(*). Of course, if duplicates are eliminated as in COUNT(DISTINCT CAT), the attribute is obviously important.

# Empty Aggregations

- If the input set is empty, most aggregations yield a null value, only `COUNT` returns 0.

  This is counter-intuitive at least for the `SUM`. One would expect that the `SUM` over the empty set is 0, but in SQL it returns `NULL`. (One reason for this behaviour might be that the `SUM` aggregation function cannot detect a difference between the empty input set because there was no qualifying tuple and the empty input set because all qualifying tuples had a null value in this argument.)

- Since it may happen that no row satisfies the `WHERE`-condition, programs must be prepared to process the resulting null value.

  Alternative: Use e.g. `NVL(SUM(POINTS),0)` in Oracle to replace the null.

# Overview

1. Subqueries, Nonmonotonic Constructs

2. Aggregations I: Aggregation Functions

3. Aggregations II: GROUP BY, HAVING

4. UNION, Conditional Expressions

5. Sorting Output: ORDER BY

6. SQL-92 Joins, Outer Join in Oracle

# GROUP BY (1)

- The above SQL constructs can produce a single aggregated output row only.

- The `GROUP BY` clause allows one to aggregate in groups rather than aggregate all tuples.

- Compute the average points for each homework:

```
SELECT    ENO, AVG(POINTS)
FROM      RESULTS
WHERE     CAT = 'H'
GROUP BY ENO
```

| ENO | AVG(POINTS) |
|-----|-------------|
| 1   | 8           |
| 2   | 8.5         |

# GROUP BY (2)

- The GROUP BY clause splits the resulting table after evaluation of FROM and WHERE into groups that have the same value in the GROUP BY columns.

| SID | CAT | ENO | POINTS |
|-----|-----|-----|--------|
| 101 | H   | 1   | 10     |
| 102 | H   | 1   | 9      |
| 103 | H   | 1   | 5      |
| 101 | H   | 2   | 8      |
| 102 | H   | 2   | 9      |

- The aggregation is then done over every group.
  So there will be one output row for every group.

# GROUP BY (3)

- This construction can never produce empty groups. So it is impossible that a COUNT(*) results in the value 0.

    The value 0 can be produced with COUNT(A) where the attribute A is null. If a query must produce groups with count 0, probably an outer join is needed (see below).

- On the other hand, simple aggregations (without GROUP BY) will always produce exactly one output row, and it is possible that their input set is empty (then COUNT(*) can be 0).

    A GROUP BY query can result in none, one, or many output rows.

# GROUP BY (4)

- Since the `GROUP BY` attributes have a unique value for every group, they can be used in the `SELECT`-list.

    Other attributes can be used under `SELECT` only inside aggregations.

- E.g. this is illegal:

```
SELECT    E.ENO, E.TOPIC, AVG(R.POINTS)   Wrong!
FROM      EXERCISES E, RESULTS R
WHERE     E.CAT='H' AND R.CAT='H' AND E.ENO=R.ENO
GROUP BY E.ENO
```

    E.TOPIC does not appear under `GROUP BY`, therefore it cannot be used
    in the `SELECT`-list outside an aggregation function. This is especially
    strange since `ENO` is a key of `EXERCISES`, so that `TOPIC` is actually unique
    in the groups. But the SQL rule is purely syntactic.

# GROUP BY (5)

- Thus, one must group by `E.ENO` and `E.TOPIC`:

```
SELECT    E.ENO, E.TOPIC, AVG(R.POINTS)
FROM      EXERCISES E, RESULTS R
WHERE     E.CAT='H' AND R.CAT='H' AND E.ENO=R.ENO
GROUP BY E.ENO, E.TOPIC
```

| E.ENO | E.TOPIC | AVG(POINTS) |
|-------|---------|-------------|
| 1 | Rel. Algeb. | 8 |
| 2 | SQL | 8.5 |

- Adding `E.TOPIC` to the `GROUP BY` attributes does not change the groups, but now one can print it.

# GROUP BY (6)

- Exercise: Is there any semantical difference between

```
SELECT    TOPIC, AVG(POINTS/MAXPT)
FROM      EXERCISES E, RESULTS R
WHERE     E.CAT='H' AND R.CAT='H' AND E.ENO=R.ENO
GROUP BY TOPIC
```

and the query which additionally groups by E.ENO, but does not print it?
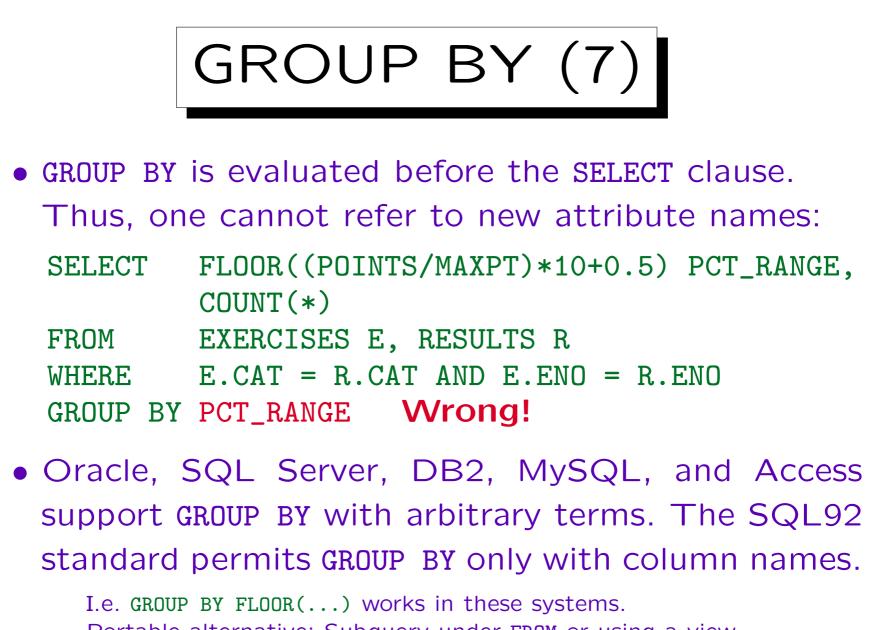
```
SELECT    TOPIC, AVG(POINTS/MAXPT)
FROM      EXERCISES E, RESULTS R
WHERE     E.CAT='H' AND R.CAT='H' AND E.ENO=R.ENO
GROUP BY TOPIC, E.ENO
```

# GROUP BY (7)

- GROUP BY is evaluated before the SELECT clause.

  Thus, one cannot refer to new attribute names:

  ```
  SELECT    FLOOR((POINTS/MAXPT)*10+0.5) PCT_RANGE,
            COUNT(*)
  FROM      EXERCISES E, RESULTS R
  WHERE     E.CAT = R.CAT AND E.ENO = R.ENO
  GROUP BY PCT_RANGE     Wrong!
  ```

- Oracle, SQL Server, DB2, MySQL, and Access support GROUP BY with arbitrary terms. The SQL92 standard permits GROUP BY only with column names.

  > I.e. GROUP BY FLOOR(...) works in these systems.
  > Portable alternative: Subquery under FROM or using a view.

# GROUP BY (8)

- The sequence of attributes in the `GROUP BY` clause is not important.

  > `GROUP BY A, B` means that two tuples $t$, $u$ belong into the same group if $t.A = u.A$ and $t.B = u.B$.
  >
  > `GROUP BY B, A` means that two tuples $t$, $u$ belong into the same group if $t.B = u.B$ and $t.A = u.A$.

- Note that it makes no sense to group by a key (if only one table is listed under `FROM`): Then every group will consist of only a single row.

- In the same way, `GROUP BY` is not useful if there can be only a single group.

# GROUP BY (9)

Warning:

- Many students mix up "GROUP BY" and "ORDER BY":
  - ◇ GROUP BY is important for the query result.
  - ◇ ORDER BY is only cosmetic (for a nice printout).

- GROUP BY usually internally sorts the tuples (so that tuples with the same values are adjacent).

- But then GROUP BY does the grouping, whereas the sort for the ORDER BY is done at the very end.

- Sometimes, the DBMS may evaluate the GROUP BY in more efficient ways without sorting.

# Syntax (1)

SELECT-Expression:

# Syntax (2)

Grouping:



- E.g. `GROUP BY TITLE, C.CRN`

- Oracle, SQL Server, DB2, Access, and MySQL support the more general "Term" instead of "Attribute-Reference". Of course, no aggregation functions are permitted under `GROUP BY`.

# HAVING (1)

- Aggregations cannot be used in the `WHERE`-clause.

- But sometimes aggregations are needed to filter output rows, not only for computing output values.

- For this reason, SQL has a second kind of condition, the `HAVING` clause. The purpose of the `HAVING` clause is to eliminate whole groups.

- Aggregation operators can be used in the `HAVING`-condition. But as under `SELECT`, outside aggregations, only `GROUP BY` attributes can be used.

# HAVING (2)

- Which students got at least 18 homework points?

```
SELECT    FIRST, LAST
FROM      STUDENTS S, RESULTS R
WHERE     S.SID=R.SID AND R.CAT='H'
GROUP BY  S.SID, FIRST, LAST
HAVING    SUM(POINTS) >= 18
```

| FIRST | LAST |
|---------|-------|
| Ann | Smith |
| Michael | Jones |

- The WHERE condition refers to single tuple combinations, the HAVING condition to entire groups.

# Evaluation

1. All combinations of rows from tables under `FROM` are considered.

2. The `WHERE`-condition selects a subset of these.

3. The remaining joined tuples are split into groups having equal values for the `GROUP BY`-attributes.

4. Groups of tuples which do not satisfy the condition in the `HAVING`-clause are eliminated.

5. One output tuple for every group is produced by evaluating the terms in the `SELECT`-clause.

# Syntax: Restrictions

- An aggregation is done if

    ◇ an aggregation function is used in the SELECT-list,

    ◇ or the GROUP BY or HAVING-clause is present.

- If an aggregation is done, then: Only GROUP BY attributes can be used under SELECT or HAVING outside aggregation functions.

    > Inside aggregation functions, i.e. as their arguments, all attributes can be used. E.g. AVG(A)/B: The attribute A appars inside an aggregation function, B outside.

- HAVING without GROUP BY is legal, but uncommon: The query could only return 0 or 1 output rows.

# WHERE vs. HAVING

- Normally, the restrictions uniquely define whether a condition must be put under `WHERE` or under `HAVING`.

  Only if a condition contains only `GROUP BY`-attributes, but no aggregations, it would be allowed in both clauses.

- If both is possible, it is much more efficient to put it under `WHERE`. E.g. this query is legal, but slow and needs lots of memory:

```
SELECT    FIRST, LAST
FROM      STUDENTS S, RESULTS R
GROUP BY S.SID, R.SID, FIRST, LAST
HAVING    S.SID = R.SID AND SUM(POINTS) >= 18
```

# Aggregation Subqueries (1)

- Who has the best result for Homework 1?

```
SELECT  S.FIRST, S.LAST, R.POINTS
FROM    STUDENTS S, RESULTS R
WHERE   S.SID=R.SID AND R.CAT='H' AND R.ENO=1
AND     R.POINTS = (SELECT MAX(POINTS)
                    FROM   RESULTS
                    WHERE  CAT='H' AND ENO=1)
```

- For an aggregation query without `GROUP BY`, it is guaranteed that it will return exactly one row. Thus `ANY/ALL` is not necessary here.

# Aggregation Subqueries (2)

- Since in SQL92, DB2, SQL Server, and Access a subquery returning a single data element can be used as a term, subqueries are also allowed in the `SELECT`-clause. Oracle 8.0 does not support this.

- This can replace `GROUP BY`. E.g. print for every student the sum of the homework points (null if none):

```
SELECT FIRST, LAST, (SELECT SUM(POINTS)
                     FROM   RESULTS R
                     WHERE  R.SID = S.SID
                     AND    R.CAT = 'H')
FROM    STUDENTS S
```

# Nested Aggregations (1)

- Nested aggregations require a subquery under `FROM`.

- What is the average number of homework points?
  (counting only students who submitted homeworks)

```
SELECT AVG(X.HW_PT)
FROM    (SELECT    SID, SUM(POINTS) AS HW_PT
         FROM      RESULTS
         WHERE     CAT = 'H'
         GROUP BY SID) X
```

| X | |
|-----|-------|
| SID | HW_PT |
| 101 | 18 |
| 102 | 18 |
| 103 | 5 |

| AVG(X.HW_PT) |
|--------------|
| 13.67 |

# Nested Aggregations (2)

- Oracle also supports nested aggregations written in this way:

```
SELECT    AVG(SUM(POINTS))    Only Oracle!
FROM      RESULTS
WHERE     CAT = 'H'
GROUP BY SID
```

This is completely non-standard (not supported in SQL92, DB2, SQL Server, Access).

Since it is much shorter than the equivalent standard query, it might be handy to use this when writing ad-hoc queries. However, in application programs, one should not create unnecessary portability problems.

# Aggregating Different Sets (1)

- Subqueries under `FROM` make it possible to aggregate over different sets:

```
SELECT FIRST, LAST, H.PT AS HOMEWORK, M.PT AS MID
FROM    STUDENTS S,
        (SELECT    SID, SUM(POINTS) AS PT
         FROM      RESULTS
         WHERE     CAT = 'H'
         GROUP BY SID) H,
        (SELECT    SID, SUM(POINTS) AS PT
         FROM      RESULTS
         WHERE     CAT = 'M'
         GROUP BY SID) M
   WHERE  S.SID = H.SID AND S.SID = M.SID
```

# Aggregating Different Sets (2)

- This is also possible with conditional expressions, e.g. in Oracle:

```
SELECT FIRST, LAST,
       SUM(DECODE(R.CAT, 'H', R.POINTS, 0)) HW
       SUM(DECODE(R.CAT, 'M', R.POINTS, 0)) MID
FROM   STUDENTS S, RESULTS R
WHERE  S.SID = R.SID
```

- E.g. the conditional expression

```
DECODE(R.CAT, 'H', R.POINTS, 0)
```

returns `R.POINTS` if `R.CAT = 'H'` and `0` otherwise.

# Maximizing Aggregations (1)

- Who has the best results in the homeworks (maximal sum of homework points)?

```
SELECT    FIRST, LAST, SUM(POINTS) AS TOTAL
FROM      STUDENTS S, RESULTS R
WHERE     S.SID = R.SID AND R.CAT = 'H'
GROUP BY S.SID, FIRST, LAST
HAVING    SUM(POINTS) >= ALL(SELECT    SUM(POINTS)
                             FROM      RESULTS
                             WHERE     CAT = 'H'
                             GROUP BY SID)
```

- Alternative solution with view: See next slide.

# Maximizing Aggregations (2)

- Total number of HW points for every student:

```
CREATE VIEW HW_TOTALS AS
    SELECT   SID, SUM(POINTS) AS TOTAL
    FROM     RESULTS
    WHERE    CAT = 'H'
    GROUP BY SID
```

- Then one can use this as follows:

```
SELECT S.FIRST, S.LAST, H.TOTAL
FROM   STUDENTS S, HW_TOTALS H
WHERE  S.SID = H.SID
AND    H.TOTAL = (SELECT MAX(TOTAL)
                  FROM   HW_TOTALS)
```

# Exercise: Possible Errors (1)

- What do you think about this query? Its task is to list all students who have solved at least two homeworks.

```
SELECT FIRST, LAST
FROM   STUDENTS S
WHERE  2 <= (SELECT COUNT(S.SID)
                FROM   RESULTS R
                WHERE  R.SID = S.SID
                AND    R.CAT = 'H')
```
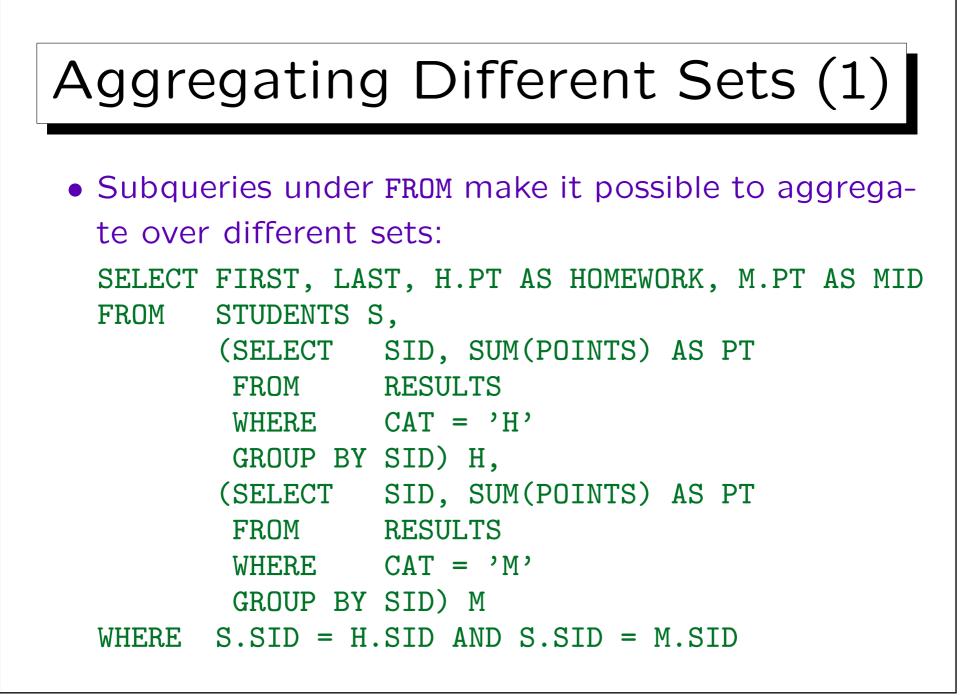
# Exercise: Possible Errors (2)

- And what about this query? Again, the task is to list students who have solved at least two homeworks.

```
SELECT FIRST, LAST
FROM   STUDENTS S, RESULTS R
WHERE  S.SID = R.SID
AND    R.CAT = 'H'
AND    COUNT(R.ENO) >= 2
```

# Exercise: Possible Errors (3)

- And what about this query? Here the task is to list the number of homeworks per student.

```
SELECT     S.SID, S.FIRST, S.LAST, SUM(R.ENO)
FROM       STUDENTS S, RESULTS R
WHERE      S.SID = R.SID
AND        R.CAT = 'H'
GROUP BY S.SID, S.FIRST, S.LAST, R.ENO
```

# Overview

# UNION (1)

- In SQL it is possible to combine the results of two queries by UNION.

  $R \cup S$ is the set of all tuples contained in $R$, in $S$, or in both.

- UNION is needed since otherwise there is no way to construct one result column that contains values drawn from different tables/columns.

  This is necessary e.g. when subclasses are represented by different tables. For instance, there may be one table GRADUATE_COURSES and another table UNDERGRADUATE_COURSES.

- UNION is also very useful for case analysis (to code an **if ... then ... else ...**).

# UNION (2)

- The subqueries which are operands to `UNION` must return tables with the same number of columns. The data types of corresponding columns must be compatible.

    The attribute names do not have to be equal. Oracle and SQL Server use the attribute names from the first operand in the result. DB2 uses artificial column names (1, 2, . . . ) if the input column names differ.

- SQL distinguishes between

    ◇ `UNION`: ∪ with duplicate elimination, and

    ◇ `UNION ALL`: concatenation (retains duplicates).

    Duplicate elimination is quite expensive.

# UNION (3)

- Print for every student his/her total number of homework points (0 if no homework submitted).

```
SELECT     S.FIRST, S.LAST, SUM(R.POINTS) AS TOTAL
FROM       STUDENTS S, RESULTS R
WHERE      S.SID = R.SID AND R.CAT = 'H'
GROUP BY S.SID, S.FIRST, S.LAST
UNION ALL
SELECT     S.FIRST, S.LAST, 0 AS TOTAL
FROM       STUDENTS S
WHERE      S.SID NOT IN (SELECT SID
                         FROM   RESULTS
                         WHERE  CAT = 'H')
```

# UNION (4)

- Assign student grades based on Homework 1:

```
SELECT S.SID, S.FIRST, S.LAST, 'A' GRADE
FROM    STUDENTS S, RESULTS R
WHERE   S.SID=R.SID AND R.CAT='H' AND R.ENO=1
AND     R.POINTS >= 9
    UNION ALL
SELECT S.SID, S.FIRST, S.LAST, 'B' GRADE
FROM    STUDENTS S, RESULTS R
WHERE   S.SID=R.SID AND R.CAT='H' AND R.ENO=1
AND     R.POINTS >= 7 AND R.POINTS < 9
    UNION ALL
...
```

# Other Set Operations in SQL

- SQL-86 contained only UNION [ALL].

- The SQL-92 standard also contains EXCEPT (set difference, −) and INTERSECT (∩).

  SQL-86, SQL Server and Access support only UNION [ALL]. MySQL does not support any of these operations. DB2 supports all SQL-92 set operators. In Oracle 8.0, the − operator is called MINUS instead of EXCEPT. ALL for MINUS and INTERSECT is not supported in Oracle.

- These operations add nothing to the expressivity to the language.

  Queries containing EXCEPT/MINUS and INTERSECT can be transformed into equivalent SQL-queries without these constructs, but queries containing UNION in general cannot. So only UNION is really important.

# UNION: Syntax

**Table Expression:**

```
──────┬──────────────┬──→┌─────────────────┐───────────┬──────┬──→
      │              │   │  SELECT Query   │           │      │
      │              │   └─────────────────┘           │      │
      │              │                                 │      │
      │              │   ┌─────────────────┐           │      │
      │              └→(─→│ Table Expression │→─)───────┘      │
      │                   └─────────────────┘                 │
      │                                                       │
      │              ┌──────────┐                             │
      └──────←───────┤  UNION   ├──────←──────────────────────┘
                     └──────────┘
                     ┌──────────┐
              ←───────┤ UNION ALL ├──────←
                     └──────────┘
```

- MySQL does not support union. SQL-86 contains `UNION` and `UNION ALL`.

- SQL-92 and DB2 support in addition `INTERSECT`, `INTERSECT ALL`, `EXCEPT`, and `EXCEPT ALL`. Oracle 8 supports `UNION`, `UNION ALL`, `INTERSECT` and `MINUS`.

- In Access, it is not possible to put parentheses around the entire query.

# Union vs. Join

Exercise:

- Two alternatives for respresenting the homework, midterm, and final results of the students are:

| Results_1 | | | |
|-----------|-----|-----|-----|
| STUDENT   | H   | M   | F   |
| Jim Ford  | 95  | 60  | 75  |
| Ann Lloyd | 80  | 90  | 95  |

| Results_2 | | |
|-----------|-----|-----|
| STUDENT   | CAT | PCT |
| Jim Ford  | H   | 95  |
| Jim Ford  | M   | 60  |
| Jim Ford  | F   | 75  |
| Ann Lloyd | H   | 80  |
| Ann Lloyd | M   | 90  |
| Ann Lloyd | F   | 95  |

- Write SQL queries to translate between the two.

# Conditional Expressions (1)

- Whereas using `UNION` is the portable way to make a case analysis, sometimes a conditional expression suffices, and is more efficient.

  Conditional expressions look differently in each DBMS.

- E.g. Oracle has expressions of the form:

$$\texttt{DECODE}(X,\ X_1,\ Y_1,\ X_2,\ Y_2,\ \ldots,\ Z)$$

- This is evaluated by comparing $X$ first to $X_1$, then to $X_2$, and so on. If $X_i$ is the first value with $X = X_i$, then $Y_i$ is returned. If no $X_i$ matches, $Z$ is returned.

# Conditional Expressions (2)

- E.g. print the exercise category in full for the results of Ann Smith (Oracle Version):

```
SELECT    DECODE(CAT, 'H', 'Homework',
                      'M', 'Midterm Exam',
                      'F', 'Final Exam',
                      'Unknown Category'),
          ENO, POINTS
FROM      STUDENTS S, RESULTS R
WHERE     S.SID = R.SID
AND       S.FIRST = 'Ann' AND S.LAST = 'Smith'
ORDER BY DECODE(CAT, 'H', 1, 'M', 2, 'F', 3, 4)
```

# Conditional Expressions (3)

- In the SQL-92 standard (and e.g. DB2), this is written as follows:

```
SELECT CASE WHEN CAT='H' THEN 'Homework'
            WHEN CAT='M' THEN 'Midterm Exam'
            WHEN CAT='F' THEN Final Exam'
            ELSE 'Unknown Category' END
       ENO, POINTS
FROM   STUDENTS S, RESULTS R
WHERE  S.SID = R.SID
AND    S.FIRST = 'Ann' AND S.LAST = 'Smith'
```

- Oracle 8i (not 8.0) supports a similar syntax, but requires a comma between the WHEN clauses.

# Conditional Expressions (4)

- The SQL-92 standard (and DB2, but not Oracle 8i) supports also the following abbreviation which is very similar to Oracle's `DECODE`:

```
SELECT CASE CAT WHEN 'H' THEN 'Homework',
                WHEN 'M' THEN 'Midterm Exam',
                WHEN 'F' THEN Final Exam',
                ELSE 'Unknown Category' END,
        ENO, POINTS
FROM    STUDENTS S, RESULTS R
WHERE   S.SID = R.SID
AND     S.FIRST = 'Ann' AND S.LAST = 'Smith'
```

# Conditional Expressions (5)

- A typical application of condtional expressions is to replace a null value by something else.

- In Oracle $\text{NVL}(X, Y)$ is equivalent to

$$\text{DECODE}(X, \text{NULL}, Y, X)$$

  I.e. if $X$ is not null, then $X$ is the result.
  If $X$ is null, then $Y$ is the result.

- $\text{COALESCE}(X, Y)$ is the same in standard SQL-92. There it abbreviates

$$\text{CASE WHEN } X \text{ IS NOT NULL THEN } X \text{ ELSE } Y \text{ END}$$

# Conditional Expressions (6)

- E.g. list the email address of all students, and write "(none)" if the column is null:

```
SELECT FIRST, LAST, NVL(EMAIL, '(none)')
FROM   STUDENTS
```

- Finally note that conditional expressions are normal terms, so they can be input for other datatype functions or e.g. aggregation functions.

# Overview

1. Subqueries, Nonmonotonic Constructs

2. Aggregations I: Aggregation Functions

3. Aggregations II: GROUP BY, HAVING

4. UNION, Conditional Expressions

5. Sorting Output: ORDER BY

6. SQL-92 Joins, Outer Join in Oracle

# Sorting Output (1)

- Output that is longer than a few lines should be sorted in some understandable way.

  It is much easier to search a specific value in a sorted table. Without "ORDER BY" the sequence of output rows means nothing (it depends on the algorithms used in the DBMS and may change between versions).

- However, it is important to understand that developing the logic of the query and nicely formatting the output are two separate things.

  Whereas sorting is the only formatting command that found its way into the SQL standard, DBMS tools usually offer more options. E.g. to have a pagebreak when the value in a specific column changes, to show negative values in red ink, etc. However, sorting may also be important when an application program retrieves the data.

# Sorting Output (2)

- E.g. print the students who solved homework 1.

  Order the list alphabetically by last name:

```
SELECT     S.FIRST, S.LAST
FROM       STUDENTS S, RESULTS R
WHERE      S.SID = R.SID
AND        R.CAT = 'H' AND R.ENO = 1
ORDER BY S.LAST
```

| FIRST   | LAST   |
|---------|--------|
| Michael | Jones  |
| Ann     | Smith  |
| Richard | Turner |

# Sorting Output (3)

- One can specify a prioritized list of sorting criteria.

    The "ORDER BY" list can contain multiple columns. The second column is only used for ordering two tuples which have the same value in the first column, and so on. Additional sorting criteria are only useful if there can still be duplicates in the previous columns.

- E.g.: Print the homework results sorted by exercise, and for each exercise by points (best result first), and if there is still a tie, alphabetically by name:

```
SELECT    R.ENO, R.POINTS, S.FIRST, S.LAST
FROM      STUDENTS S, RESULTS R
WHERE     S.SID = R.SID AND R.CAT = 'H'
ORDER BY R.ENO, R.POINTS DESC, S.LAST, S.FIRST
```

# Sorting Output (4)

- Result of the example query on the previous page:

| ENO | POINTS | FIRST | LAST |
|-----|--------|---------|--------|
| 1 | 10 | Ann | Smith |
| 1 | 9 | Michael | Jones |
| 1 | 5 | Richard | Turner |
| 2 | 9 | Michael | Jones |
| 2 | 8 | Ann | Smith |

- E.g. the first two tuples have the same value in the highest priority sort criterion (ENO), and the second criterion (POINTS DESC) determines their sequence.

  It does not matter that according to the criterion of third priority (LAST) the sequence would be the other way round.

# Sorting Output (5)

- According to the SQL-92 standard, one can only sort by columns that appear in the output.

  > E.g. it is impossible to print a list of student names ordered by total points without printing these points. But tools like SQL*Plus can suppress output columns from the query result.

- However, in all five systems (Oracle 8, DB2, SQL Server, Access, MySQL) one can sort by any term that would be allowed in the SELECT-clause.

  > In these systems, it is not necessary that the term really appears in the SELECT-clause. E.g. one can sort by UPPER(LAST), but print LAST. With DISTINCT, one can only sort by result columns (in Oracle one can still use them in terms and MySQL has no restriction).

# Sorting Output (6)

- Sometimes it is necessary to add columns to database tables to get a sort value, e.g.

  ◇ The results should be printed in the sequence: Homeworks, Midterm, Final (not alphabetically).

  ◇ The "University of Pittsburgh" should appear in a list of universities under "P", not under "U".

- If the student names were stored as a single string in the form "FIRST LAST", it would be (more or less) impossible to sort by last name.

  Important DB design question: What do I want to do with the data?

# Sorting Output (7)

- "DESC" means descending (inverse order from high to low values), the default is "ASC" (ascending).

- It is also possible to refer to columns by number, e.g.: ORDER BY 2, 4 DESC, 1

  > Column numbers refer to the sequence in the SELECT-list. They were important in earlier SQL versions, where one could not explicitly name the result columns. Today, one probably should use column names.

- Null values are all listed first or all listed last in the sort sequence (depending on the DBMS).

  > In Oracle, one can specify NULLS FIRST or NULLS LAST.

# Sorting Output (8)

- The effect of "`ORDER BY`" is purely cosmetic. It does not change the set of output tuples in any way.

- Thus, "`ORDER BY`" can only be applied at the very end of the query. It cannot be used in subqueries.

- Even when multiple `SELECT`-expressions are combined with `UNION`, the `ORDER BY` can only be placed at the very end (it refers to all result tuples).

# Sorting Output (9)

SQL Query:

Table Expression

ORDER BY → Order Specification

# Sorting Output (10)

Order Specification:



- Most DBMS permit "Term" instead of "Attribute Reference" (except if `DISTINCT` or `UNION` etc. are specified). Then basically the same restrictions apply as for terms in the `SELECT`-list (there might be additional restrictions for the use of aggregation functions).

# Overview

1. Subqueries, Nonmonotonic Constructs

2. Aggregations I: Aggregation Functions

3. Aggregations II: GROUP BY, HAVING

4. UNION, Conditional Expressions

5. Sorting Output: ORDER BY

6. SQL-92 Joins, Outer Join in Oracle

# Example Database (again)

**STUDENTS**

| SID | FIRST | LAST | EMAIL |
|-----|---------|--------|--------|
| 101 | Ann | Smith | $\cdots$ |
| 102 | Michael | Jones | (null) |
| 103 | Richard | Turner | $\cdots$ |
| 104 | Maria | Brown | $\cdots$ |

**EXERCISES**

| CAT | ENO | TOPIC | MAXPT |
|-----|-----|-------------|-------|
| H | 1 | Rel. Algeb. | 10 |
| H | 2 | SQL | 10 |
| M | 1 | SQL | 14 |

**RESULTS**

| SID | CAT | ENO | POINTS |
|-----|-----|-----|--------|
| 101 | H | 1 | 10 |
| 101 | H | 2 | 8 |
| 101 | M | 1 | 12 |
| 102 | H | 1 | 9 |
| 102 | H | 2 | 9 |
| 102 | M | 1 | 10 |
| 103 | H | 1 | 5 |
| 103 | M | 1 | 7 |

# Joins in SQL-92 (1)

- An important and useful operation of relational algebra is the join (in several variants).

- In SQL-86, one cannot directly specify a join. One writes a cartesian product (`FROM`) and then does a selection (`WHERE`). This is still the usual case.

- E.g. the natural join of `RESULTS` and `EXERCISES` is:

```
SELECT R.CAT AS CAT, R.ENO AS ENO, SID,
       POINTS, TOPIC, MAXPT
FROM   RESULTS R, EXERCISES E
WHERE  R.CAT = E.CAT AND R.ENO = E.ENO
```

# Joins in SQL-92 (2)

- In SQL-92 one can write e.g.

```
SELECT SID, ENO, (POINTS/MAXPT)*100
FROM    RESULTS R NATURAL JOIN EXERCISES E
WHERE   CAT = 'H'
```

- Because of the keywords "NATURAL JOIN" the system automatically adds the join condition

```
R.CAT = E.CAT AND R.ENO = E.ENO
```

- SQL-92 permits to use joins in the FROM-clause and even on the outer query level (like UNION).

  So one can write quite a lot in "relational algebra style".

# Joins in SQL-92 (3)

- Current systems support the standard only partially:

  ◇ SQL-92 joins are not supported in Oracle 8i.

    But Oracle 9i supports nearly the complete set.

  ◇ Some types of joins are supported in DB2, SQL

    Server, and Access, but the above "natural join"

    is not. A join with explicit condition is possible:

```
SELECT SID, R.ENO, (POINTS/MAXPT)*100
FROM   RESULTS R INNER JOIN EXERCISES E
       ON R.CAT = E.CAT AND R.ENO = E.ENO
WHERE  R.CAT = 'H'
```

# Joins in SQL-92 (4)

- With the explicit join condition, the query is not shorter than the equivalent one with the standard `WHERE` condition.

- The power of SQL is not increased by adding the new join constructs.

    Every query with the new join constructs can be translated in an equivalent one that does not use these constructs.

- The reason why joins where added to SQL is probably the "outer join": For the outer join, the equivalent formulation in SQL-86 is significantly longer.

# Outer Join: Repetition

- The usual join eliminates tuples without partner:

$$
\begin{array}{|c|c|}
\hline
A & B \\
\hline
a_1 & b_1 \\
a_2 & b_2 \\
\hline
\end{array}
\bowtie
\begin{array}{|c|c|}
\hline
B & C \\
\hline
b_2 & c_2 \\
b_3 & c_3 \\
\hline
\end{array}
=
\begin{array}{|c|c|c|}
\hline
A & B & C \\
\hline
a_2 & b_2 & c_2 \\
\hline
\end{array}
$$

- The left outer join guarantees that tuples from the left table will appear in the result:

$$
\begin{array}{|c|c|}
\hline
A & B \\
\hline
a_1 & b_1 \\
a_2 & b_2 \\
\hline
\end{array}
\bowtie
\begin{array}{|c|c|}
\hline
B & C \\
\hline
b_2 & c_2 \\
b_3 & c_3 \\
\hline
\end{array}
=
\begin{array}{|c|c|c|}
\hline
A & B & C \\
\hline
a_1 & b_1 & \\
a_2 & b_2 & c_2 \\
\hline
\end{array}
$$

Rows from the left table are filled with "null" if necessary.
There are also a right outer join and a full outer join.

# Outer Join in SQL-92 (1)

- E.g. number of submissions per homework. If there is no submission, the number 0 should be printed:

```
SELECT    E.ENO, COUNT(SID)
FROM      EXERCISES E LEFT OUTER JOIN RESULTS R
          ON E.CAT = R.CAT AND E.ENO = R.ENO
WHERE     E.CAT = 'H'
GROUP BY E.ENO
```

- All exercises are present in the result of the left outer join. In exercises without solutions, the attributes of SID and POINTS are filled with null values.

- COUNT(SID) does not count rows where SID is null.

# Outer Join in SQL-92 (2)

- Equivalent query without outer join (12 vs. 5 lines):

```
SELECT    E.ENO, COUNT(*)
FROM      EXERCISES E, RESULTS R
WHERE     E.CAT = 'H' AND R.CAT = 'H'
AND       E.ENO = R.ENO
GROUP BY E.ENO
UNION ALL
SELECT    E.ENO, 0
FROM      EXERCISES E
WHERE     E.CAT = 'H'
AND       E.ENO NOT IN (SELECT R.ENO
                        FROM   RESULTS R
                        WHERE  R.CAT = 'H')
```

# Outer Join in SQL-92 (3)

- E.g. print for every student the number of home-works he/she has solved (including 0).

- The following query does not work:
  Students without homework are not listed.

```
SELECT    FIRST, LAST, COUNT(ENO)     Wrong!
FROM      STUDENTS S LEFT OUTER JOIN RESULTS R
          ON S.SID = R.SID
WHERE     R.CAT = 'H'
GROUP BY S.SID, FIRST, LAST
```

- The outer join is constructed before the WHERE-condition is evaluated.

# Outer Join in SQL-92 (4)

- In general, one must be careful not to eliminate possible join partners after the outer join is done.

- One has to select the homework results before the outer join is done:

```
SELECT    FIRST, LAST, COUNT(R.ENO)
FROM      STUDENTS S LEFT OUTER JOIN
                (SELECT SID, ENO
                  FROM   RESULTS
                  WHERE  CAT = 'H') R
          ON S.SID = R.SID
GROUP BY S.SID, FIRST, LAST
```

# Outer Join in SQL-92 (5)

- One can also put the condition on the right table into the join condition:

```
SELECT    FIRST, LAST, COUNT(R.ENO)
FROM      STUDENTS S LEFT OUTER JOIN RESULTS R
          ON S.SID = R.SID AND R.CAT = 'H'
GROUP BY S.SID, FIRST, LAST
```

- The SQL-92 permits any WHERE-condition that refers only to the tuple variables on the left and right side of the join. (But don't abuse this.)

  It seems that DB2 and Access permit no subqueries in the ON-clause.
  More complex conditions must be enclosed in parentheses in Access.

# Outer Join in SQL-92 (6)

- Conditions on the left table make little sense in the condition of the left outer join.

- E.g. consider this query:

```
SELECT E.CAT, E.ENO, R.SID, R.POINTS
FROM   EXERCISES E LEFT OUTER JOIN RESULTS R
       ON E.CAT = 'H' AND R.CAT = 'H'
          AND E.ENO = R.ENO
```

- Exercise: Will E.CAT = 'M' appear in the output?

  ☐ yes      ☐ no

# Outer Join in SQL-92 (7)

- MySQL has no subqueries, but sometimes one can use the outer join instead.

- E.g. students who did not submit any homework:

```
SELECT S.SID. S.FIRST, S.LAST
FROM   STUDENTS S LEFT OUTER JOIN RESULTS R
       ON S.SID = R.SID AND R.CAT = 'H'
WHERE  R.CAT IS NULL
```

- Of course, instead of `R.CAT` one can test any attribute of `RESULTS` for the null value.

  The test for the null value checks whether the current `STUDENTS` tuple did not find a join partner.

# Join Syntax in SQL-92 (1)

- SQL-92 has the following join types:

  ◇ [INNER] JOIN: Usual Join.

  ◇ LEFT [OUTER] JOIN: Preserves rows from left table.

  ◇ RIGHT [OUTER] JOIN: Preserves right table tuples.

  ◇ FULL [OUTER] JOIN: All input tuples are preserved.

  ◇ CROSS JOIN: Cartesian product ×.

  ◇ UNION JOIN: This is a union that fills the columns of the other table with null values.

- The brackets mean that INNER/OUTER are optional.

# Join Syntax in SQL-92 (2)

- The join condition can be specified as follows:

  ◇ The keyword `NATURAL` in front of the join name.

  ◇ "`ON` ⟨`Condition`⟩" follows the join.

  ◇ "`USING` $(A_1, \ldots, A_n)$" follows the join.

    `USING` lists join attributes (e.g. for specifying the natural join). Attributes with the names $A_1, \ldots, A_n$ must appear in both tables and the join condition is $R.A_1 = S.A_1 \wedge \cdots \wedge R.A_n = S.A_n$. `NATURAL` is equivalent to specifying `USING` with all common attribute names.

- Only one of these constructs can be used.

- `CROSS JOIN` and `UNION JOIN` have no join condition.

# Join Syntax in SQL-92 (3)

- According to the standard, the `NATURAL` join and the join with `USING` produce a table with only one copy of the common attributes.

- Furthermore, the common attributes are listed first and cannot be referenced with a tuple variable.

```
SELECT *
FROM    RESULTS R NATURAL JOIN EXERCISES E
```

- The result columns are `CAT`, `ENO`, `R.SID`, `R.POINTS`, `E.TOPIC`, `E.MAXPT` (in this sequence).

    It is illegal to refer to `R.CAT` or `E.CAT`, only `CAT` can be used (and the same for `ENO`).

# Join Syntax in SQL-92 (4)

- Oracle 9i supports the SQL-92 joins.

    Including the natural join, but except UNION JOIN (which was removed in SQL:1999). Oracle 8i did not support any SQL-92 joins.

- Inner and outer join with ON work also in DB2, SQL Server, Access, and MySQL.

    In Access and MySQL, the keyword INNER is not optional.

- USING and NATURAL work only in Oracle 9i.

    NATURAL exists also in MySQL, but MySQL does not merge the common columns. This violates the SQL-92 standard.

# Join Syntax in SQL-92 (5)

- CROSS JOIN is supported only in Oracle 9i, SQL Server and MySQL, not Access and DB2.

    One can write a comma for CROSS JOIN, so it is not very useful.

- UNION JOIN is supported in none of the five systems.

    However, in SQL-92 (and e.g. Oracle, DB2, SQL Server, not Access), one can write a subquery containing UNION or UNION ALL also in the FROM-clause. So with a bit more keystrokes, one can simulate the union join. By the way, it is a bit strange that e.g. "FROM A NATURAL JOIN B" is legal in SQL-92, but "FROM A UNION B" is not. Also, SQL-92 permits to write "FROM (SELECT * FROM A UNION SELECT * FROM B) X", but the same with "NATURAL JOIN" instead of "UNION" is a syntax error [Date/Darwen, 1997, p. 148].

# Join Syntax in SQL-92 (6)

- In the `FROM` clause, one can also combine joins and the declaration of further tuple variables (separated by "," as usual).

- One can also join the result of joining two tables with a third one (and so on). The syntax is:

```
SELECT ...
FROM    R LEFT JOIN S ON R.A=S.B
             LEFT JOIN T ON S.C=T.D
```

- It is also possible to use parentheses, but then one has to declare a new tuple variable after the (...).

# Outer Join in Oracle (1)

- In Oracle, the outer join is traditionally specified under `WHERE` (no longer necessary in 9i).

- Instead of the join condition $R.A = S.B$ one writes

  ◇ $R.A = S.B$ **(+)** for the left outer join of $R$ and $S$,

  ◇ $R.A$ **(+)** $= S.B$ for the right outer join of $R$ and $S$.

  I.e. the special marker "(+)" is appended to attributes of the table which can be replaced with nulls.

    I.e. this protects the tuples of the other table (not marked with "(+)").
    There are many syntactic restrictions which ensure that this is really
    an outer join. If the join is done on several attributes, all must be
    marked. It is possible to write also $S.B$ **(+)** $= c$ with a constant $c$ or
    e.g. $R.A = S.B$ **(+)** $+ 1$.

# Outer Join in Oracle (2)

- E.g. number of submissions per exercise (can be 0):

```
SELECT    E.CAT, E.ENO, COUNT(SID)
FROM      EXERCISES E, RESULTS R
WHERE     E.CAT = R.CAT(+) AND E.ENO = R.ENO(+)
GROUP BY E.CAT, E.ENO
```

- As in the SQL-92 outer join, the outer join is constructed before any other conditions in the WHERE-clause are applied.

  > No matter in what sequence the conditions are written. But as shown above, one can use a subquery under FROM to do a selection before the outer join.