course:

**Database Systems** (A7B36DBS) lecture 6:

# Query formalisms for relational model – relational algebra

Doc. RNDr. Irena Holubova, Ph.D.

Acknowledgement: The slides were kindly lent by Doc. RNDr. Tomas Skopal, Ph.D., Department of Software Engineering, Charles University in Prague

#### **Today's lecture outline**

- relational algebra
  - relational operations
  - equivalent expressions
  - relational completeness

#### Map of the Lecture



estimatedEndDate

offeredPrice

wardedSupplier

Lissuad Contr

streetNumber city country

+tenderAddres

0..1



ame

IOUSE

Vindows

Printer

Phone

aptop

Bing

key

Database design

Formulation of the task



startDate > ...;

#### Database query

- query = delimiting particular set of data instances
  - a single query may be expressed by multiple expressions of the query language – equivalent expressions
- query extent (power of the query language)
  - in classical models, only subset of the database is expected as a query result
    - i.e., values actually present in the database tables
  - in **extended models**, also derived data can be returned
    - i.e., computations, statistics, aggregations derived from the data

### **Query language formalisms**

- as the "table data" model is based on the relational model, there can be used well-known formalisms
  - relational algebra (this lecture)
    - operations on relations used as query constructs
  - relational calculus (next lecture)
    - database extension of the first-order logic used as a query language

### Relational algebra (RA)

- RA is a set of operations (unary or binary) on relations (having schemas); their results are also relations (having schemas)
  - for completeness, with a relation (table content) R\* we always consider also a schema R(A) consisting of name and (typed) attributes, i.e., a tuple <R\*, R(A)>
- a schema will be named by any unique user-defined identifier
  - for the relation resulting from an operation we mostly do not need to define a name for the relation and the schema
    - it either enters another operation or it is the final result
  - if we need to "store" (or label) the result, e.g., for decomposition of complex query, we use

#### ResultName := <*expression consisting of relational operations*>

### Relational algebra (RA)

- if it is clear from the context, we use just R<sub>1</sub> operation R<sub>2</sub> instead of <R<sub>1</sub>, R<sub>1</sub>(A<sub>1</sub>)> operation <R<sub>2</sub>, R<sub>2</sub>(A<sub>2</sub>)>
- for binary operations we use infix notation
- for unary operations we use postfix notation
- the operation result can be used recursively as an operand of another operation, i.e., a tree of operations can be defined for a more complex query



#### RA – attribute renaming

attribute renaming (unary operation)

$$\begin{array}{l} \mathbf{R}^{*} < \mathbf{a}_{i} \rightarrow \mathbf{b}_{i}, \ \mathbf{a}_{j} \rightarrow \mathbf{b}_{j}, \ \dots > = \\ < \mathbf{R}^{*}, \ \mathbf{R}_{x}((\mathbf{A} - \{\mathbf{a}_{i}, \ \mathbf{a}_{j}, \ \dots\}) \cup \{\mathbf{b}_{i}, \ \mathbf{b}_{j}, \ \dots\}) > \end{array}$$

 only attributes in the schema are renamed, no data manipulation (i.e., the result is the same relation and the same schema, just of different attribute names)

#### RA – set operations

- set operations (binary, infix notation)
  - union:  $(R_1, R_1(A)) \cup (R_2, R_2(A)) = (R_1 \cup R_2, R_2(A))$
  - intersection:  $\langle R_1, R_1(A) \rangle \cap \langle R_2, R_2(A) \rangle = \langle R_1 \cap R_2, R_2(A) \rangle$
  - subtraction:  $\langle R_1, R_1(A) \rangle \langle R_2, R_2(A) \rangle = \langle R_1 R_2, R_2(A) \rangle$
  - Cartesian product:  $\langle R_1, R_1(A) \rangle \times \langle R_2, R_2(B) \rangle$

 $= \langle \mathsf{R}_1 \times \mathsf{R}_2, \mathsf{R}_x(\{\mathsf{R}_1\} \times \mathsf{A} \cup \{\mathsf{R}_2\} \times \mathsf{B}) \rangle$ 

- union, intersection and subtraction require compatible schemas of the operands!!!
  - it is also the schema of the result
  - i.e., we cannot, e.g., unify two different schemas

#### **RA – Cartesian product**

- a Cartesian product produces a new schema consisting of attributes from both source schemas
  - if the attribute names are ambiguous, we use a prefix notation, e.g., R<sub>1</sub>.a, R<sub>2</sub>.a
- if both the operands are the same, we first need to rename the attributes of one operand, i.e.,

 $\langle R_1, R_1(\{a,b,c\}) \rangle \times R_1 \langle a \rightarrow d, b \rightarrow e, c \rightarrow f \rangle$ 

#### Example – set operations

- FILM(FILM\_NAME, ACTOR\_NAME)
- AMERICAN\_FILM = {('Titanic', 'DiCaprio'), ('Titanic', 'Winslet'), ('Top Gun', 'Cruise')}
- NEW\_FILM = {('Titanic', 'DiCaprio'), ('Titanic', 'Winslet'), ('Samotáři', 'Macháček')}
- CZECH\_FILM = {('Pelíšky', 'Donutil'), ('Samotáři', 'Macháček')}

```
ALL_FILMS := AMERICAN_FILM \cup CZECH_FILM =
{(`Titanic', `DiCaprio'), (`Titanic', `Winslet'), (`Top Gun', `Cruise'),
(`Pelíšky', `Donutil'), (`Samotáři', `Macháček')}
```

OLD\_AMERICAN\_AND\_CZECH\_FILM := (AMERICAN\_FILM \colored CZECH\_FILM) - NEW\_FILM = {(`Top Gun', `Cruise'), (`Pelíšky', `Donutil')}

NEW\_CZECH\_FILM := **NEW\_FILM**  $\cap$  **CZECH\_FILM** = {(`Samotáři', `Macháček')}

### **RA** – projection

projection (unary operation)

**<R\*[C], R(A)>** = <{ $U[C] \in R*$ }, R(C)>, where C  $\subseteq$  A

- u[C] = values only in attributes from C
- possible duplicities are removed

#### **RA** – selection

selection (unary)

<R\*( $\phi$ ), R(A)> = <{ $\upsilon | \upsilon \in R^* \text{ and } \phi(\upsilon)$ }, R(A)>

- selection of those elements from R\* that match a condition φ
- condition φ is a Boolean expression
  - using and, or, not on atomic formulas t<sub>1</sub> O t<sub>2</sub> or t<sub>1</sub> O a
    - $\Theta \in \{<_{I}>_{I}=_{I}\geq_{I}\leq_{I}\neq\}$
    - t<sub>i</sub> are names of attributes

#### RA – natural join

natural join (binary)

<R\*, R(A)> \* <S\*, S(B)> =
 <{u | u[A]  $\in$  R\* and u[B]  $\in$  S\*}, R<sub>x</sub>(A  $\cup$  B)>

- joining elements of relations A, B using identity on <u>all shared</u> attributes
- if  $A \cap B = \emptyset$ , natural join corresponds to Cartesian product
  - no shared attributes
  - everything in A is joined with everything in B
- could be expressed using Cartesian product, selection and projection

#### Example – selection, projection, natural join

FILM(FILM\_NAME, ACTOR\_NAME)
FILM = {('Titanic', 'DiCaprio'), ('Titanic', 'Winslet'), ('Top Gun', 'Cruise')}
ACTOR(ACTOR\_NAME, BIRTH\_YEAR)
ACTOR = {('DiCaprio',1974), ('Winslet',1975), ('Cruise', 1962), ('Jolie', 1975)}

ACTOR\_YEAR := ACTOR[BIRTH\_YEAR] =

{(1974), (1975), (1962)}

YOUNG\_ACTOR := ACTOR(BIRTH\_YEAR > 1970) [ACTOR\_NAME] = {(`DiCaprio'), (`Winslet'), ('Jolie')}

FILM\_ACTOR := FILM \* ACTOR =

{('Titanic', 'DiCaprio', 1974), ('Titanic', 'Winslet', 1975), ('Top Gun', 'Cruise', 1962)}

#### RA – inner Θ-join

#### theta

inner Θ-join (binary)

## <R\*, R(A)>[t<sub>1</sub> $\Theta$ t<sub>2</sub>]<S\*, S(B)> = <{u | u[A] $\in$ R\*, u[B] $\in$ S\*, u.t<sub>1</sub> $\Theta$ u.t<sub>2</sub>}, A $\cup$ B>

- generalization of natural join

### RA – left 🛛 -semi-join

Ieft inner Θ-semi-join (binary)

#### <R\*, R(A)><t<sub>1</sub> $\Theta$ t<sub>2</sub>]<S\*, S(B)> = (R[t<sub>1</sub> $\Theta$ t<sub>2</sub>]S)[A]

- join restricted to the "left side"
  - only attributes of A in the resulting schema
- right semi-join similar
  - only attributes of B in the resulting schema

#### Inner vs. outer join

- in practice, it is useful to introduce **null meta-values** (NULL) of attributes
- outer join appends series of NULL values to those elements, that were not joined (i.e., they do not appear in inner join)
  - left outer join

 $\mathbf{R} *_{\mathsf{L}} \mathbf{S} = (\mathbf{R} * \mathbf{S}) \cup (\mathbf{R} \times (\mathsf{NULL}, \mathsf{NULL}, \ldots))$ 

• right outer join

 $\mathbf{R} *_{\mathbf{R}} \mathbf{S} = (\mathbf{R} * \mathbf{S}) \cup ((\mathbf{NULL}, \mathbf{NULL}, ...) \times \underline{\mathbf{S}})$ 

where  $\mathbb{R}$ , resp.  $\underline{S}$  consist of *n*-tuples not joined with S, resp. R

full outer join

 $\mathbf{R} *_{\mathbf{F}} \mathbf{S} = (\mathbf{R} *_{\mathsf{L}} \mathbf{S}) \cup (\mathbf{R} *_{\mathsf{R}} \mathbf{S})$ 

- the above joins are defined as natural joins, outer  $\Theta$ -joins are defined similarly
- the reason for outer join is a complete information on elements of a relation being joined
  - some are joined regularly, some only with NULLs

#### **RA** – relation division

relation division (binary)

<R\*, R(A)> ÷ <S\*, S(B  $\subset$  A)> =
<{t |  $\forall s \in S^*$  (t  $\oplus s$ )  $\in R^*$ }, A - B}

- used in situations where objects with all properties are needed
  - kind of universal quantifier in RA
- ⊕ is concatenation operation
  - relation elements <a<sub>1</sub>, a<sub>2</sub>, ...> and <b<sub>1</sub>, b<sub>2</sub>, ...> become <a<sub>1</sub>, a<sub>2</sub>, ..., b<sub>1</sub>, b<sub>2</sub>, ...>
- returns those elements from R\* that, when projected on A–B, are duplicates and, when projected on B, is equal to S\*
- alternative definition:  $\mathbf{R}^* \div \mathbf{S}^* = \mathbf{R}^*[A-B] ((\mathbf{R}^*[A-B] \times S^*) \mathbf{R}^*)[A-B]$

#### **Example – relation division**

FILM(FILM\_NAME, ACTOR\_NAME) ACTOR(ACTOR\_NAME, BIRTH\_YEAR)

What are the films where **all** the actors appeared? ACTOR\_ALL\_FILM := **FILM ÷ ACTOR[ACTOR\_NAME])** = {('Titanic')}

FILM_NAME	ACTOR_NAME		ACTOR_NAME	BIRTH_YEAR
Titanic	DiCaprio –		DiCaprio	1974
Titanic	Winslet	V,	Zane	1966
The Beach	DiCaprio	M	Winslet	1075
Enigma	Winslet	T	WINSIEt	1975
The Kiss	Zane			
Titanic	Zane			

#### **RA query evaluation**

- logical order of operation evaluation
  - depth-first traversal of a syntactic tree
  - e.g., (((S1 op1 S2) op2 (S3 op4)) op5 S4 op6 S5)
  - syntactic tree construction (query parsing) is driven by operation op6 priorities, parentheses, or associativity conventions operation precedence (priority) projection R[] (highest) 1. 2. selection R() op1 Cart. product X join, division \*,÷ subtraction 5.
    - 6. union, intersection  $\cup$ ,  $\cap$  (lowest)

#### Example – query evaluation

To which destination can fly Boeings? (such that all passengers in the flight fit the plane)

(Flight[Passengers, Destination] [Passengers <= Capacity] (Plane(Plane = 'Boeing\*')[Capacity]))[Destination]



#### **Equivalent expressions**

- a single query may be defined by multiple expressions
  - by replacing "redundant" operations with the basic ones (e.g., division, natural join)
  - by use of commutativity, distributivity and associativity of (some) operations
- selection
  - selection cascade  $(...((R(\phi_1))(\phi_2))...)(\phi_n) \equiv R(\phi_1 \land \phi_2 \land ... \land \phi_n)$
  - commutativity of selection  $(R(\phi_1))(\phi_2) \equiv (R(\phi_2))(\phi_1)$
- projection
  - projection cascade  $(...(R[A_1])[A_2])...)[A_n] \equiv R[A_n]$ , where  $A_n \subseteq A_{n-1} \subseteq ... \subseteq A_2 \subseteq A_1$
- join and Cartesian product
  - commutativity  $\mathbf{R} \times \mathbf{S} \equiv \mathbf{S} \times \mathbf{R}, \mathbf{R} [\Theta] \mathbf{S} \equiv \mathbf{S} [\Theta] \mathbf{R}, \text{ etc.}$
  - associativity  $\mathbf{R} \times (\mathbf{S} \times \mathbf{T}) \equiv (\mathbf{R} \times \mathbf{S}) \times \mathbf{T}, \mathbf{R} [\Theta] (\mathbf{S} [\Theta] \mathbf{T}) \equiv (\mathbf{R} [\Theta] \mathbf{S}) [\Theta] \mathbf{T}, \text{ etc.}$

### **Relational completeness**

- not all the mentioned operations are necessary for expression of every query
  - the minimal set consists of the following operations B = {union, Cartesian product, subtraction, selection, projection, attribute renaming}
- relational algebra query language is a set of expressions that result from composition of operations in B over a database schema
- if two expressions denote the same query they are equivalent
- query language that is able to express all queries of RA is relational complete
- Questions:

- How can we prove that a particular language is relational complete?
- Is SQL relational complete?

#### **RA – properties**

#### RA = declarative query language

- i.e., non-procedural, however, the structure of the expression suggests the sequence of operations
- the result is **always finite** relation
  - "safely" defined operations
- operation properties
  - associativity, commutativity
    - cart. product, join