Lecture 7

# **Relational Algebra**

**Martin Svoboda**

martin.svoboda@fel.cvut.cz

10. 4. 2018

**Czech Technical University in Prague**, Faculty of Electrical Engineering

# Lecture Outline

**Relational algebra**

- **Operations**: syntax, semantics and examples
  - Selection, projection, attribute renaming
  - Cartesian product, natural join, theta join, …
  - Division
  - Outer join
- **Relational completeness**

# Relational Model

**Relational model**

- Logical model where all data is represented in terms of **tuples** (rows) that are grouped into **relations** (tables)

**Schema of a relation**

- $S(a_1 : T_1, \ldots, a_n : T_n)$
  - $S$ is a relation name
  - $a_i$ are attribute names, $T_i$ are optional domains (data types)

**Relation** = data

- <u>Set</u> of **tuples**
- Unordered, no duplicities, without missing values (null), first normal form

# Relational Model

**Relation** structure revisited

- *I.e. formal definition for the purpose of this lecture*

$$\langle R, A_R \rangle$$

- $R$ = **set of tuples** = actual data
  - Tuple $t = \{(a_1, v_1), \dots, (a_n, v_n)\}$, where:
    - $a_i \in A_R$ is an **attribute** name
    - $t(a_i) = v_i \in T_i$ is a **value** this attribute is associated with
  - I.e. each tuple acts as a function
    - $t : A_R \to \bigcup_{i=1,\dots,n} T_i$
    - $t(a_i) = v_i \in T_i$
- $A_R$ = **set of attributes** = schema without relation name

# Relation Structure: Example

Sample relation of actors

- **Actor(name, surname, year)**

| name | surname | year |
|------|---------|------|
| Ivan | Trojan | 1964 |
| Jiří | Macháček | 1966 |
| Jitka | Schneiderová | 1973 |

$$
\Big\langle
$$
$$
\{ \ \{(name, Ivan), (surname, Trojan), (year, 1964)\},
$$
$$
\{(name, Jiří), (surname, Macháček), (year, 1966)\},
$$
$$
\{(name, Jitka), (surname, Schneiderová), (year, 1973)\} \ \},
$$
$$
\{name, surname, year\}
$$
$$
\Big\rangle
$$

# Query Languages

Formal **query languages** based on the **relational model**

- **Relational algebra**
    - Algebraic expressions with relations and operations on them
    - E.g. names and surnames of all actors born in 1970 or earlier
    
    $\pi_{\text{name,surname}}(\sigma_{\text{year} \leq 1970}(\text{Actor}))$
    
    $\text{Actor}(\text{year} \leq 1970)[\text{name, surname}]$

- **Relational calculi**
    - Expressions based on the **first-order predicate logic**
    - Domain relational calculus
        - E.g. $\{\, (n, s) \mid \exists\, y : \text{Actor}(n, s, y) \land y \leq 1970 \,\}$
    - Tuple relational calculus
        - E.g. $\{\, t[\text{name, surname}] \mid \text{Actor}(t) \land t.\text{year} \leq 1970 \,\}$

# Query Languages: Terminology

**Query language**

- **Set of all syntactically well-formed query expressions** with respect to a certain grammar
- E.g. relational algebra, SQL, ...

**Query**

- **Actual data** we are attempting to retrieve
  - I.e. result of the evaluation of a given query expression
- E.g. relation, table, ...

**Query expression**

- **Expression in a given language describing the intended query**
- Multiple expressions for the same query often exist
  - Then they are mutually equivalent

# Sample Query

**First names of all actors born in *1960* or later**

$$\pi_{\mathsf{name}}(\sigma_{\mathsf{year} \geq 1960}(\mathsf{Actor}))$$   $\mathsf{Actor}(\mathsf{year} \geq 1960)[\mathsf{name}]$

| name | surname | year |
|------|---------|------|
| Ivan | Trojan | 1964 |
| Jiří | Macháček | 1966 |
| Jitka | Schneiderová | 1973 |
| Zdeněk | Svěrák | 1936 |
| Jitka | Čvančarová | 1978 |

| name |
|------|
| Ivan |
| Jiří |
| Jitka |

# Relational Algebra

**Inductive construction** of RA expressions

- Basic expressions
    - **Relation name**: $R$
    - **Constant relation**
- General expressions are formed using smaller subexpressions
    - **Projection**: $\pi_{a_1,\ldots,a_n}(E)$
    - **Selection**: $\sigma_\varphi(E)$
    - **Attribute renaming**: $\rho_{b_1/a_1,\ldots,b_n/a_n}(E)$
    - **Union**: $E_R \cup E_S$
    - **Difference**: $E_R \setminus E_S$
    - **Cartesian product**: $E_R \times E_S$
    - …

# Projection

**Projection**: preserves only attributes we are interested in

$$\pi_{a_1,\ldots,a_n}(E) \quad \text{or} \quad E[a_1,\ldots,a_n]$$

- $\langle R, A_R \rangle = [\![E]\!]$ is a relation $R$ with attributes $A_R$
- $a_1,\ldots,a_n$ is a set of **attributes to be preserved**, each $a_i \in A_R$, all the other attributes are to be removed

$$[\![\ \pi_{a_1,\ldots,a_n}(E)\ ]\!] = \langle\ \{t[a_1,\ldots,a_n] \mid t \in R\}\ ,\ \{a_1,\ldots,a_n\}\ \rangle$$

- $t[a_1,\ldots,a_n] = \{(a,v) \mid (a,v) \in t, a \in \{a_1,\ldots,a_n\}\}$ is a restriction of a tuple $t$ to attributes $a_1,\ldots,a_n$
- **Duplicate tuples** in the result are (of course) **suppressed**!

# Projection: Example

**First names of all actors**

$\pi_{\text{name}}(\text{Actor})$    Actor[name]

| name | surname | year |
|------|---------|------|
| Ivan | Trojan | 1964 |
| Jiří | Macháček | 1966 |
| Jitka | Schneiderová | 1973 |
| Zdeněk | Svěrák | 1936 |
| Jitka | Čvančarová | 1978 |

| name |
|------|
| Ivan |
| Jiří |
| Zdeněk |
| Jitka |

# Selection

**Selection**: preserves only tuples we are interested in

$$\sigma_\varphi(E) \quad \text{or} \quad E(\varphi)$$

- $\langle R, A_R \rangle = [\![E]\!]$
- $\varphi$ is a condition (**Boolean expression**) to be satisfied
  - **Connectives**: $\wedge$ (and), $\vee$ (or), $\neg$ (negation)
  - Two forms of **atomic formulae**: $a\Theta b$ or $a\Theta v$
  - $a, b \in A_R$ are **attributes**, $v$ is a value **constant**
  - $\Theta \in \{<, \leq, =, \neq, \geq, >\}$ is a **comparison** operator

$$[\![ \, \sigma_\varphi(E) \, ]\!] = \langle \, \{t \,|\, t \in R, \varphi \text{ is satisfied by } t\} \, , \, A_R \, \rangle$$

# Selection: Example

**Actors born in *1960* or later having a first name other than *Jitka***

$$\sigma_{\text{year} \geq 1960 \,\wedge\, \text{name} \neq \text{Jitka}}(\text{Actor})$$

$$\text{Actor}(\text{year} \geq 1960 \,\wedge\, \text{name} \neq \text{Jitka})$$

| name | surname | year |
|------|---------|------|
| Ivan | Trojan | 1964 |
| Jiří | Macháček | 1966 |
| Jitka | Schneiderová | 1973 |
| Zdeněk | Svěrák | 1936 |
| Jitka | Čvančarová | 1978 |

| name | surname | year |
|------|---------|------|
| Ivan | Trojan | 1964 |
| Jiří | Macháček | 1966 |

# Attribute Renaming

**Rename**: changes names of certain attributes

$$\rho_{b_1/a_1,\ldots,b_n/a_n}(E) \quad \text{or} \quad E\langle a_1 \to b_1, \ldots, a_n \to b_n \rangle$$

- $\langle R, A_R \rangle = [\![E]\!]$
- $a_1, \ldots, a_n$ are **current attributes,** each $a_i \in A_R$, $b_1, \ldots, b_n$ are **new attributes** (distinct)

$$[\![\, \rho_{b_1/a_1,\ldots,b_n/a_n}(E) \,]\!] = \langle\, \{t[b_1/a_1, \ldots, b_n/a_n] \mid t \in R\},$$
$$(A_R \setminus \{a_1, \ldots, a_n\}) \cup \{b_1, \ldots, b_n\} \,\rangle$$

- $t[b_1/a_1, \ldots, b_n/a_n] =$
  $\{(a, v) \mid (a, v) \in t, a \notin \{a_1, \ldots, a_n\}\} \cup$
  $\{(b_i, v) \mid (a_i, v) \in t, i \in \{1, \ldots, n\}\}$

# Attribute Renaming: Example

**Actors with renamed attributes of first and last names**

$\rho_{\text{fname}/\text{name}, \text{lname}/\text{surname}}(\text{Actor})$

$\text{Actor}\langle\text{name} \to \text{fname}, \text{surname} \to \text{lname}\rangle$

| name | surname | year |
|------|---------|------|

| fname | lname | year |
|-------|-------|------|
| Ivan | Trojan | 1964 |
| Jiří | Macháček | 1966 |
| Jitka | Schneiderová | 1973 |
| Zdeněk | Svěrák | 1936 |
| Jitka | Čvančarová | 1978 |

# Set Operations

**Union**, **intersection**, **difference**: standard set operations

$$E_R \cup E_S \qquad [\![\, E_R \cup E_S \,]\!] = \langle\, R \cup S\, ,\, A\, \rangle$$

$$E_R \cap E_S \qquad [\![\, E_R \cap E_S \,]\!] = \langle\, R \cap S\, ,\, A\, \rangle$$

$$E_R \setminus E_S \qquad [\![\, E_R \setminus E_S \,]\!] = \langle\, R \setminus S\, ,\, A\, \rangle$$

- $\langle R, A \rangle = [\![E_R]\!]$ and $\langle S, A \rangle = [\![E_S]\!]$
- Both the relations must be **compatible**
  - I.e. they must have the same attributes

# Set Operations: Difference: Example

**Movies that do not have a good rating**

AllMovies \ GoodMovies

| title | rating |
|---|---|
| Vratné lahve | 76 |
| Samotáři | 84 |
| Medvídek | 53 |
| Štěstí | 72 |

\

| title | rating |
|---|---|
| Samotáři | 84 |
| Kolja | 86 |

=

| title | rating |
|---|---|
| Vratné lahve | 76 |
| Medvídek | 53 |
| Štěstí | 72 |

# Cartesian Product

**Cartesian product** (**cross join**): yields all combinations of tuples among two relations, i.e. unconditionally joins two relations

$$E_R \times E_S$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$
- Both the relations must have **disjoint attributes**
  - ***Dot naming convention*** *based on names of relations is often used in practice, but this approach is not always applicable*
    - $R.a$ and $S.a$ for ambiguous attributes $a$

$$[\![\, E_R \times E_S \,]\!] = \langle\ \{t_1 \cup t_2 \mid t_1 \in R, t_2 \in S\}\ ,\ A_R \cup A_S\ \rangle$$

- Resulting relations are **flat**
  - *I.e. our Cartesian product differs to the one in the set theory*
- Cardinality of the result: $|R|.|S|$

# Cartesian Product: Example

**All possible combinations of movies and actors**

Movie × Actor

| title | rating |
|-------|--------|
| Vratné lahve | 76 |
| Samotáři | 84 |
| Medvídek | 53 |

×

| actor |
|-------|
| Ivan Trojan |
| Jiří Macháček |

=

| title | rating | actor |
|-------|--------|-------|
| Vratné lahve | 76 | Ivan Trojan |
| Vratné lahve | 76 | Jiří Macháček |
| Samotáři | 84 | Ivan Trojan |
| Samotáři | 84 | Jiří Macháček |
| Medvídek | 53 | Ivan Trojan |
| Medvídek | 53 | Jiří Macháček |

# Natural Join

**Natural join**: joins two relations based on the equality of values on all the attributes they mutually share

$$E_R \bowtie E_S \quad \text{or} \quad E_R * E_S$$

- $\langle R, A_R \rangle = [\![ E_R ]\!]$ and $\langle S, A_S \rangle = [\![ E_S ]\!]$

$$[\![ \; E_R \bowtie E_S \; ]\!] = \langle$$

$$\{ t_1 \cup t_2 \mid t_1 \in R, t_2 \in S, \forall\, a \in A_R \cap A_S : t_1(a) = t_2(a) \} \,,$$

$$A_R \cup A_S$$

$$\rangle$$

- When there are no shared attributes (i.e. $A_R \cap A_S = \emptyset$), $\bowtie$ corresponds to $\times$

# Natural Join: Example

**Movie characters with full actor names**

Cast ⋈ Actor      Cast ∗ Actor

| title | actor |
|-------|-------|
| Vratné lahve | 2 |
| Vratné lahve | 4 |
| Samotáři | 1 |
| Medvídek | 1 |
| Medvídek | 2 |

⋈

| actor | name |
|-------|------|
| 1 | Ivan Trojan |
| 2 | Jiří Macháček |
| 3 | Jitka Schneiderová |

=

| title | actor | name |
|-------|-------|------|
| Vratné lahve | 2 | Jiří Macháček |
| Samotáři | 1 | Ivan Trojan |
| Medvídek | 1 | Ivan Trojan |
| Medvídek | 2 | Jiří Macháček |

# Natural Join: Inference

$$E_R \bowtie E_S \ \equiv$$

$$\pi_{r_1,\ldots,r_m,a_1,\ldots,a_n,s_1,\ldots,s_o} \Big($$
$$\sigma_{x_1=a_1 \wedge \cdots \wedge x_n=a_n} \big( \ \rho_{x_1/a_1,\ldots,x_n/a_n}( \ E_R \ ) \ \times \ E_S \ \big)$$
$$\Big)$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$
  - $a_1, \ldots, a_n$ are all the attributes shared by $R$ and $S$
  - $x_1, \ldots, x_n$ are unused attributes, i.e. each $x_i \notin A_R$, $x_i \notin A_S$
  - $r_1, \ldots, r_m$ are all the attributes from $A_R \setminus \{a_1, \ldots, a_n\}$
  - $s_1, \ldots, s_o$ are all the attributes from $A_S \setminus \{a_1, \ldots, a_n\}$

# Theta Join

**Theta join** (Θ**-join**): joins two relations based on a certain condition

$$E_R \bowtie_\varphi E_S \quad \text{or} \quad E_R[\varphi]E_S$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$
- **Disjoint attributes**, i.e. $A_R \cap A_S = \emptyset$
- $\varphi$ is a **condition to be satisfied**
  - Works the same way as conditions in selections

$$[\![ \; E_R \bowtie_\varphi E_S \; ]\!] = \langle$$

$$\{t_1 \cup t_2 \mid t_1 \in R, t_2 \in S, \varphi \text{ is satisfied by } t_1 \cup t_2\},$$

$$A_R \cup A_S$$

$$\rangle$$

# Theta Join

**Inference**

- $E_R \bowtie_\varphi E_S \equiv \sigma_\varphi( E_R \times E_S )$

# Theta Join: Example

**Suitable combinations of movies and actors based on years**

Movie $\bowtie_{\text{filmed} \geq \text{born}}$ Actor     Movie[filmed $\geq$ born]Actor

| title | filmed |
|-------|--------|
| Vratné lahve | 2006 |
| Ecce homo Homolka | 1970 |

$\bowtie_{\varphi}$

| actor | born |
|-------|------|
| Trojan | 1964 |
| Macháček | 1966 |
| Schneiderová | 1973 |

$=$

| title | filmed | actor | born |
|-------|--------|-------|------|
| Vratné lahve | 2006 | Trojan | 1964 |
| Vratné lahve | 2006 | Macháček | 1966 |
| Vratné lahve | 2006 | Schneiderová | 1973 |
| Ecce homo Homolka | 1970 | Trojan | 1964 |
| Ecce homo Homolka | 1970 | Macháček | 1966 |

# Semijoin

**Left / right** (natural) **semijoin**: yields tuples from the left / right relation that can be naturally joined with the other relation

$$E_R \ltimes E_S \quad \text{or} \quad E_R <* E_S \quad / \quad E_R \rtimes E_S \quad \text{or} \quad E_R *> E_S$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$

**Left semijoin**:

$$[\![\ E_R \ltimes E_S\ ]\!] = \langle$$

$$\{t_1 \mid t_1 \in R, \exists\, t_2 \in S : \forall\, a \in A_R \cap A_S : t_1(a) = t_2(a)\},$$

$$A_R$$

$$\rangle$$

**Right semijoin**: analogously

# Semijoin

**Inference**

- $E_R \ltimes E_S \equiv \pi_{r_1,\ldots,r_n}( E_R \bowtie E_S )$
    - where $r_1, \ldots, r_n$ are all attributes from the left relation
- Analogously for the right semijoin

# Semijoin: Example

**Movie characters that have actor details available**

$$\text{Cast} \ltimes \text{Actor} \qquad \text{Cast} <\!* \text{Actor}$$

| title | actor |
|-------|-------|
| Vratné lahve | 2 |
| Vratné lahve | 4 |
| Samotáři | 1 |
| Medvídek | 1 |
| Medvídek | 2 |

$\ltimes$

| actor | name |
|-------|------|
| 1 | Ivan Trojan |
| 2 | Jiří Macháček |
| 3 | Jitka Schneiderová |

$=$

| title | actor |
|-------|-------|
| Vratné lahve | 2 |
| Samotáři | 1 |
| Medvídek | 1 |
| Medvídek | 2 |

# Antijoin

**Left / right antijoin**: yields tuples from the left / right relation that <u>cannot</u> be naturally joined with the other relation

$$E_R \triangleright E_S \ / \ E_R \triangleleft E_S$$

- $\langle R, A_R \rangle = [\![ E_R ]\!]$ and $\langle S, A_S \rangle = [\![ E_S ]\!]$

**Left antijoin**:

$$[\![ \ E_R \triangleright E_S \ ]\!] = \langle$$

$$\{ t_1 \mid t_1 \in R, \neg \exists \, t_2 \in S : \forall \, a \in A_R \cap A_S : t_1(a) = t_2(a) \} \ ,$$

$$A_R$$

$$\rangle$$

**Right antijoin**: analogously

# Antijoin

**Inference**

- $E_R \triangleright E_S \equiv E_R \setminus (E_R \bowtie E_S)$

- Analogously for the right antijoin

# Antijoin: Example

**Movie characters that do not have actor details available**

Cast ▷ Actor

| title | actor |
|-------|-------|
| Vratné lahve | 2 |
| Vratné lahve | 4 |
| Samotáři | 1 |
| Medvídek | 1 |
| Medvídek | 2 |

▷

| actor | name |
|-------|------|
| 1 | Ivan Trojan |
| 2 | Jiří Macháček |
| 3 | Jitka Schneiderová |

=

| title | actor |
|-------|-------|
| Vratné lahve | 4 |

# Theta Semijoin

**Left / right theta semijoin** ($\Theta$-**semijoin**): yields tuples from a given relation that can be joined using a certain condition

$$E_R \ltimes_\varphi E_S \quad \text{or} \quad E_R \langle \varphi \rangle E_S \quad / \quad E_R \rtimes_\varphi E_S \quad \text{or} \quad E_R [\varphi\rangle E_S$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$
- Disjoint attributes, condition $\varphi$ to be satisfied

**Left $\Theta$-semijoin**:

$$[\![\, E_R \ltimes_\varphi E_S \,]\!] = \langle$$

$$\{t_1 \mid t_1 \in R, \exists\, t_2 \in S : \varphi \text{ is satisfied by } t_1 \cup t_2\}\,, \; A_R$$

$$\rangle$$

**Right $\Theta$-semijoin**: analogously

# Theta Semijoin

**Inference**

- $E_R \ltimes_\varphi E_S \equiv \pi_{r_1,\ldots,r_n}( E_R \bowtie_\varphi E_S )$
  - where $r_1, \ldots, r_n$ are all attributes from the left relation
- Analogously for the right $\Theta$-semijoin

# Division

**Division**: returns restrictions of tuples from the first relation such that all combinations of these restricted tuples with tuples from the second relation are present in the first relation

$$E_R \div E_S$$

- $\langle R, A_R \rangle = [\![ E_R ]\!]$ and $\langle S, A_S \rangle = [\![ E_S ]\!]$
- **Assumption on attributes**: $A_S \subset A_R$ (proper subset)

$$[\![ \; E_R \div E_S \; ]\!] = \langle \; \{ t \,|\, \forall\, t_2 \in S : (t \cup t_2) \in R \} \;,\; A_R \setminus A_S \; \rangle$$

- Division allows for the **simulation of a universal quantifier**

# Division: Example

**Movies in which all the actors played**

Cast ÷ Actor

| title | actor |
|-------|-------|
| Vratné lahve | Jiří Macháček |
| Vratné lahve | Zdeněk Svěrák |
| Samotáři | Ivan Trojan |
| Medvídek | Ivan Trojan |
| Medvídek | Jiří Macháček |

÷

| actor |
|-------|
| Ivan Trojan |
| Jiří Macháček |

=

| title |
|-------|
| Medvídek |

# Division: Inference

$$E_R \div E_S \quad \equiv$$

$$\pi_{r_1,\ldots,r_m}(\ E_R\ )\ \setminus$$

$$\pi_{r_1,\ldots,r_m}\Big(\ \big(\ \pi_{r_1,\ldots,r_m}(\ E_R\ )\ \times\ E_S\ \big)\ \setminus\ E_R\ \Big)$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$
- $A_R = \{r_1, \ldots, r_m\} \cup \{s_1, \ldots, s_n\}$ and $A_S = \{s_1, \ldots, s_n\}$
  - I.e. $s_1, \ldots, s_n$ are all the attributes shared by $R$ and $S$, $r_1, \ldots, r_m$ are all the remaining attributes in $R$

# Join Operations

**Inner** joins (and antijoins)

- **Cartesian product**: $E_R \times E_S$
- **Natural join**: $E_R \bowtie E_S$
- **Theta join**: $E_R \bowtie_\varphi E_S$
- Left / right **semijoin**: $E_R \ltimes E_S$, $E_R \rtimes E_S$
- Left / right **antijoin**: $E_R \rhd E_S$, $E_R \lhd E_S$
- Left / right **theta semijoin**: $E_R \ltimes_\varphi E_S$, $E_R \rtimes_\varphi E_S$

**Outer** joins – require an <u>extended relational model with **null** values</u>

- Left / right / full **outer join**: $E_R \ltimes\!\!\bowtie E_S$, $E_R \bowtie\!\!\rtimes E_S$, $E_R \ltimes\!\!\bowtie\!\!\rtimes E_S$

# Outer Join

**Left** / **right** / **full outer join**: natural join of two relations extended by tuples of the first / second / both relations that cannot be joined

$$E_R \text{ ⟗} E_S \quad / \quad E_R \text{ ⟗} E_S \quad / \quad E_R \text{ ⟗} E_S$$

$$E_R *_\mathsf{L} E_S \quad / \quad E_R *_\mathsf{R} E_S \quad / \quad E_R *_\mathsf{F} E_S$$

- $\langle R, A_R \rangle = [\![E_R]\!]$ and $\langle S, A_S \rangle = [\![E_S]\!]$
- $A_R = \{r_1, \dots, r_m\}$, $A_S = \{s_1, \dots, s_n\}$

$$E_R \text{ ⟕} E_S \equiv (E_R \bowtie E_S) \cup ((E_R \triangleright E_S) \times \{(\mathsf{null}, \dots, \mathsf{null})\}_{s_1,\dots,s_n})$$

$$E_R \text{ ⟖} E_S \equiv (E_R \bowtie E_S) \cup (\{(\mathsf{null}, \dots, \mathsf{null})\}_{r_1,\dots,r_m} \times (E_R \triangleleft E_S))$$

$$E_R \text{ ⟗} E_S \equiv (E_R \text{ ⟕} E_S) \cup (E_R \text{ ⟖} E_S)$$

# Outer Join: Example

**Movie characters with full actor names if possible**

Cast ⋈ Actor    Cast ∗_L Actor

| title | actor |
|-------|-------|
| Vratné lahve | 2 |
| Vratné lahve | 4 |
| Samotáři | 1 |
| Medvídek | 1 |

⋈

| actor | name |
|-------|------|
| 1 | Ivan Trojan |
| 2 | Jiří Macháček |
| 3 | Jitka Schneiderová |

=

| title | actor | name |
|-------|-------|------|
| Vratné lahve | 2 | Jiří Macháček |
| Vratné lahve | 4 | null |
| Samotáři | 1 | Ivan Trojan |
| Medvídek | 1 | Ivan Trojan |

# Observations

Relational algebra

- **Declarative** query language
  - Query expressions describe **what data to retrieve**, not (necessarily) how such data should be retrieved
- Both *inputs* and *outputs* of queries are **relations**
- Only values actually present in the database can be returned
  - I.e. **derived data cannot be returned** (such as various calculations, statistics, aggregations, …)

# Observations

**Query evaluation**

- Construction of a **syntactic tree** (query parsing)
  - Based on an inductive structure of a given query expression
  - I.e. based on parentheses (often omitted), operation priorities, associativity conventions, …
- Nodes
  - **Leaf nodes** correspond to individual input relations
  - **Inner nodes** correspond to individual operations
- Evaluation
  - Node can be evaluated when all its child nodes are evaluated, i.e. when all operands of a given operation are available
  - **Root node** represents the result of the entire query

# Observations

**Equivalent expressions**

- Query expressions that **define the same query** (regardless the input relations)
- Various causes
    - **Inference of extended operations** using the basic ones
    - **Commutativity, distributivity or associativity** of (some) operations
    - …
- Examples
    - Commutativity of selection: $(E(\varphi_1))(\varphi_2) \equiv (E(\varphi_2))(\varphi_1)$
    - Selection cascade: $(E(\varphi_1))(\varphi_2) \equiv E(\varphi_1 \wedge \varphi_2)$
    - …

# Observations

**Basic operations**

- **Not all the introduced operations are actually necessary** in order to form expressions of all the possible queries
- The minimal set of required operations:
    - **Projection**: $\pi_{a_1,\ldots,a_n}(E)$
    - **Selection**: $\sigma_\varphi(E)$
    - **Attribute renaming**: $\rho_{b_1/a_1,\ldots,b_n/a_n}(E)$
    - **Union**: $E_R \cup E_S$
    - **Difference**: $E_R \setminus E_S$
    - **Cartesian product**: $E_R \times E_S$

# Observations

**Relational completeness**

- Query language that is able to express all queries of RA is relational complete
- Questions
  - *How to prove that a given language is relational complete?*
  - *Is SQL relational complete?*

# Conclusion

Relational algebra

- **Declarative query language** for the **relational model**
- Operations
  - Basic: **projection**, **selection**, **attribute renaming**, **union**, **difference**, **Cartesian product**
  - Extended: intersection, natural join, theta join, semijoin, antijoin, division, outer join
  - …
- **Relational completeness**
- Motivation
  - Evaluation of SQL queries