

# THE MATHEMATICAL THEORY OF RELATIONAL DATABASES

Joachim Borya

Johannes Kepler Universität

November 30th, 2021

SQL

History

Examples

Time varying  
relations and  
algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

Relations and  
query languages

Syntax

Semantics

- ▶ What is SQL?
  - ▶ Successor of SEQUEL (*Structured English Query Lanugage*)
  - ▶ First database system using SQL: IBM System R (1974)
  - ▶ For users in the between of IT specialists and other people with technical backgrounds.

# Examples

- ▶ Simple SQL database: `sqlite`

## SQL

History

**Examples**

Time varying  
relations and  
algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

Relations and  
query languages

Syntax

Semantics

# Time varying relations and algebra

The mathematical  
theory of relational  
databases

Joachim Borya

SQL

History

Examples

**Time varying  
relations and  
algebra**

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

Relations and  
query languages

Syntax

Semantics

- ▶ Mathematical framework for tables in a relational database
- ▶ Questions:
  - ▶ What is a relation?
  - ▶ What does time-varying mean?
  - ▶ How can relations be manipulated?
    - ▶ "Cartesian product"  $\otimes$
    - ▶ Projection  $\pi$
    - ▶ Selection  $\sigma$
    - ▶ Join  $\bowtie$

# Time varying relations and algebra

- ▶ A *relation schema* is a finite ordered set  $R = (\{A_i\}_{i=1}^n, \leq)$  of *attribute names*.
- ▶ Each attribute name corresponds to a (*simple normal*) *domain*  $\text{dom}(A_i) \subseteq U \in \{\Omega^*, \mathbb{N}\}$ ,  $i \in [n]$ .
- ▶ The set of tuples  $\text{Tup}(R)$  on a relation scheme  $R$  is

$$\left\{ t : R \rightarrow \bigcup_{i=1}^n \text{dom}(A_i) : t(A_i) \in \text{dom}(A_i), i \in [n] \right\}.$$

Because  $R$  is ordered we can associate  $\text{Tup}(R)$  with  $\times_{i=1}^n \text{dom}(A_i)$ .

- ▶ A (*normalized*) *relation*  $r$  on a relation scheme  $R$  is a finite subset of  $\text{Tup}(R)$ . We also write  $r \in \text{Rel}(R)$  and vice versa  $R = \text{Att}(r)$ .

# Normalization of an unnormalized relation

Pet DB		
Owner	Name	Species
Alice	Banjo	Dog
Alice	Bo	Iguana
Alice	Bailey	Cat
Bob	Aj	Cat
Bob	Angus	Cat
Charlie	Callie	Dog
Dave	Flower	Spider

```
{
  "Alice": {
    "Dog": ["Banjo"],
    "Iguana": ["Bo"],
    "Cat": ["Bailey"]
  },
  "Bob": {
    "Cat": ["Aj", "Angus"]
  },
  "Charlie": {
    "Dog": ["Callie"]
  },
  "Dave": {
    "Spider": ["Flower"]
  }
}
```

The mathematical theory of relational databases

Joachim Borya

SQL

History

Examples

Time varying relations and algebra

Basic definitions

Set operations

Projections and normal forms

Selection

Join

Relations and query languages

Syntax

Semantics

# Set operations

## Definition

Let  $R = \{A_1, \dots, A_n\}$  and  $S = \{B_1, \dots, B_m\}$  be relation schemes and  $r \in \text{Rel}(R), s \in \text{Rel}(S)$  relations, then we call  $r$  and  $s$  *union compatible* if  $n = m$  and  $\text{dom}(A_i) = \text{dom}(B_i)$  for every  $i \in [n]$ .

For union-compatible relations we can execute every set operation as usual.

- ▶  $r \cup s := \{t \in \text{Tup}(R) : t \in r \vee t \in s\}$
- ▶  $r \cap s := \{t \in \text{Tup}(R) : t \in r \wedge t \in s\}$
- ▶  $r \setminus s := \{t \in \text{Tup}(R) : t \in r \wedge t \notin s\}$

## SQL

History

Examples

Time varying  
relations and  
algebra

Basic definitions

**Set operations**

Projections and normal  
forms

Selection

Join

Relations and  
query languages

Syntax

Semantics

# Cartesian product

## Definition

Let  $R = \{A_i\}_{i=1}^n$ ,  $S = \{B_i\}_{i=1}^m$  be relation schemes,  $r \in \text{Rel}(R)$ ,  $s \in \text{Rel}(S)$  be relations and  $t_r \in r$ ,  $t_s \in s$  be tuples. We define

$$t_r \circ t_s := (t_r(A_1), \dots, t_r(A_n), t_s(B_1), \dots, t_s(B_m))$$

and

$$r \otimes s := \{t_r \circ t_s : t_r \in r, t_s \in s\}.$$

The relation  $r \otimes s \in \text{Rel}(T)$  is the *cartesian product* of the relations  $r, s$  on the relation schema

$$T = \{C_i\}_{i=1}^k := \{A_1, \dots, A_n, B_1, \dots, B_m\}.$$

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

## Relations and query languages

Syntax

Semantics



# Projection

## Definition

Let  $R = \{A_i\}_{i=1}^n$  be a relation scheme and  $r \in \text{Rel}(R)$  a relation. Let  $S := \{A_{i_k}\}_{k=1}^m \subseteq R$  be a subset of attribute names. The *projection* of  $r$  on  $S$  is defined as

$$\pi_S(r) = \{(t(A_{i_1}), \dots, t(A_{i_m})) : t \in r\}.$$

We immediately see that

$$\pi_S(r) \subseteq \pi_{\{A_{i_1}\}}(r) \otimes \cdots \otimes \pi_{\{A_{i_m}\}}(r).$$

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

## Relations and query languages

Syntax

Semantics

# Functional dependency

## Definition

Let  $R = \{A_i\}_{i=1}^n$  be a relation scheme and  $r \in \text{Rel}(R)$  a relation. Let  $S_1 = \{A_{i_k}\}_{k=1}^m$  and  $S_2 = \{A_{j_k}\}_{k=1}^l$  subsets of  $R$ . We call  $S_2$  *functionally dependent* on  $S_1$  w.r.t.  $r$  if

$$\{((t(A_{i_k}))_{k=1}^m, (t(A_{j_k}))_{k=1}^l) : t \in r\} \subseteq \pi_{S_1}(r) \times \pi_{S_2}(r)$$

is a function. In this case we write  $S_1 \xrightarrow{r} S_2$ . We also define

$$\text{Rel}_{S_1 \rightarrow S_2}(R) = \{r \in \text{Rel}(R) : S_1 \xrightarrow{r} S_2\}.$$

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

## Relations and query languages

Syntax

Semantics

# Keys, prime attributes and second normal form

## Definition

Let  $R = \{A_i\}_{i=1}^n$  be a relation schema,  $r \in \text{Rel}(R)$  a relation and  $K \subset R$ . We call  $K$  a key of the relation  $r$ , if  $K \rightarrow R \setminus K$  and there is no  $K' \subset K$  with  $K' \rightarrow R \setminus K'$ . A superset of a key  $K$  is called a superkey and the attributes  $A \in K$  of a key are called prime. An attribute  $A \in R$  s.t. there is no key  $K$  with  $A \in K$  is called non-prime.

## Definition

Let  $R = \{A_i\}_{i=1}^n$  be a relation schema and  $r \in \text{Rel}(R)$  a relation. Then we say, that  $r$  is in second normal form, if for every non-prime attribute  $A \in R$ , every key  $K \subseteq R$  of  $r$  and  $S \subset K$ ,  $K \rightarrow \{A\}$  and  $S \not\rightarrow \{A\}$  hold.

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

## Relations and query languages

Syntax

Semantics

## Definition

Let  $R = \{A_i\}_{i=1}^n$  a relation schema,  $r \in \text{Rel}(R)$  a relation and  $a \in \text{dom}(A_\nu)$  for some  $A_\nu \in R$ . Then we define the *selection* as

$$\sigma_{A_\nu=a}(r) := \{t \in r : t(A_\nu) = a\}.$$

- ▶  $\sigma_{A=a}(\sigma_{B=b}(r)) = \sigma_{B=b}(\sigma_{A=a}(r))$
- ▶  $\sigma_{A=a}(r \gamma s) = \sigma_{A=a}(r) \gamma \sigma_{A=a}(s)$  for  $\gamma \in \{\cup, \cap, \setminus\}$
- ▶  $\sigma_{A=a \wedge B=b}(r) = \sigma_{A=a}(r) \cap \sigma_{B=b}(r)$

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

**Selection**

Join

## Relations and query languages

Syntax

Semantics

## Definition

Let  $R = \{A_i\}_{i=1}^n$ ,  $S = \{B_i\}_{i=1}^m$  be relation schemes,  
 $r \in \text{Rel}(R)$ ,  $s \in \text{Rel}(S)$  be relations and  $A_\nu \in R$ ,  $B_\mu \in S$   
attribute names. The relation

$$r \overset{A_\nu=B_\mu}{\bowtie} s := \{t_r \circ t_s : t_r \in r \wedge t_s \in s \wedge t_r(A_\nu) = t_s(B_\mu)\}$$

is called (*equi-*)join of  $r$  and  $s$  on  $A_\nu$  and  $B_\mu$ .

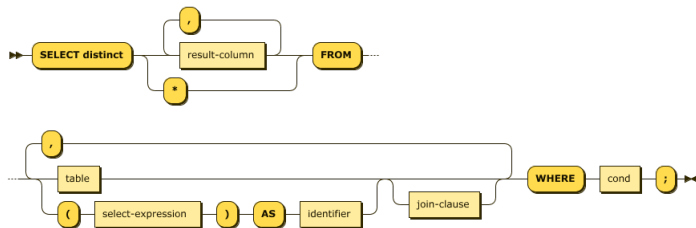
- Selections and joins are exchangeable.

$$r \overset{A_\nu=B_\mu}{\bowtie} s = \sigma_{A_\nu=B_\mu}(r \otimes s)$$

$$r \overset{A_\nu=A}{\bowtie} \{(a)\} = \sigma_{A_\nu=a}(r)$$

for  $a \in \text{dom}(A_\nu)$  and  $\{(a)\} \in \text{Rel}(\{A\})$

# select-expression



- ▶ All valid expressions have this form.
- ▶ The query is restricted to produce distinct tuples.
- ▶ The part `identifier` is an ad-hoc replacement for a `table-name`.
- ▶ The set of all such expression will be called SPJ.

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

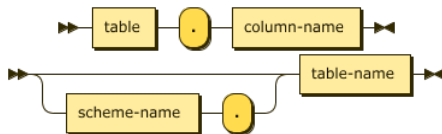
Join

## Relations and query languages

Syntax

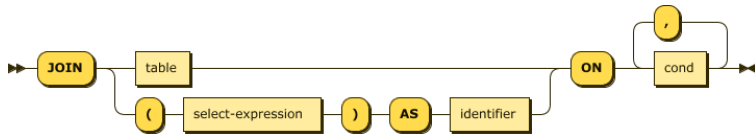
Semantics

# result-column, table



- ▶ The table of a column has to be referenced as well (technically SQL automatically looks this up).
- ▶ The part `scheme-name` has no importance in our considerations.

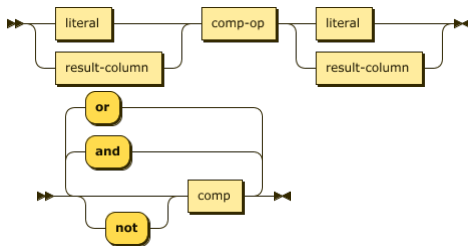
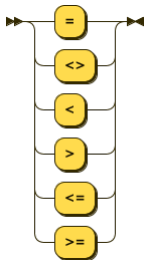
# join-clause



- ▶ This clause provides the rhs of a join.



# comp-op, comp, cond



- ▶ The part `literal` denotes a C-style integer or string literal.
- ▶ The part `cond` is a simple boolean expression involving literals, column names and comparison operators.

# Meaning of a SPJ expression

- ▶ What are the objects we work with? Answer:

$$\mathcal{R} := \left\{ r \subset \prod_{i=1}^n D_i : n \in \mathbb{N}, \forall i \in [n] : D_i \subset \mathbb{N} \vee D_i \subset \Omega^* \right\}$$

- ▶ In reality only finite relations are important.

# Mapping of tables and relations

- ▶ How do we connect the concepts of tables and relations?
- ▶ Let  $\{R_i\}_{i=1}^m$  be relation schemes, i.e.

$$R_i := \{A_{i1}, \dots, A_{in_i}\} \text{ for } i \in [m]$$

and  $r_1 \in \text{Rel}(R_1), \dots, r_m \in \text{Rel}(R_m)$  relations.

- ▶ Consider a set of variables e.g. identifiers

$$V := \{\text{tid}_1, \dots, \text{tid}_m, \text{tid}_1.\text{att}_1, \dots, \text{tid}_1.\text{att}_{n_1}, \\ \dots, \text{tid}_m.\text{att}_1, \dots, \text{tid}_m.\text{att}_{n_m}\}$$

and a map  $\mathcal{D} : V \rightarrow \{r_i\}_{i=1}^m \cup \bigcup \{R_i\}_{i=1}^m$  given by

$$\mathcal{D}(\text{tid}_1) := r_1, \dots, \mathcal{D}(\text{tid}_m) := r_m,$$

$$\mathcal{D}(\text{tid}_1.\text{att}_1) := A_{11}, \dots, \mathcal{D}(\text{tid}_1.\text{att}_{n_1}) := A_{1n_1}, \\ \dots, \mathcal{D}(\text{tid}_m.\text{att}_1) := A_{m1}, \dots, \mathcal{D}(\text{tid}_m.\text{att}_{n_m}) := A_{mn_m}$$

# Semantic function definition

- ▶ Goal: Construction of a map  $\llbracket \cdot \rrbracket_D : \text{SPJ} \rightarrow \mathcal{R}$ .

# Semantic function

for table expressions  $t$

$$\llbracket \text{SELECT } * \text{ FROM } t \rrbracket_{\mathcal{D}} := \mathcal{D}(t).$$

for select-expressions  $S$

$$\llbracket \text{SELECT } * \text{ FROM } (S) \text{ AS subquery} \rrbracket_{\mathcal{D}} =$$

$$\llbracket \text{SELECT } * \text{ FROM subquery} \rrbracket_{\mathcal{D}^*} \text{ with}$$

$$\mathcal{D}^*(E) := \begin{cases} \llbracket E \rrbracket_{\mathcal{D}}, & \text{if } E = \text{subquery} \\ A, & \text{if } E = \text{subquery}.* \\ \mathcal{D}(E), & \text{else} \end{cases}$$

The attribute names of  $\llbracket \text{subquery} \rrbracket_{\mathcal{D}}$  are inherited from the tables involved in the construction of the subquery but with the prefix subquery. Note that  $A \in \text{Att}(\llbracket \text{subquery} \rrbracket_{\mathcal{D}})$ .

# Example

```
QS := SELECT adult.name, adult.address  
FROM (S)  
) AS adult;
```

```
S = SELECT owner.name, owner.address  
FROM owner  
WHERE owner.age > 17;
```

## Meaning

$$\text{Att}(\llbracket S \rrbracket_{\mathcal{D}}) = \{N, A\}$$

$$\mathcal{D}^* := \mathcal{D} \cup [\text{adult} \mapsto \llbracket S \rrbracket_{\mathcal{D}}, \text{adult.name} \mapsto N, \text{adult.address} \mapsto A]$$
$$\llbracket Q_S \rrbracket_{\mathcal{D}} = \llbracket \text{SELECT adult.name, adult.address FROM adult} \rrbracket_{\mathcal{D}^*}$$

# Semantic function

for table or select-expressions  $F_1, \dots, F_n$

$$\llbracket \text{SELECT } * \text{ FROM } F_1, \dots, F_n \rrbracket_{\mathcal{D}} =$$
$$\bigotimes_{i=1}^n \llbracket \text{SELECT } * \text{ FROM } F_i \rrbracket_{\mathcal{D}}$$

for result-column expressions  $c_1, \dots, c_n$

$$\llbracket \text{SELECT } c_1, \dots, c_n \text{ FROM } F \rrbracket_{\mathcal{D}} =$$
$$\pi_{(\mathcal{D}(c_1), \dots, \mathcal{D}(c_n))}(\llbracket \text{SELECT } * \text{ FROM } F \rrbracket_{\mathcal{D}})$$

It is only defined this way if

$$\mathcal{D}(c_1), \dots, \mathcal{D}(c_n) \in \text{Att}(\llbracket \text{SELECT } * \text{ FROM } F \rrbracket_{\mathcal{D}}).$$

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

## Relations and query languages

Syntax

Semantics

# Semantic function

for a list of table or select-expressions  $F$

$$\llbracket \text{SELECT } * \text{ FROM } F \text{ WHERE } c = a \rrbracket_{\mathcal{D}} = \sigma_{\mathcal{D}(c)=\bar{a}}(\llbracket \text{SELECT } * \text{ FROM } F \rrbracket_{\mathcal{D}})$$

It is only defined this way if  $\mathcal{D}(c) \in \text{Att}(F)$  and  $a \in \text{dom}(\mathcal{D}(c))$ .

- ▶  $\llbracket \text{SELECT } * \text{ FROM } F \text{ WHERE } c < a \rrbracket_{\mathcal{D}} = \bigcup_{i=1}^{\infty} \sigma_{\mathcal{D}(c)=\bar{a}-i}(\llbracket \text{SELECT } * \text{ FROM } F \rrbracket_{\mathcal{D}})$
- ▶  $\llbracket \text{SELECT } * \text{ FROM } F \text{ WHERE } c \leq a \rrbracket_{\mathcal{D}} = \bigcup_{i=0}^{\infty} \sigma_{\mathcal{D}(c)=\bar{a}-i}(\llbracket \text{SELECT } * \text{ FROM } F \rrbracket_{\mathcal{D}})$

## SQL

History

Examples

## Time varying relations and algebra

Basic definitions

Set operations

Projections and normal  
forms

Selection

Join

## Relations and query languages

Syntax

Semantics



# Semantic function

for table or select-expressions  $F_1, F_2$  and  
 $c_1 \in \text{Att}(F_1), c_2 \in \text{Att}(F_2)$

$$\begin{aligned} & \llbracket \text{SELECT } * \text{ FROM } F_1 \text{ JOIN } F_2 \text{ ON } c_1 = c_2 \rrbracket_{\mathcal{D}} \\ &= \llbracket \text{SELECT } * \text{ FROM } F_1 \rrbracket_{\mathcal{D}} \overset{\mathcal{D}(c_1)=\mathcal{D}(c_2)}{\bowtie} \llbracket \text{SELECT } * \text{ FROM } F_2 \rrbracket_{\mathcal{D}} \\ &= \sigma_{\mathcal{D}(c_1)=\mathcal{D}(c_2)}(\llbracket \text{SELECT } * \text{ FROM } F_1 \rrbracket_{\mathcal{D}} \otimes \llbracket \text{SELECT } * \text{ FROM } F_2 \rrbracket_{\mathcal{D}}) \\ &= \llbracket \text{SELECT } * \text{ FROM } F_1, F_2 \text{ WHERE } c_1 = c_2 \rrbracket_{\mathcal{D}} \end{aligned}$$

# Example

```
SELECT inventory.desc, inventory.stock
FROM inventory
JOIN (SELECT * FROM product WHERE product.type = "fruit")
AS fruits
ON inventory.desc = fruits.desc
WHERE inventory.stock < 8;
```

# Example

```
[[ SELECT inventory.desc, inventory.stock
FROM inventory
JOIN (SELECT * FROM product WHERE product.type =
"fruit") AS fruits
ON inventory.desc = fruits.desc
WHERE inventory.stock < 8 ]]D
```

```
=  $\pi(\mathcal{D}(\text{inventory.desc}), \mathcal{D}(\text{inventory.stock}))$ (
[[ SELECT * FROM inventory
JOIN (SELECT * FROM product WHERE product.type =
"fruit") AS fruits
ON inventory.desc = fruits.desc
WHERE inventory.stock < 8 ]]D
)
```

# Example





```
=  $\pi(\mathcal{D}(\text{inventory.desc}), \mathcal{D}(\text{inventory.stock}))$ (  
 $\sigma_{\mathcal{D}(\text{inventory.stock}) < 8}$ (  
[[ SELECT * FROM inventory  
JOIN (SELECT * FROM product WHERE product.type =  
"fruit") AS fruits  
ON inventory.desc = fruits.desc ]]D  
))
```

```
=  $\pi(\mathcal{D}(\text{inventory.desc}), \mathcal{D}(\text{inventory.stock}))$ (  
 $\sigma_{\mathcal{D}(\text{inventory.stock}) < 8}$ ( [[ SELECT * FROM inventory ]]D  
 $\bowtie_{\mathcal{D}(\text{inventory.desc}) = \mathcal{D}(\text{fruits.desc})}$   
[[ SELECT *  
FROM (SELECT * FROM product WHERE product.type =  
"fruit") AS fruits ]]D  
))
```

# Example

$$\begin{aligned} &= \pi(\mathcal{D}(\text{inventory.desc}), \mathcal{D}(\text{inventory.stock}))(\sigma_{\mathcal{D}(\text{inventory.stock}) < 8}(\llbracket \text{SELECT } * \text{ FROM inventory } \rrbracket_{\mathcal{D}} \\ &\bowtie_{\mathcal{D}(\text{inventory.desc}) = \mathcal{D}(\text{fruits.desc})} \llbracket \text{SELECT } * \text{ FROM product WHERE product.type = } \\ &\text{"fruit"} \rrbracket_{\mathcal{D}} \\ &)) \end{aligned}$$
$$\begin{aligned} &= \pi(\mathcal{D}(\text{inventory.desc}), \mathcal{D}(\text{inventory.stock}))(\sigma_{\mathcal{D}(\text{inventory.stock}) < 8}(\mathcal{D}(\text{inventory}) \\ &\bowtie_{\mathcal{D}(\text{inventory.desc}) = \mathcal{D}(\text{fruits.desc})} \sigma_{\mathcal{D}(\text{product.type} = \text{"fruit"})}(\mathcal{D}(\text{product}))) \\ &)) \end{aligned}$$
$$\begin{aligned} &= \pi(\mathcal{D}(\text{inventory.desc}), \mathcal{D}(\text{inventory.stock}))(\bigcup_{i=1}^{\infty} \sigma_{\mathcal{D}(\text{inventory.stock}) = 8-i}(\mathcal{D}(\text{inventory}) \\ &\bowtie_{\mathcal{D}(\text{inventory.desc}) = \mathcal{D}(\text{fruits.desc})} \sigma_{\mathcal{D}(\text{product.type} = \text{"fruit"})}(\mathcal{D}(\text{product}))) \\ &)) \end{aligned}$$

# Sources

-  E. F. Codd, "Further normalization of the database relational model", 1971.
-  Gottfried Vossen, "Datenbankmodelle, Datenbanksprachen und Datenbankmanagementsysteme", 2008.
-  David Maier, "The theory of relational databases", 1983.
-  [https://sqlite.org/lang\\_select.html](https://sqlite.org/lang_select.html)

Syntax diagrams created with Railroad Diagram Generator  
<https://www.bottlecaps.de/rr/ui>