

Database Theory

VU 181.140, SS 2018

3. Codd's Theorem

Reinhard Pichler

Institute of Logic and Computation
DBAI Group
TU Wien

20 March, 2018



What is Codd's Theorem?

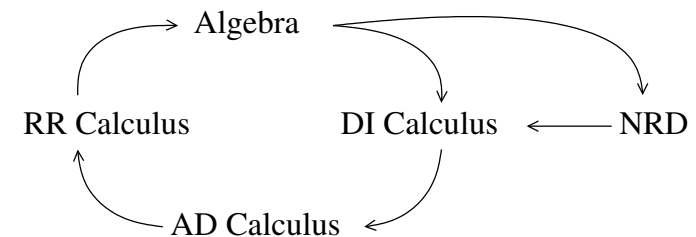
- Several query languages have been considered for relational databases:
 - some are more convenient for query specification
 - some are easier to optimize
 - some are more succinct
- (Original) **Codd's Theorem** shows that the most important of the above languages, the so-called **domain-independent, relational calculus** and relational algebra, are **equally expressive**.
- By a set of translations we prove here a **stronger** result. The following languages are equally expressive:
 - domain-independent relational calculus,
 - relational calculus under the active domain semantics,
 - range-restricted relational calculus,
 - relational algebra,
 - non-recursive Datalog with negation,

Outline

3. Codd's Theorem

- 3.1 What is Codd's Theorem?
- 3.2 Domain Independent Relational Calculus
- 3.3 From the Algebra to DI Calculus
- 3.4 Active Domain Semantics
- 3.5 Range Restricted Queries
- 3.6 From Range Restricted Queries to Algebra
- 3.7 Non-recursive Datalog with Negation

Proof Strategy



Domain Independence

Some queries are “unsafe” and must be avoided:

- **Domain Independence:** Given a query and a database, the query must evaluate to the same result on the database no matter what the domain is assumed to be.
- Idea: exclude “unsafe” queries, i.e., in particular, queries that may yield an infinite answer.
- Let $Q_B(\mathcal{A})$ denote the result of evaluating query Q on database \mathcal{A} assuming domain B .

Definition (domain-independence)

A query Q is *domain-independent* iff there do **not** exist

- a database instance \mathcal{A} and
- two sets B, C that contain all constants that appear in \mathcal{A} or in Q (also known as the **active domain**), such that $Q_B(\mathcal{A}) \neq Q_C(\mathcal{A})$.

Undecidability of Domain Independence

Definition

For a domain-independent query Q and a database \mathcal{A} , we may define $Q(\mathcal{A}) := Q_B(\mathcal{A})$ for an arbitrary domain B that contains the active domain (the choice is irrelevant due to domain-independence).

- We would like to require domain-independence in queries.
- Domain-independence is an **undecidable** property for FO queries.
- Possible solutions:
 - semantic restriction on quantification
 - syntactic restriction on queries

Queries Violating Domain Independence

Example (Unsafe Queries)

- $\{x \mid \neg R(x)\}$
 - $R = \emptyset, 1 \in B, 1 \notin C: 1 \in Q_B(R), 1 \notin Q_C(R)$.
- $\{x \mid R(x) \vee R(y)\}$
- $\{y \mid R(x)\}$
- $\{x \mid R(x) \vee \neg R(x)\}$

Remark. Over infinite domains, these queries may yield an **infinite result**.

Theorem

For every relational algebra expression there exists a domain-independent, relational calculus query.

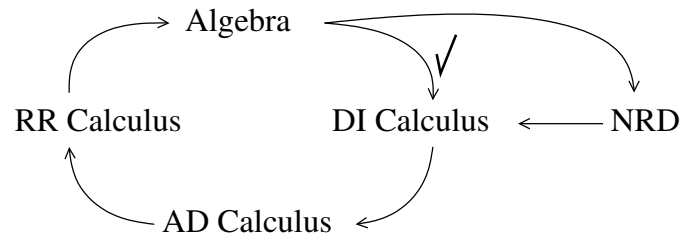
Proof.

Almost by definition:

$$\begin{aligned}
 R &:= \{\vec{x} \mid R(\vec{x})\} \\
 \sigma_{x=y}(\{\vec{x} \mid \varphi(\vec{x})\}) &:= \{\vec{x} \mid \varphi(\vec{x}) \wedge x = y\} \\
 \pi_{\vec{y}}(\{\vec{x} \mid \varphi(\vec{x})\}) &:= \{\vec{y} \mid \exists \vec{z} \varphi(\vec{x})\} \quad (\vec{x} = \vec{y}\vec{z}) \\
 \{\vec{x} \mid \varphi(\vec{x})\} \times \{\vec{y} \mid \psi(\vec{y})\} &:= \{\vec{x}\vec{y} \mid \varphi(\vec{x}) \wedge \psi(\vec{y})\} \\
 \{\vec{x} \mid \varphi(\vec{x})\} \cup \{\vec{x} \mid \psi(\vec{x})\} &:= \{\vec{x} \mid \varphi(\vec{x}) \vee \psi(\vec{x})\} \\
 \{\vec{x} \mid \varphi(\vec{x})\} - \{\vec{x} \mid \psi(\vec{x})\} &:= \{\vec{x} \mid \varphi(\vec{x}) \wedge \neg\psi(\vec{x})\}
 \end{aligned}$$

It can be shown by an easy induction argument that the resulting relational calculus query is domain-independent. \square

Current Status



Proof.

Assume a domain-independent query Q . We claim that Q itself is the desired query, such that its evaluation under standard semantics and under active domain semantics coincides.

Consider an arbitrary database \mathcal{A} .

By domain-independence, $Q_B(\mathcal{A}) = Q_{\text{adom}(Q, \mathcal{A})}(\mathcal{A})$ under standard semantics for every domain B with $\text{adom}(Q, \mathcal{A}) \subseteq B$.

Likewise, under active domain semantics, $Q_B(\mathcal{A}) = Q_{\text{adom}(Q, \mathcal{A})}(\mathcal{A})$ for every domain B with $\text{adom}(Q, \mathcal{A}) \subseteq B$.

Finally, $Q_{\text{adom}(Q, \mathcal{A})}(\mathcal{A})$ yields the same result under standard semantics and under active domain semantics. Hence, for arbitrary domain B with $\text{adom}(Q, \mathcal{A}) \subseteq B$, the set $Q_B(\mathcal{A})$ is the same under standard semantics and under active domain semantics. \square

Active Domain Interpretation

Definition

Active domain semantics for relational calculus query q over DB \mathcal{A} : variables range over the **active domain**, i.e. over values occurring in the database \mathcal{A} and the query q , denoted by $\text{adom}(q, \mathcal{A})$.

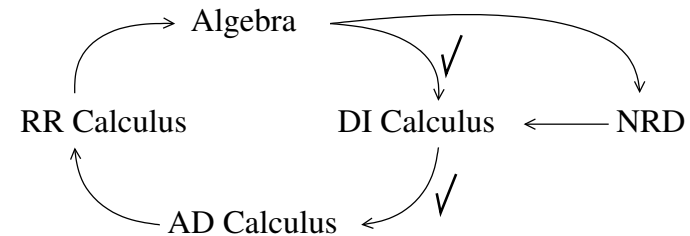
Example

- Assume a database \mathcal{A} with $R = \{\langle 1, 1 \rangle\}$ and $\text{dom}(R) = \{1, 2\}$.
- Consider the Boolean query $q = \{\langle \rangle \mid \forall x, y. R(x, y)\}$:
 - q is false in \mathcal{A} under the standard semantics, but
 - q is true in \mathcal{A} under the active domain semantics.

Theorem

For every domain-independent relational calculus query there is an equivalent relational calculus query under the active domain semantics.

Current Status



Range-restricted Queries Capture AD Calculus (continued)

Proof.

We now translate φ into an rr query φ' such that answering φ under the active domain semantics is equivalent to evaluating φ' under the standard semantics. (We write $tr(\psi)$ to denote the translation of a formula ψ .)

- 1 Turn φ into SRNF.
- 2 We build φ' from φ as follows:

$$\varphi' := D(x_1) \wedge \dots \wedge D(x_n) \wedge tr(\varphi),$$

where $\{x_1, \dots, x_n\} = free(\varphi)$ and

$$tr(A) := A \quad (A \text{ is an atom})$$

$$tr(\varphi \wedge \psi) := tr(\varphi) \wedge tr(\psi)$$

$$tr(\varphi \vee \psi) := tr(\varphi) \vee tr(\psi)$$

$$tr(\neg\varphi) := \neg tr(\varphi)$$

$$tr(\exists x \varphi) := (\exists x) (D(x) \wedge tr(\varphi))$$

□

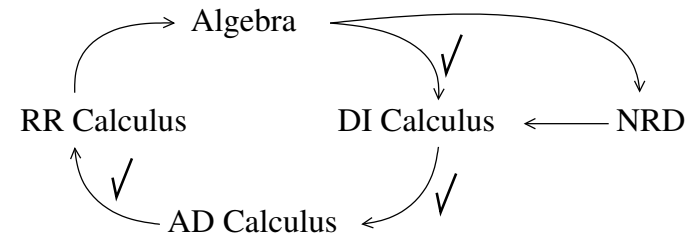
Theorem

For every range-restricted relational calculus query there exists an equivalent relational algebra expression.

Proof idea.

- 1 Put the query into SRNF.
- 2 Put the query into **Relational-Algebra NF (RANF)**:
 - RANF: Each subformula is **range-restricted**. (Exception: In a subformula $\pi = \varphi_1 \wedge \dots \wedge \varphi_k \wedge \neg\psi$, π has to be RR and the φ_i and ψ have to be in RANF, but $\neg\psi$ does not need to be RR).
- 3 Translate the RANF formula into relational algebra. This can be done inductively from the leaves to the root of the parse tree of the formula.

Current Status



From SRNF to RANF

Given a formula in SRNF.

- The formula is in RANF if each subformula is range-restricted.
- Possible obstacles: subformulae of the form $\varphi \vee \psi$ or $\neg\varphi$.
- Only these remove possibly relevant variables from rr.
- Solution: relativize using the active domain relation D :

$$\varphi \vee \psi \vdash \underbrace{(\varphi \wedge \bigwedge_{x \in (free(\psi) - free(\varphi))} D(x)) \vee (\psi \wedge \bigwedge_{x \in (free(\varphi) - free(\psi))} D(x))}_{rr(*)=free(*)=free(\varphi) \cup free(\psi)}$$

$$\neg\varphi \vdash (\neg\varphi \wedge \bigwedge_{x \in free(\varphi)} D(x))$$

- Shorter (=better) RANF queries can be achieved by rewriting the input formula using equivalences we already know.

From RANF to Relational Algebra

RANF formulae can be translated to relational algebra using the following rules:

$$\begin{aligned}
 \text{Alg}(\varphi \wedge \psi) &:= \text{Alg}(\varphi) \bowtie \text{Alg}(\psi) \\
 \text{Alg}(\varphi \wedge \neg\psi) &:= \text{Alg}(\varphi) - \text{Alg}(\psi) \\
 &\dots \varphi \text{ and } \psi \text{ have the same schema} \\
 \text{Alg}(\exists y \varphi(\vec{x}, y)) &:= \pi_{\vec{x}} \text{Alg}(\varphi(\vec{x}, y)) \\
 \text{Alg}(\varphi \wedge x \vartheta t) &:= \sigma_{x \vartheta t} \text{Alg}(\varphi) \\
 &\dots \vartheta \text{ is either } = \text{ or } \neq \text{ and } t \text{ is a term.} \\
 \text{Alg}(R(x_1, \dots, x_k)) &:= \rho_{A_1 \dots A_k \rightarrow x_1 \dots x_k} R \\
 &\dots \text{ relation } R \text{ has schema } R(A_1, \dots, A_k)
 \end{aligned}$$

From RR Queries to the Algebra

Example

Range-restricted but not in RANF (i.e., not locally range-restricted):

$$\{(x, y) \mid \exists z : P(x, y, z) \vee (R(x, y) \wedge \overbrace{((S(z) \wedge \neg T(x, z)) \vee T(y, z))}^{rr(*)=\{z\}})\}$$

We transform this formula using the rewrite rule

$$\varphi \wedge (\psi_1 \vee \psi_2) \vdash (\varphi \wedge \psi_1) \vee (\varphi \wedge \psi_2)$$

into RANF:

$$\{(x, y) \mid \exists z : P(x, y, z) \vee (R(x, y) \wedge S(z) \wedge \neg T(x, z)) \vee (R(x, y) \wedge T(y, z))\}$$

From RR Queries to the Algebra

Example

Let D be the active domain relation with schema $D(D)$ and let R have schema $R(A, B)$. The formula

$$\exists x \left(D(x) \wedge \exists y (D(y) \wedge \neg R(x, y)) \right)$$

corresponds to the RANF formula

$$\exists x \exists y \left((D(x) \wedge D(y)) \wedge \neg R(x, y) \right).$$

An equivalent relational algebra expression looks as follows.

$$\begin{aligned}
 &\text{Alg}(\exists x \exists y ((D(x) \wedge D(y)) \wedge \neg R(x, y))) \\
 \vdash &\pi_{\emptyset} (\text{Alg}((D(x) \wedge D(y)) - \text{Alg}(R(x, y)))) \\
 \vdash &\pi_{\emptyset} ((\rho_x D \bowtie \rho_y D) - \rho_{xy} R)
 \end{aligned}$$

From RR Queries to the Algebra

Example (in RANF)

$$\{(x, y) \mid \exists z : \overbrace{P(x, y, z)}^{e_1=\rho_{xyz} P} \vee \underbrace{\overbrace{((\rho_{xy} R \wedge \rho_z S) \wedge \neg \overbrace{T(x, z)}^{e_{21}=\rho_{xz} T}) \vee \overbrace{(R(x, y) \wedge \overbrace{T(y, z)}^{\rho_{yz} T})}^{e_3=(\cdot) \bowtie (\cdot)})}_{e_2=(\cdot) - ((\rho_y D) \bowtie (\cdot))}}_{(\cdot) \bowtie (\cdot)}\}$$

Equivalent relational algebra expression: $\pi_{xy}(e_1 \cup e_2 \cup e_3)$.

From RR Queries to the Algebra

Example

"Select all professors (P) who only give lectures (L) in the field of computer science (C)." The schemata are $P(P)$, $L(P, C)$, $C(C)$ and the active domain is given by a relation D with schema $D(D)$.

$$\begin{aligned} & \{x \mid P(x) \wedge \forall y(L(x, y) \rightarrow C(y))\} \\ \text{to SRNF} & \vdash \{x \mid P(x) \wedge \neg \exists y(L(x, y) \wedge \neg C(y))\} \\ \text{to RANF} & \vdash \{x \mid P(x) \wedge \neg \exists y(L(x, y) \wedge \underbrace{(D(y) \wedge \neg C(y))}_{(\rho_C D) - C})\} \\ & \underbrace{\hspace{10em}}_{L \bowtie ((\rho_C D) - C)} \\ & \underbrace{\hspace{10em}}_{\pi_P(L \bowtie ((\rho_C D) - C))} \\ & \underbrace{\hspace{10em}}_{P - \pi_P(L \bowtie ((\rho_C D) - C))} \end{aligned}$$

From RR Queries to the Algebra

Example

Schema $R(AB)$, $S(C)$, $T(AC)$; active domain $D(D)$.

$$\begin{aligned} & \{ \langle x, y, z \rangle \mid R(x, y) \wedge (S(z) \wedge \neg T(x, z)) \} \\ \vdash & \{ \langle x, y, z \rangle \mid R(x, y) \wedge \underbrace{((D(x) \wedge S(z)) \wedge \neg T(x, z))}_{(D \times S) - T} \} \\ & \underbrace{\hspace{10em}}_{R \bowtie ((D \times S) - T)} \end{aligned}$$

From RR Queries to the Algebra

Example

The schemata are $L(P, C)$, $C(C)$ and the domain is given by $D(D)$.

$$\begin{aligned} \{ \langle x, y \rangle \mid L(x, y) \wedge \neg C(y) \} & \equiv \{ \langle x, y \rangle \mid \underbrace{L(x, y) \wedge (D(y) \wedge \neg C(y))}_{L \bowtie ((\rho_C D) - C)} \} \\ & \equiv \{ \langle x, y \rangle \mid \underbrace{L(x, y) \wedge \neg (D(x) \wedge C(y))}_{L - ((\rho_C D) \times C)} \} \end{aligned}$$

D	D
1	
2	
3	
4	

L	P	C
	1	2
	3	4

C	C
	2

(ρ _C D) - C	C
	1
	3
	4

(ρ _C D) × C	P	C
	1	2
	2	2
	3	2
	4	2

L ⋈ ((ρ _C D) - C)	P	C
	3	4

From RR Queries to the Algebra

Example

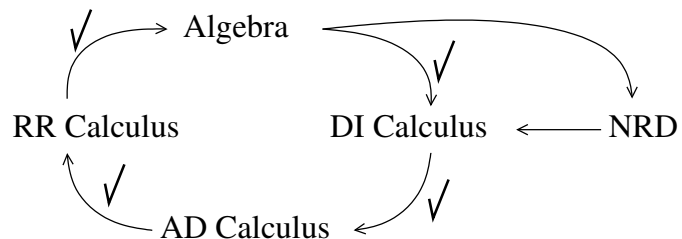
Schema $R(AB)$, $S(AC)$, $T(BC)$; active domain $D(D)$.

$$\begin{aligned} & \{ \langle x, y, z \rangle \mid R(x, y) \wedge (S(x, z) \vee T(y, z)) \} \\ \vdash & \{ \langle x, y, z \rangle \mid \underbrace{(R(x, y) \wedge S(x, z))}_{R \bowtie S} \vee \underbrace{(R(x, y) \wedge T(y, z))}_{R \bowtie T} \} \\ & \underbrace{\hspace{10em}}_{(R \bowtie S) \cup (R \bowtie T)} \\ \text{or} & \{ \langle x, y, z \rangle \mid R(x, y) \wedge (S(x, z) \vee T(y, z)) \} \\ \vdash & \{ \langle x, y, z \rangle \mid \underbrace{R(x, y) \wedge ((S(x, z) \wedge D(y)) \vee (T(y, z) \wedge D(x)))}_{R \bowtie ((S \times \rho_B D) \cup (T \times \rho_A D))} \} \end{aligned}$$

This is correct because

$$\begin{aligned} R(x, y) \wedge (\varphi \vee \psi) & \equiv R(x, y) \wedge D(x) \wedge D(y) \wedge (\varphi \vee \psi) \equiv \\ R(x, y) \wedge ((D(x) \wedge D(y) \wedge \varphi) \vee (D(x) \wedge D(y) \wedge \psi)). \end{aligned}$$

Current Status

From Algebra to $nr\text{-Datalog}^-$

Proof.

$$P_E := \{Q_E(\vec{x}) :- R(\vec{x})\} \text{ if } E \text{ is a relation } R$$

$$P_{\sigma_{x=y}(E)} := \{Q_{\sigma_{x=y}(E)}(\vec{x}) :- Q_E(\vec{x}) \wedge x = y\} \cup P_E$$

$$P_{\pi_{\vec{y}}(E)} := \{Q_{\pi_{\vec{y}}(E)}(\vec{y}) :- Q_E(\vec{x})\} \cup P_E$$

$$P_{E_1 \times E_2} := \{Q_{E_1 \times E_2}(\vec{x}, \vec{y}) :- Q_{E_1}(\vec{x}), Q_{E_2}(\vec{y})\} \cup P_{E_1} \cup P_{E_2}$$

$$P_{E_1 \cup E_2} := \{Q_{E_1 \cup E_2}(\vec{x}) :- Q_{E_1}(\vec{x})\} \cup \\ \{Q_{E_1 \cup E_2}(\vec{x}) :- Q_{E_2}(\vec{x})\} \cup P_{E_1} \cup P_{E_2}$$

$$P_{E_1 - E_2} := \{Q_{E_1 - E_2}(\vec{x}) :- Q_{E_1}(\vec{x}), \text{not } Q_{E_2}(\vec{x})\} \cup P_{E_1} \cup P_{E_2}$$

Clearly, E and P_E have the same answer, i.e. for any database DB and any constant tuple \vec{c} , we have $\vec{c} \in M[E](DB)$ iff $P_E^* \wedge DB^* \models Q_E(\vec{c})$. \square

Non-recursive Datalog with negation

Definition

Non-recursive Datalog with negation ($nr\text{-Datalog}^-$) prohibits cycles in the dependency graph $DEP(P)$ of a program P .

- non-recursiveness means that negation is trivially stratified!

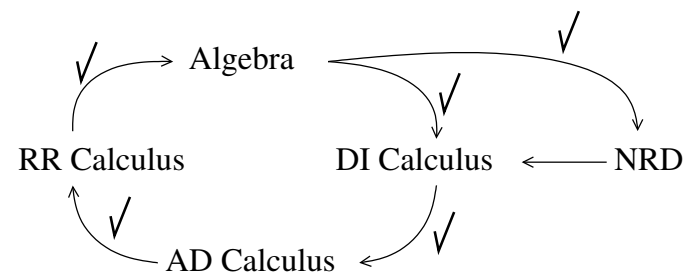
Theorem

For every relational algebra query there exists an equivalent $nr\text{-Datalog}^-$ query.

Proof idea.

- Inductively, for each algebra expression E , we construct the program P_E that defines the predicate Q_E of the same arity as E .
- We proceed by a bottom-up traversal of the syntax tree of algebra expression E , introducing a new predicate for each subexpression of E . Hence, the resulting program is clearly non-recursive.

Current Status



From nr-Datalog[¬] to DI Calculus

Theorem

For every query in non-recursive Datalog with negation there exists an equivalent domain-independent relational calculus query.

Proof idea.

- Let us assume an nr-Datalog[¬] program P .
- For every predicate R occurring in P we define a formula $\varphi_{P,R}$.
- $\varphi_{P,R}$ captures the query defined by the program P "restricted to predicate R ", i.e., for any database DB the following are equal:
 - the answer to $\{\vec{x} \mid \varphi_{P,R}(\vec{x})\}$, where \vec{x} are the free variables of $\varphi_{P,R}$;
 - the set of constant tuples \vec{c} with $P^* \wedge DB^* \models R(\vec{c})$.

Proof.

Suppose that the following are the rules with head predicate Q :

$$\begin{aligned} Q(\vec{x}) & :- \alpha_{0,0}, \dots, \alpha_{0,n_0} \\ & \vdots \\ Q(\vec{x}) & :- \alpha_{m,0}, \dots, \alpha_{m,n_m} \end{aligned}$$

Then we set

$$\varphi_{P,Q} = \bigvee_{i \in \{0, \dots, m\}} (\exists \vec{y}_i : \bigwedge_{j \in \{0, \dots, n_i\}} \beta_{i,j}),$$

where \vec{y}_i are the existential variables of the i th rule (i.e. the variables not occurring in the rule head), and

$$\beta_{i,j} = \begin{cases} \varphi_{P,W}(v) & \text{if } \alpha_{i,j} \text{ is an atom } W(v). \\ \neg \varphi_{P,W}(v) & \text{if } \alpha_{i,j} \text{ is an atom not } W(v). \end{cases} \quad \square$$

From nr-Datalog[¬] to DI Calculus (continued)

Proof

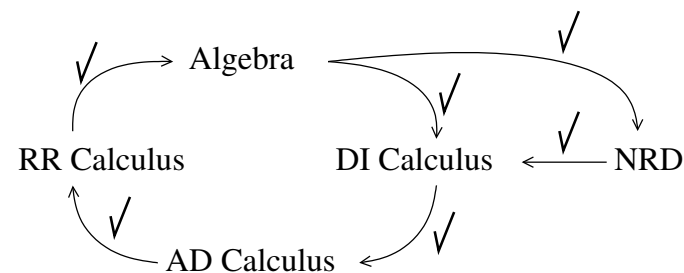
W.l.o.g., we assume:

- P has no constants and no multiple occurrences of variables in rule heads: this can be simulated using (fresh) variables and $=$.
- for every pair $H_1(\vec{x}_1), H_2(\vec{x}_2)$ of head atoms in P , $H_1 = H_2$ implies $\vec{x}_1 = \vec{x}_2$ (can be achieved by variable renaming)

Inductive definition of $\varphi_{P,R}$:

- (Base case) If R is a predicate that does not appear in the head of any rule in P , then $\varphi_{P,R} = R(\vec{x})$.
- (Inductive step) Choose a predicate Q in P , s.t. for all predicates W occurring in the body of some rule with head predicate Q , the formula $\varphi_{P,W}$ has already been defined.

Current Status



Aggregation-free SQL

Aggregation-free SQL is just a different format for expressing relational algebra queries:

Theorem

Relational algebra and aggregation-free SQL queries are equally expressive.

Learning objectives

Understanding:

- the notion of domain-independence,
- the active domain semantics,
- the notion of range restricted queries,
- Codd's Theorem: equivalence of various relational query languages in terms expressive power.

More details: Serge Abiteboul, Richard Hull, Victor Vianu: Foundations of Databases. Chapter 5. Addison-Wesley 1995,