

# Database Theory

VU 181.140, WS 2024

## 3. Codd's Theorem

Matthias Lanzinger  
(Based on slides of Reinhard Pichler)

Institute of Logic and Computation  
DBAI Group  
TU Wien

22 October, 2024



# 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

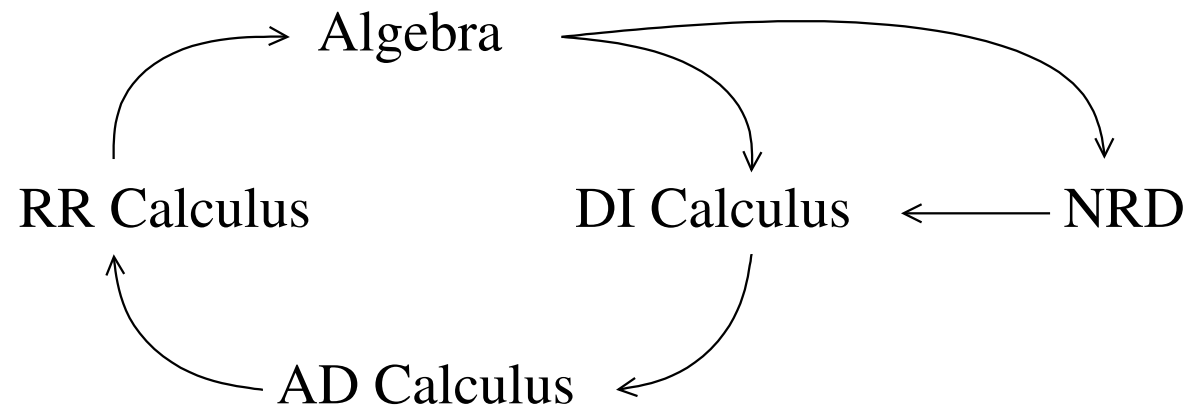
# What is Codd's Theorem?

- 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**.

# 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,

# 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  $Q_B(D)$  denote the result of evaluating query  $Q$  on database  $D$  assuming domain  $B$ .

# 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  $Q_B(D)$  denote the result of evaluating query  $Q$  on database  $D$  assuming domain  $B$ .

## Definition (domain-independence)

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

- a database instance  $D$  and
- two sets  $B, C$  that contain all constants that appear in  $D$  or in  $Q$  (also known as the **active domain**),

such that  $Q_B(D) \neq Q_C(D)$ .

# Queries Violating Domain Independence

## Example (Unsafe Queries)

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

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



# Undecidability of Domain Independence

## Definition

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

# Undecidability of Domain Independence

## Definition

For a domain-independent query  $Q$  and a database  $D$ , we may define  $Q(D) := Q_B(D)$  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

## Theorem

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

## Proof.

Almost by definition:

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

## Theorem

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

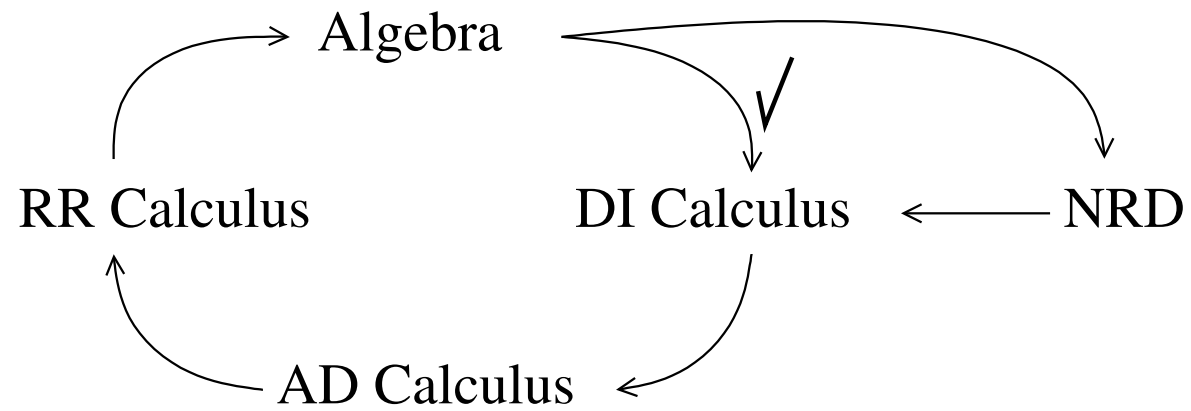
## Proof.

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. □

# Current Status



# Active Domain Interpretation

## Definition

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

# Active Domain Interpretation

## Definition

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

## Example

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

# Active Domain Interpretation

## Definition

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

## Example

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



## 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.

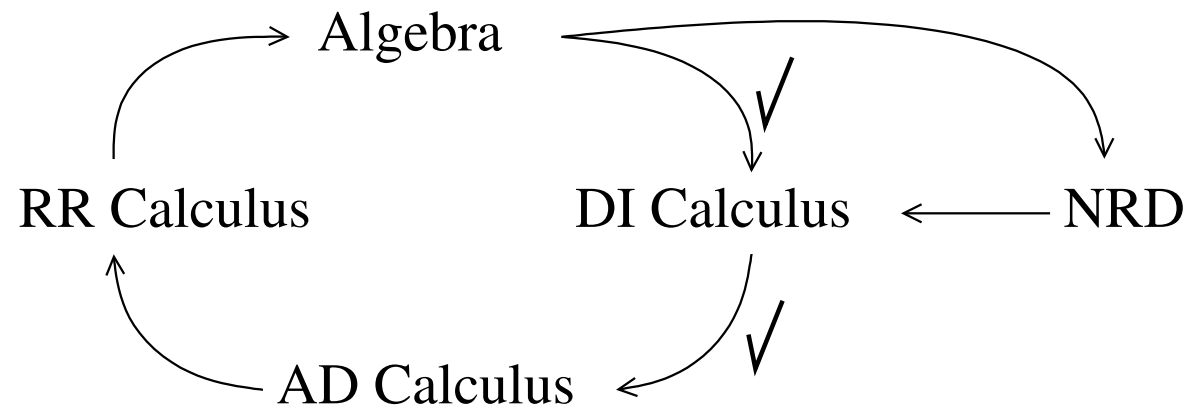
Consider an arbitrary database  $D$ .

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

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

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

# Current Status



## Range-restricted Queries: a sufficient condition for DI

**Preprocessing.** A formula is turned into **safe-range normal form (SRNF)** by the following steps:

- 1 Variable substitution: no distinct pair of quantifiers may employ the same variable and no variable may occur both bound and free.

**Example:**  $(\exists x \varphi(x)) \vee (\exists x \psi(x)) \vdash (\exists x \varphi(x)) \vee (\exists x' \psi(x'))$ .

- 2 Remove universal quantifiers:  $\forall x \varphi \vdash \neg \exists x \neg \varphi$ .
- 3 Remove implications:  $\varphi \Rightarrow \psi \vdash \neg \varphi \vee \psi$ .
- 4 Remove double negation:  $\neg \neg \varphi \vdash \varphi$ .
- 5 Flatten and/or, e.g.:  $(\varphi \wedge \psi) \wedge \pi \vdash \varphi \wedge \psi \wedge \pi$ .

# Range-restriction Check

## Definition

Given: Input formula  $\pi$  in SRNF.

$$\begin{aligned}
 rr(x = a) &:= \{x\} \\
 rr(R(t_1, \dots, t_n)) &:= \text{Vars}(\{t_1, \dots, t_n\}) \\
 rr(\varphi \wedge \psi) &:= rr(\varphi) \cup rr(\psi) \\
 rr(\varphi \vee \psi) &:= rr(\varphi) \cap rr(\psi) \\
 rr(\varphi \wedge x = y) &:= \begin{cases} rr(\varphi) & \dots \{x, y\} \cap rr(\varphi) = \emptyset \\ rr(\varphi) \cup \{x, y\} & \dots \text{otherwise} \end{cases} \\
 rr(\neg\varphi) &:= \emptyset \\
 rr(\exists x \varphi) &:= \begin{cases} rr(\varphi) - \{x\} & \dots x \in rr(\varphi) \\ \text{fail} & \dots \text{otherwise} \end{cases}
 \end{aligned}$$

If  $\text{free}(\pi) \subseteq rr(\pi)$ , then  $\pi$  is **range-restricted (RR)**.

# Range-restriction Examples

## Example (in SRNF)

$$\begin{array}{c}
 \overbrace{\hspace{15em}}^{rr(\cdot)=\{x,y\}} \\
 \overbrace{\hspace{10em}}^{rr(\cdot)=\{x,y,z\}} \\
 \overbrace{\hspace{10em}}^{rr(\cdot)=\{x,y,z\}} \\
 \overbrace{\hspace{10em}}^{rr(\cdot)=\{z\}} \\
 \overbrace{\hspace{10em}}^{rr(\cdot)=\{z\}} \\
 \overbrace{\hspace{10em}}^{rr(\cdot)=\{z\}} \quad \overbrace{\hspace{10em}}^{rr(\cdot)=\emptyset} \\
 \exists z : \overbrace{P(x,y,z)}^{rr(\cdot)=\{x,y,z\}} \vee (\overbrace{R(x,y)}^{rr(\cdot)=\{x,y\}} \wedge ((\overbrace{S(z)}^{rr(\cdot)=\{z\}} \wedge \overbrace{\neg T(x,z)}^{rr(\cdot)=\emptyset})) \vee T(y,z)))
 \end{array}$$

$rr(*) = \text{free}(*) = \{x, y\} \Rightarrow \text{range-restricted!}$

# Range-restricted Queries Capture AD Calculus

## Theorem

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

## Range-restricted Queries Capture AD Calculus

### Theorem

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

### Preparation of the proof.

#### Idea.

- define a predicate  $\alpha$  that captures exactly the active domain,
- use it to effectively restrict the quantification to the active domain (relativization)

For an  $n$ -ary predicate  $R$ , let  $dom(R, i)$  denote the formula  $\exists y_1, \dots, \exists y_{i-1}, \exists y_{i+1}, \dots, \exists y_n R(y_1, \dots, y_{i-1}, x, y_{i+1}, \dots, y_n)$ .

Let  $\alpha(x)$  be the disjunction of

- $dom(R, i)$  for all  $n$ -ary predicates of the schema and all  $1 \leq i \leq n$ , and
- $x = a$  for every constant of  $\varphi$ .

## Range-restricted Queries Capture AD Calculus (continued)

### Proof.

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

- 1 Turn  $\varphi$  into SRNF.
- 2 We build  $\varphi'$  from  $\varphi$  as follows:

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

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

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

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

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

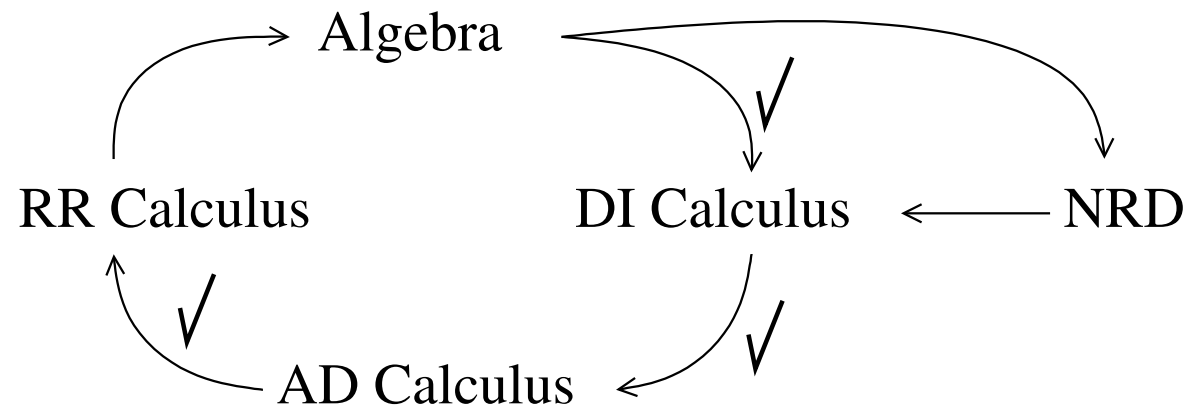
$$tr(\neg\varphi) := \neg tr(\varphi)$$

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





# Current Status



## Theorem

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

## Theorem

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

## Proof idea.

- 1 Start off with a query in SRNF.
- 2 Put the query into **Relational-Algebra NF (RANF)**:
  - RANF: Each subformula is **range-restricted**. (Exception: In a subformula  $\pi = \varphi_1 \wedge \dots \wedge \varphi_k \wedge \neg\psi$ ,  $\pi$  has to be RR and the  $\varphi_i$  and  $\psi$  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.

## 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  $\alpha$ :

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

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

- 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:

$$\begin{aligned} \text{Alg}(\varphi \wedge \psi) &:= \text{Alg}(\varphi) \bowtie \text{Alg}(\psi) \\ &(\text{Alg}(\varphi) \times \text{Alg}(\psi) \text{ if } \varphi \text{ and } \psi \text{ have disjoint schemas}) \end{aligned}$$

$$\begin{aligned} \text{Alg}(\varphi \vee \psi) &:= \text{Alg}(\varphi) \cup \text{Alg}(\psi) \\ \dots &\varphi \text{ and } \psi \text{ have the same schema} \end{aligned}$$

$$\begin{aligned} \text{Alg}(\varphi \wedge \neg\psi) &:= \text{Alg}(\varphi) - \text{Alg}(\psi) \\ \dots &\varphi \text{ and } \psi \text{ have the same schema} \end{aligned}$$

$$\text{Alg}(\exists y \varphi(\vec{x}, y)) := \pi_{\vec{x}} \text{Alg}(\varphi(\vec{x}, y))$$

$$\begin{aligned} \text{Alg}(\varphi \wedge x \vartheta t) &:= \sigma_{x \vartheta t} \text{Alg}(\varphi) \\ \dots &\vartheta \text{ is either } = \text{ or } \neq \text{ and } t \text{ is a term.} \end{aligned}$$

$$\begin{aligned} \text{Alg}(R(x_1, \dots, x_k)) &:= \rho_{A_1 \dots A_k \rightarrow x_1 \dots x_k} R \\ \dots &\text{relation } R \text{ has schema } R(A_1, \dots, A_k) \end{aligned}$$

## From RR Queries to the Algebra

### Example

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

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

corresponds to the RANF formula

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

An equivalent relational algebra expression looks as follows.

## From RR Queries to the Algebra

### Example

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

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

corresponds to the **RANF formula**

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

An **equivalent relational algebra expression** looks as follows.

$$\begin{aligned} & Alg(\exists x \exists y ((\alpha(x) \wedge \alpha(y)) \wedge \neg R(x, y))) \\ \vdash & \pi_{\emptyset} (Alg((\alpha(x) \wedge \alpha(y)) - Alg(R(x, y))) \\ \vdash & \pi_{\emptyset} ((\rho_{E \rightarrow x} \alpha \times \rho_{E \rightarrow y} \alpha) - \rho_{A, B \rightarrow x, y} R) \end{aligned}$$

## From RR Queries to the Algebra

### Example

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

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



## 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:

# 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 (in RANF)

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

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  $\alpha$  with schema  $\alpha(E)$ .

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

# From RR Queries to the Algebra

## Example

Schema  $R(AB)$ ,  $S(C)$ ,  $T(AC)$ ; active domain  $\alpha(E)$ .

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

# From RR Queries to the Algebra

## Example

Schema  $R(AB)$ ,  $S(AC)$ ,  $T(BC)$ ; active domain  $\alpha(E)$ .

$$\begin{aligned}
 & \{\langle x, y, z \rangle \mid R(x, y) \wedge (S(x, z) \vee T(y, z))\} \\
 \vdash & \{\langle x, y, z \rangle \mid \underbrace{(R(x, y) \wedge S(x, z))}_{R \bowtie S} \vee \underbrace{(R(x, y) \wedge T(y, z))}_{R \bowtie T}\} \\
 & \qquad \qquad \qquad \underbrace{\hspace{15em}}_{(R \bowtie S) \cup (R \bowtie T)} \\
 \text{or} & \{\langle x, y, z \rangle \mid R(x, y) \wedge (S(x, z) \vee T(y, z))\} \\
 \vdash & \{\langle x, y, z \rangle \mid \underbrace{R(x, y) \wedge ((S(x, z) \wedge \alpha(y)) \vee (T(y, z) \wedge \alpha(x)))}_{R \bowtie ((S \times \rho_B \alpha) \cup (T \times \rho_A \alpha))}\}
 \end{aligned}$$

## From RR Queries to the Algebra

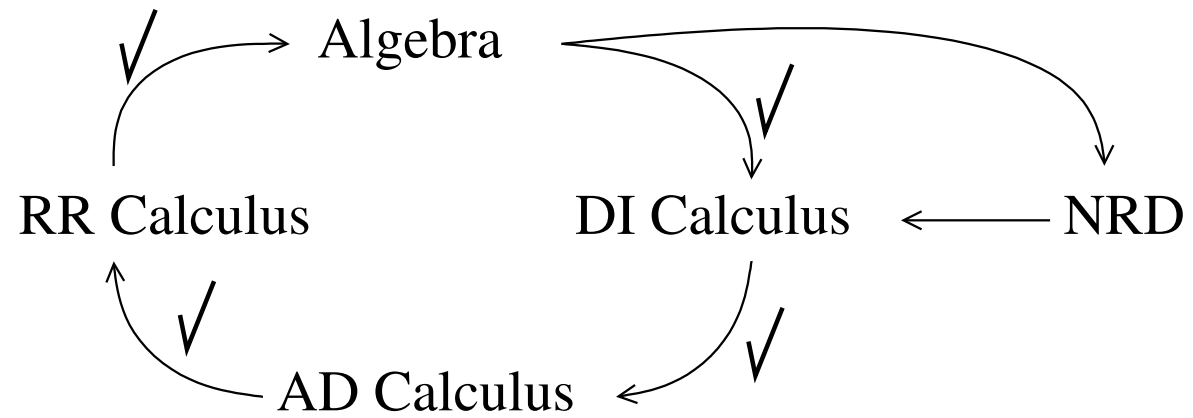
### Example

Schema  $R(AB)$ ,  $S(AC)$ ,  $T(BC)$ ; active domain  $\alpha(E)$ .

$$\begin{aligned}
 & \{\langle x, y, z \rangle \mid R(x, y) \wedge (S(x, z) \vee T(y, z))\} \\
 \vdash & \{\langle x, y, z \rangle \mid \underbrace{(R(x, y) \wedge S(x, z))}_{R \bowtie S} \vee \underbrace{(R(x, y) \wedge T(y, z))}_{R \bowtie T}\} \\
 & \qquad \qquad \qquad \underbrace{\hspace{15em}}_{(R \bowtie S) \cup (R \bowtie T)} \\
 \text{or} & \quad \{\langle x, y, z \rangle \mid R(x, y) \wedge (S(x, z) \vee T(y, z))\} \\
 \vdash & \quad \{\langle x, y, z \rangle \mid \underbrace{R(x, y) \wedge ((S(x, z) \wedge \alpha(y)) \vee (T(y, z) \wedge \alpha(x)))}_{R \bowtie ((S \times \rho_B \alpha) \cup (T \times \rho_A \alpha))}\}
 \end{aligned}$$

This is correct because  $R(x, y) \wedge (\varphi \vee \psi) \equiv R(x, y) \wedge \alpha(x) \wedge \alpha(y) \wedge (\varphi \vee \psi) \equiv R(x, y) \wedge ((\alpha(x) \wedge \alpha(y) \wedge \varphi) \vee (\alpha(x) \wedge \alpha(y) \wedge \psi))$ .

## Current Status





## Non-recursive Datalog with negation

### Definition

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

- non-recursive means that negation is trivially stratified!

## Non-recursive Datalog with negation

### Definition

Non-recursive Datalog with negation ( $nr\text{-Datalog}^-$ ) prohibits cycles in the dependency graph 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.*

## Non-recursive Datalog with negation

### Definition

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

- non-recursiveness means that negation is trivially stratified!

### Theorem

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

# From Algebra to nr-*Datalog*<sup>⊃</sup>

Proof.

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

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

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

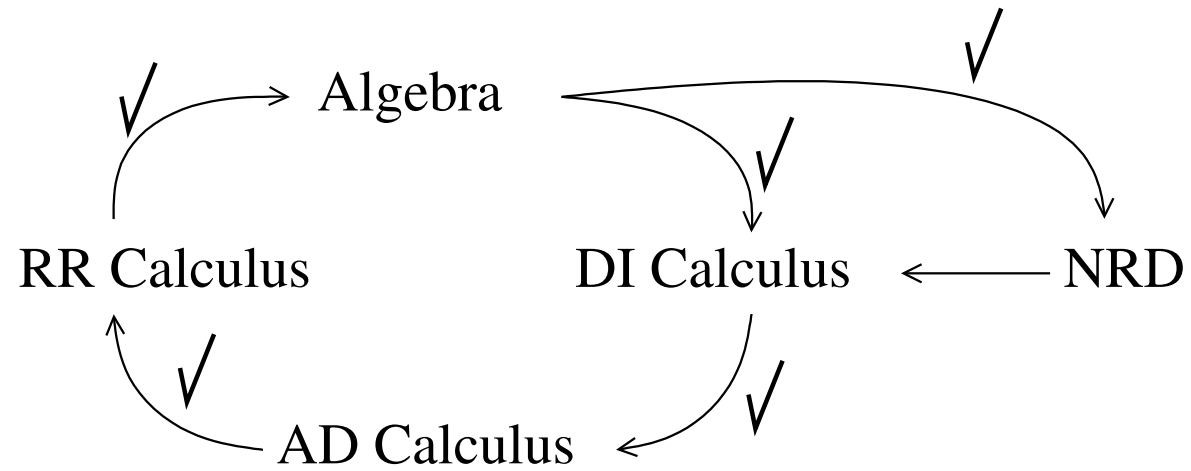
$$P_{E_1 \times E_2} := \{Q_{E_1 \times E_2}(\vec{x}, \vec{y}) :- Q_{E_1}(\vec{x}), Q_{E_2}(\vec{y})\} \cup P_{E_1} \cup P_{E_2}$$

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

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

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  $P_E^* \wedge DB^* \models Q_E(\vec{c})$ . □

# Current Status



## From nr-*Datalog*<sup>¬</sup> to DI Calculus

### Theorem

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

## From nr-*Datalog*<sup>¬</sup> 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*<sup>¬</sup> 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  $P^* \wedge DB^* \models R(\vec{c})$ .

## From nr-*Datalog*<sup>¬</sup> 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)



## From nr-*Datalog*<sup>¬</sup> 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.

## Proof.

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

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

Then we set

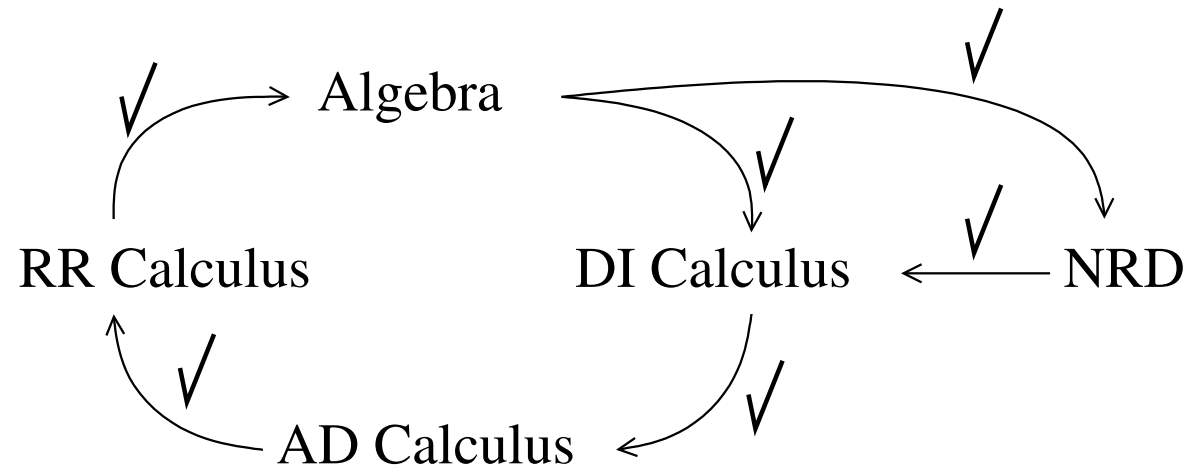
$$\varphi_{P,Q} = \bigvee_{i \in \{0, \dots, m\}} (\exists \vec{y}_i : \bigwedge_{j \in \{0, \dots, n_i\}} L'_{i,j}),$$

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

$$L'_{i,j} = \begin{cases} \varphi_{P,W}(v) & \text{if } L_{i,j} \text{ is an atom } W(v). \\ \neg \varphi_{P,W}(v) & \text{if } L_{i,j} \text{ is an atom } \neg W(v). \end{cases}$$

□

# Current Status



# 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,