

# The Semijoin Algebra

Jan Van den Bussche  
Hasselt University

joint work with Dirk Leinders, Maarten Marx, Jerzy Tyszkiewicz

# The relational algebra, RA

Projection  $\pi_{A,B,C}$

- allow repetitions:  $\pi_{A,A/B}$

Selection  $\sigma_{A=B}, \sigma_{A<B}$

Renaming  $\rho_R$

Union, intersection, difference

Equijoin  $R \bowtie_{\substack{R.A=S.B \\ R.C=S.D}} S$

- special cases: natural join, cartesian product

## RA expressions

Build up expressions for complex queries

Likes(drinker, beer), Serves(bar, beer), Visits(drinker, bar)

Losers:

$$\pi_{V.d}(V) - \pi_{V.d}\sigma_{V.d=L.d}(V \bowtie S \bowtie L)$$

Codd's theorem: RA is equivalent to first-order logic (relational calculus)

## The semijoin algebra, SA

Equi-semijoin:

$$\begin{aligned} R \bowtie_{\theta} S &:= \{t \in R \mid \exists s \in S : \theta(t, s) \text{ is true}\} \\ &= \pi_R(R \bowtie_{\theta} S) \end{aligned}$$

with  $\theta$  a conjunction of equalities

SA is RA where we replace  $\bowtie$  by  $\bowtie$

Visitors of lousy bars:

$$V \bowtie (\pi_{bar}(S) - \pi_{bar}(S \bowtie L))$$

## The guarded fragment of first-order logic, GF

[Andréka, van Benthem, Németi]

Quantifiers are restricted to the following form:

$$\begin{aligned} & \exists \bar{y} (\alpha(\bar{x}, \bar{y}) \wedge \psi(\bar{x}, \bar{y})) \\ & \forall \bar{y} (\alpha(\bar{x}, \bar{y}) \rightarrow \psi(\bar{x}, \bar{y})) \end{aligned}$$

- $\alpha$  atomic formula (single relation)
- all free variables of  $\psi$  must occur in  $\alpha$

Visitors of lousy bars:

$$\{d, ba \mid V(d, ba) \wedge \neg \exists be (S(ba, be) \wedge \exists d L(d, be))\}$$

Originally introduced in the context of modal, algebraic logic

## Codd theorem for SA

SA is equivalent to GF

From SA to GF:

$$\begin{aligned} & V \times (\pi_{bar}(S) - \pi_{bar}(S \times L)) \\ & V \times (\pi_{bar}(S) - \pi_{bar}\{ba, be \mid S(ba, be) \wedge \exists d L(d, be)\}) \\ & V \times (\pi_{bar}(S) - \{ba \mid \exists be(S(ba, be) \wedge \exists d L(d, be))\}) \\ & V \times (\{ba \mid \exists be S(ba, be)\} - \{ba \mid \exists be(S(ba, be) \wedge \exists d L(d, be))\}) \\ & V \times (\{ba \mid \exists be S(ba, be) \wedge \neg \exists be(S(ba, be) \wedge \exists d L(d, be))\}) \\ & \{d, ba \mid V(d, ba) \wedge \exists be S(ba, be) \wedge \neg \exists be(S(ba, be) \wedge \exists d L(d, be))\}) \end{aligned}$$

From GF to SA:

$$\begin{aligned} & \{d, ba \mid V(d, ba) \wedge \neg \exists be (S(ba, be) \wedge \exists d L(d, be))\} \\ & \{d, ba \mid V(d, ba) \wedge \neg \exists be (S(ba, be) \wedge be \in \pi_{be}(L))\} \\ & \{d, ba \mid V(d, ba) \wedge \neg \exists be ((ba, be) \in S \times \pi_{be}(L))\} \\ & \{d, ba \mid V(d, ba) \wedge \neg (ba \in \pi_{ba}(S \times \pi_{be}(L)))\} \\ & \{d, ba \mid V(d, ba) \wedge ba \in (\pi_{ba}(S) - \pi_{ba}(S \times \pi_{be}(L)))\} \\ & V \times (\pi_{ba}(S) - \pi_{ba}(S \times \pi_{be}(L))) \end{aligned}$$

## Consequences of $SA = GF$

Equivalence extends to fixpoint logic:  $\mu SA = \mu GF$

**Ex:** Database relations  $R(A, B)$  and  $T(B)$ ,  
relation variable  $X(A, B)$ :

$$\text{LFP } X. (R \times T) \cup (R \underset{R.B=X.A}{\times} X)$$

SA has the finite model property

Our translation  $SA \rightarrow GF$  is exponential; still:

- Satisfiability of SA-expressions is decidable  
(complete for EXPTIME)

Polynomial translation  $SA \rightarrow GF$ ?



## Guarded bisimilarity

GF is invariant under *guarded bisimilarity*,  $\simeq_g$

Databases  $A$  and  $B$ , same schema, tuple  $\bar{a}$  in  $A$ , tuple  $\bar{b}$  in  $B$

**Def.**  $(A, \bar{a}) \simeq_g (B, \bar{b})$  if player II can keep up forever in the following game:

1. initial game position is  $(A, \bar{a})$  and  $(B, \bar{b})$
2. player I chooses a tuple in one of the databases, say  $A$
3. player II responds in other database  $\Rightarrow (A, \bar{a}')$  and  $(B, \bar{b}')$

- $\bar{a}'$  and  $\bar{b}'$  must satisfy precisely same relations, predicates
  - if  $\bar{a}$  and  $\bar{a}'$  agree in  $i$ th position, then  $\bar{b}$  and  $\bar{b}'$  must too
4. if player II cannot respond correctly he loses;  
otherwise repeat from new position  $(A, \bar{a}')$  and  $(B, \bar{b}')$ .

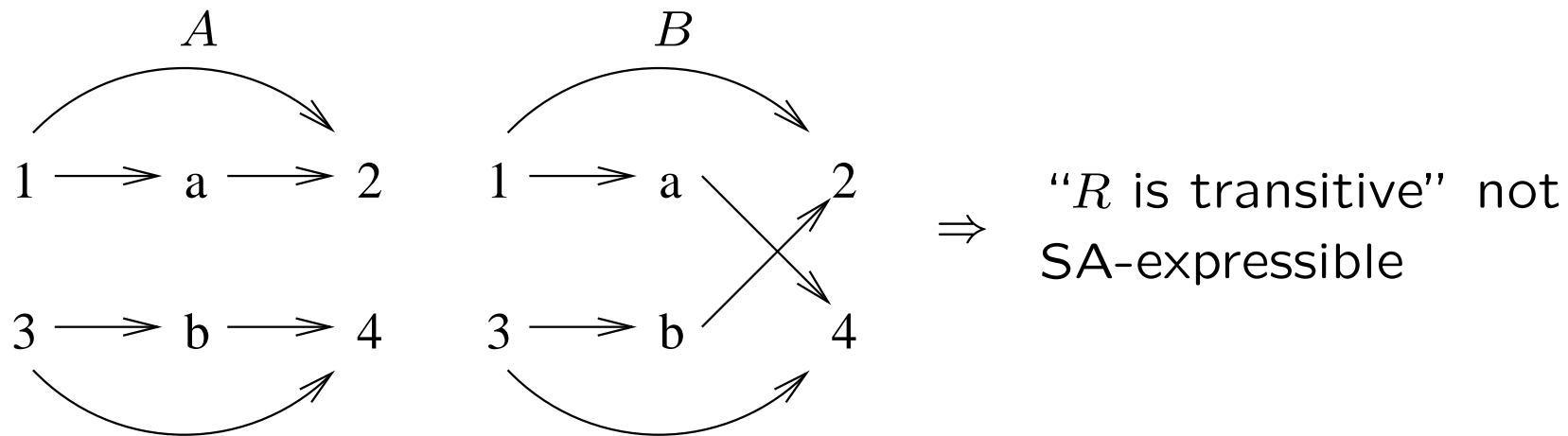
## Invariance property

If  $(A, \bar{a}) \simeq_g (B, \bar{b})$  then for all SA-expressions  $E$ :

$$\bar{a} \in E(A) \Leftrightarrow \bar{b} \in E(B)$$

Use to prove SA-inexpressibility of certain queries

**Ex.** single relation  $R$ :



## Division

$$R(A, B) \div S(C) := \{a \mid \forall b \in S : (a, b) \in R\}$$

RA-expressible, but not SA:

<i>A</i>			<i>B</i>	
<i>R</i>	<i>S</i>		<i>R</i>	<i>S</i>
1 7	7		1 7	7
1 8	8		1 8	8
2 7			2 8	9
2 8			2 9	
			3 7	
			3 9	

## Linear query processing

*Linear* RA expression: on every database, every intermediate result has linear size

**linear:**  $(\sigma R \cup \pi S) - T$

**not linear:**  $R \cap (S \bowtie T)$

**linear:**  $R \bowtie_{R.A=S.B} \pi_B(S) = R \bowtie_{R.A=S.B} S$

Every query expressible by a linear RA expression is already expressible by a SA expression

Note that SA-expressions are always linear

## Proof idea

For  $E_1 \bowtie_{\theta} E_2$  to be linear, every joining tuple pair  $(\bar{a}, \bar{b})$  must satisfy  $\forall i \exists j : a_i = b_j$  or vice versa

- if not, we could “blow up” the database by duplicating the “free” values in  $\bar{a}$  and  $\bar{b}$
- blown up database is guarded bisimilar
- since  $E_1$  and  $E_2$  can be assumed linear by induction, they will output the duplicate tuples  $\Rightarrow$  quadratic join size

Such joins can be expressed in SA

## Corollaries

Every RA-expression is either linear, or has a subexpression that has quadratic output size

Every RA-expression either produces quadratic intermediate results, or is equivalent to an SA-expression

## Set joins

We now know that division is not expressible in linear RA

Division is a restricted kind of *set join*

**Def.** Let  $P(X, Y)$  be a predicate about sets.  
For relations  $R(A, B)$  and  $S(C, D)$ :

$$R \bowtie_P^{\text{set}} S := \{a, c \mid P(\{b : R(a, b)\}, \{d : S(c, d)\})\}$$

subset join:  $\bowtie_{X \subseteq Y}^{\text{set}}$

set-equality join:  $\bowtie_{X=Y}^{\text{set}}$

standard equijoin:  $\bowtie_{X \cap Y \neq \emptyset}^{\text{set}}$

If the emptiness query for  $\bowtie_P^{\text{set}}$  can be expressed in linear RA, then  $P$  must be monotone



## Sidenote: Grouping and aggregation

It is well known that division can be linearly expressed using counting:

$$R(A, B) \div S(C) = \pi_A(\gamma_{A, \text{count}(B)}(R \bowtie_{B=C} S) \bowtie_{\text{count}(B)=\text{count}(C)} \gamma_{\text{count}(C)}(S))$$

## Theta-equi joins

$\bowtie_{\theta}$  with  $\theta$  more than just conjunction of equalities?

**Ex.** “ $R(A, B) \models A \rightarrow B$ ”:

$$R \underset{\substack{R.A=S.A \\ R.B \neq S.B}}{\bowtie} \rho_S(R)$$

Satisfiability of  $SA^{\neq}$  is undecidable

Do our linearity results extend to  $SA^{\neq}$ ? and to  $SA^{<}$ ?

## $\Omega$ -guarded bisimilarity

$\Omega$ , signature of predicates that can be used in  $\theta$

$SA^=$ : if  $\bar{a}$  and  $\bar{a}'$  agree in  $i$ th position, then  $\bar{b}$  and  $\bar{b}'$  must too

$SA^\Omega$ :  $(\bar{a}, \bar{a}')$  and  $(\bar{b}, \bar{b}')$  must satisfy precisely the same predicates from  $\Omega$

## Conclusion

SA = GF but  $SA^\Omega$  is more powerful

SA = linear RA

Division, set joins not linear RA

- theoretical explanation why these queries are hard on the query processor

Many open problems remain