



The expressive power of SQL

Submitted by :

Doris Fischer 011143252

Ayelet Kaidar 031756471

SQL vs. RA

- ◆ Aggregate functions - such as COUNT, AVG, MIN, MAX etc.
- ◆ Grouping – not only can one compute aggregates, one can also group them by values of different attributes.
- ◆ Arithmetic – SQL allows one to apply arithmetic operations to numerical values.

A simple example

- ◆ $R(\text{Src}, \text{Dest})$ is a relation with flight information, where Src stands for source and Dest stands for destination.
- ◆ To find pairs of cities (A, B) such that it's possible to fly from A to B with one stop, one would use a self-join:

A simple example (cont.)

- ◆ SQL:
SELECT R1.Src,R2.Dest
FROM R AS R1, R AS R2
WHERE R1.Dest = R2.Src
- ◆ RA:
 $\pi_{\{R1.Src,R2.Dest\}}(\sigma_{\{R1.Dest=R2.Src\}}(R1 \Join R2))$

A simple example (cont.)

- ◆ Suppose we want pairs of cities such that one makes 2 stops on the way?
- ◆ `SELECT R1.Src , R3.Dest`
`FROM R AS R1, R AS R2 , R AS R3`
`WHERE R1.Dest = R2.Src AND R2.Dest = R3.Src`

A simple example (cont.)

- ◆ RA:
$$\pi_{\{R1.Src, R3.Dest\}}(\sigma_{\{R1.Dest=R2.Src \wedge R2.Dest=R3.Src\}}(R1 \Join R2 \Join R3))$$

Transitive closure

- ◆ But often one needs a general reachability query in which no a priori bound on the number of stops is known; that is, whether it is possible to get to B from A.
- ◆ Graph-theoretically, This means computing the transitive closure of R.
- ◆ It's well known that the transitive closure of a graph is not expressible in relational algebra (RA) or calculus (RC)

Transitive closure (Cont.)

- ◆ In SQL3, one could write the reachability query as:
- ◆ WITH RECURSIVE TrCL(Src, Dest) AS

R

UNION

SELECT TrCL.Src , R.Dest

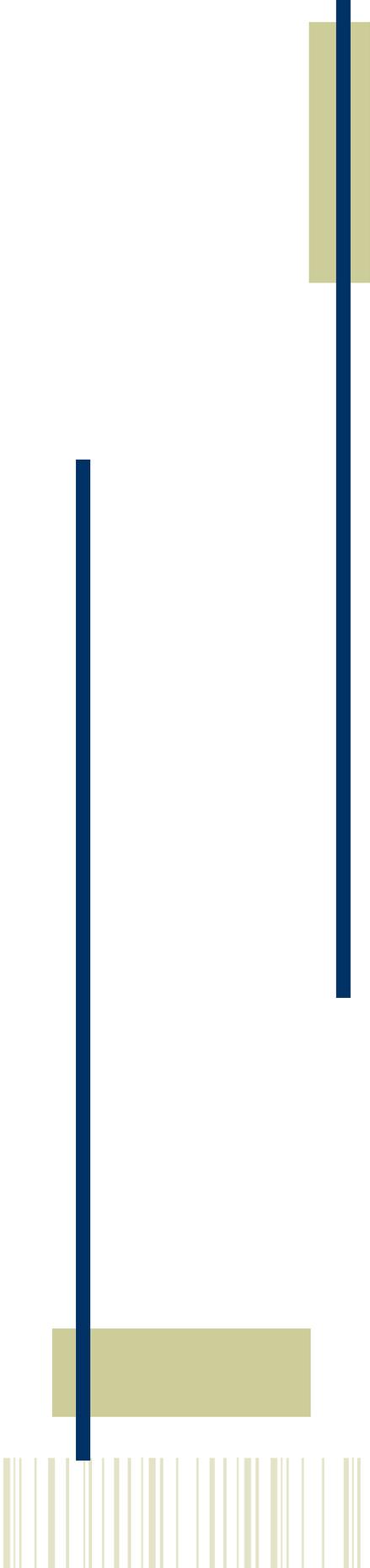
FROM TrCL, R

WHERE TrCL.Dest = R.Src

SELECT * FROM TrCL

Transitive closure (Cont.)

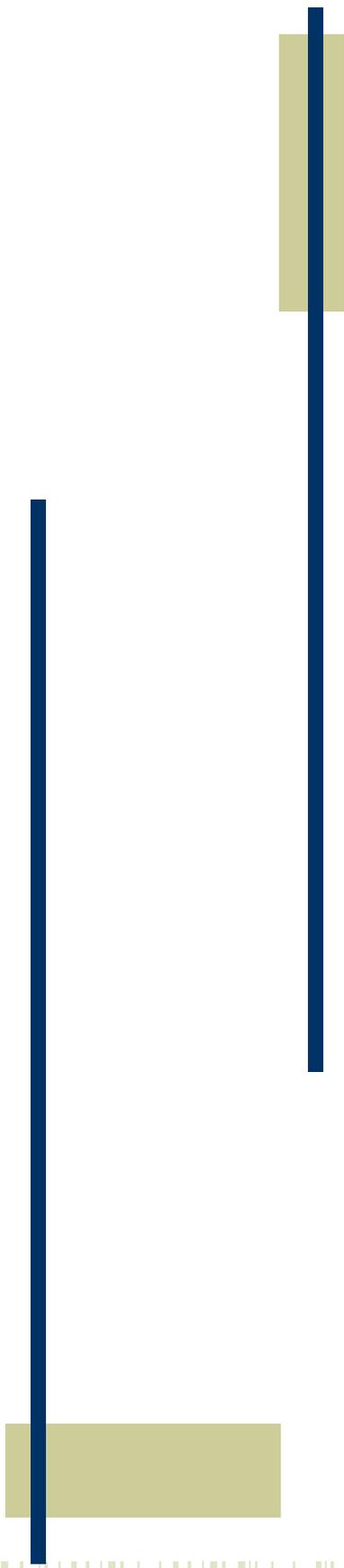
- ◆ This simply models the usual datalog rules for transitive closure :
- ◆ $\text{trcl}(x,y) :- r(x,y)$
- ◆ $\text{trcl}(x,y) :- \text{trcl}(x,z),r(z,y)$

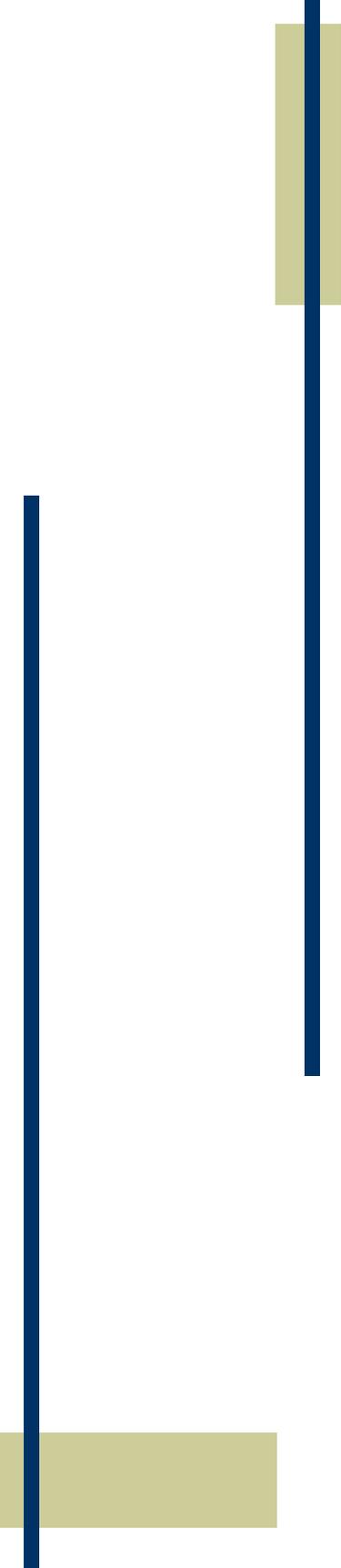
- 
- ◆ The reason for introducing recursion in this SQL standard is precisely this folk result stating that it cannot be expressed in the language.
 - ◆ But when one looks at what evidence is provided to support this claim, one notices that all the references point to papers , in which it proved that RA and RC cannot express recursive queries.
 - ◆ Why is this not sufficient?

Consider this query

- ◆

```
SELECT R1.A  
FROM R1, R2  
WHERE(SELECT COUNT(*) FROM R1) >  
(SELECT COUNT(*) FROM R2)
```
- ◆ This query tests if $|R1| > |R2|$: in that case it returns the A attributes of R1, otherwise it returns the empty set.

- 
- ◆ However, logicians proved, long ago, that first-order logic, and thus RC, cannot compare cardinalities of relations, and yet we have a very simple SQL query doing precisely that.
 - ◆ The conclusion, of course, is that SQL has more power than RA, and the main source of this additional power is its *aggregation* and *grouping* constructs, together with *arithmetic* operations on numerical attributes.

- 
- 
- ◆ But than one cannot say that the transitive closure query is not expressible in SQL simply because it is inexpressible in RC.

RA with aggregates

- ◆ In order to deal with aggregation we need to distinguish numerical and non-numerical columns.
 - ◆ We define 2 types:
 - b - non numerical type with domain Dom.
 - n - numerical type with domain Num.
- Num \belongs N,Z,Q,R etc.

RA with aggregates (Cont.)

- ◆ The type of a relation is a string over $\{b,n\}$.
For example:
the type of $S(\text{Empl}, \text{Salary})$ is bn .
- ◆ For a type t , $t.i$ is the i th position in the string. The length of t is denoted by $|t|$.
- ◆ A database scheme is a collection of relation names R_i and their types t_i .
- ◆ We write $R_i:t_i$ if the type of R_i is t_i .



RA with aggregates (Cont.)



- ◆ $\backslash\Omega$ – a collection of functions and predicates on Num.
- ◆ $\backslash\tau$ – a collection of aggregates.



Locality of queries



- ◆ A query can only look at a small portion of its input.
- ◆ If the input is a graph, “small” means a neighborhood of a small radius.

Example 1 - Reachability

- ◆ Relation R contains edges of a directed graph.
- ◆ $R:bb$ is a *purely relational* schema (contains only non-numerical types).
- ◆ Suppose e is an expression of type bb .
- ◆ Given a pair (a,b) belongs R , and a number $r > 0$:
 $N_r^R(a,b) = r$ -neighborhood of a,b in R , is the subgraph on the set of nodes in R whose distance from either a or b is at most r .

Example 1 – Reachability

(cont.)

- ◆ $N_r^R(a,b) = \{ x \mid |x-a| \leq r \text{ or } |x-b| \leq r \}$
- ◆ (a,b) isomorphic r^R (c,d) meaning the 2 neighborhoods $N_r^R(a,b)$ and $N_r^R(c,d)$ are isomorphic,
- ◆ Reachability is not local.



Example2 - RC



- ◆ RC expressions are all local.
- ◆ See proof in class lectures – winning series for the second player.

Locality of queries

- ◆ e is local if there is a number r , depending on e only, s.t.:

$(a,b) \setminus \text{isomorphic}_{r^R}(c,d) \setminus \mathcal{M}(a,b) \setminus \text{belongs to } e(R) \text{ iff } (c,d) \setminus \text{belongs to } e(R)$



ALG_agg(All, All)

- ◆ Relational algebra with grouping and aggregate functions.
- ◆ When we consider only pure relational scheme the only aggregate function will be **COUNT**.



Theorem 1 – Locality of SQL

Let e be a pure relational graph query in
 $\text{ALG_aggr}(\text{All}, \text{All})$, that is, an expression of
type bb over the scheme of one symbol $R:\text{bb}$.
Then e is local.



Conclusion so far...

- ◆ Reachability is not expressible in SQL over the base type of b , no matter what aggregate functions and arithmetic operations are available.
- ◆ Inexpressibility of many other queries can be derived from this, for example, tests for graph connectivity and acyclicity.



Question

Question: Are expressions in SQL(pure relational)-RA local?

Answer: YES.

Next steps...

1. Introduce an *aggregate logic* $\mathcal{L}_{\text{aggr}}$, as an extension of FOL and show how ALG_aggr queries are translated into it.
2. Replace aggregate terms by *counting quantifiers*, thereby translating $\mathcal{L}_{\text{aggr}}$ into a simpler logic \mathcal{L}_{C} .
3. Inductive proof of locality of \mathcal{L}_{C} formulae.

The logic $\mathcal{L}_{\text{aggr}}$

- ◆ The structures are relational databases over 2 base types with domains Dom and Num; that is vocabularies are just schemas.
- ◆ The logic is *two-sorted*.
- ◆ Dom will be referred as *first-sort* and Num as *second-sort*.

The logic $\mathcal{L}_{\text{aggr}}(\omega, \text{teta})$

- ◆ A SC structure D is a tuple $\langle A, R_1^D, \dots, \rangle$, where A is a finite subset of Dom , and R_i^D is a finite subset of $\prod_{j=1}^{|t_j|} \text{dom}_j(D)$

where $\text{dom}_j(D) = A$ for $t_{i,j} = b$,
and $\text{dom}_j(D) = \text{Num}$ for $t_{i,j} = n$.

The logic $\mathcal{L}_{aggr}(\omega, \tau)$

cont.

- ◆ A variable of sort i is a term of sort i , $i=1,2,\dots$
- ◆ If $R:t$ is in SC, and u is a tuple of terms of type t , then $R(u)$ is a formula.
- ◆ Formulae are closed under Boolean connectives and quantification (respecting sorts). If x is a first-sort variable, $\exists x$ is interpreted as $\exists x$ x belongs to A ; if k is a second-sort variable $\exists k$ is interpreted as $\exists k$ k belongs to Num .

The logic $\mathcal{L}_{\text{aggr}}(\omega, \tau)$ cont.

- ◆ If P is an n -ary predicate in ω and $\tau_1 \dots \tau_n$ are second-sort terms, then $P(\tau_1 \dots \tau_n)$ is a formula.
- ◆ If f is an n -ary function in ω and $\tau_1 \dots \tau_n$ are second-sort terms, then $f(\tau_1 \dots \tau_n)$ is a second-sort term.

The logic $\mathcal{L}_{\text{aggr}}(\omega, \tau)$

cont.

- ◆ If \mathcal{F} is an aggregate in τ ,
 $\phi(\bar{x}, \bar{y})$ is a formula and
 $\tau(\bar{x}, \bar{y})$ a second-sort term,
then $\tau'(x)$
 $= \text{Aggr}_{\mathcal{F}}(\bar{y}.(\phi(\bar{x}, \bar{y})))$
is a second-sort term
with free variables x .

Theorem 2

Let $e:t$ be an expression of $\text{ALG}_{\text{aggr}}(\omega, \text{teta})$. Then there is a formula $\phi_e(x)$ of $\mathcal{L}_{\text{aggr}}(\omega, \text{teta})$, with x of type t , such that for any SC-database D ,

$$e(D) = \{a \mid D \models \phi_e(a)\}$$

Theorem 2 – proof

- ◆ Numerical selection:
$$e' = \sigma_{\{i_1 \dots i_k\}}(e), P \in \omega,$$
$$\phi_{\{e'\}}(\bar{x}) \text{ defined } \phi_{\{e\}}(\bar{x}) \text{ and } P(x_{i_1} \dots x_{i_k}).$$
- ◆ Function application:
$$e' = \text{Apply}[f]_{\{i_1 \dots i_k\}}(e), f: \text{Num}^k \rightarrow \text{Num}$$
$$\in \omega,$$
$$\phi_{\{e'\}}(\bar{x}, q) \text{ defined } \phi_{\{e\}}(\bar{x})$$
$$\text{ and } (q = f(x_{i_1} \dots x_{i_k})).$$

Theorem 2 – proof (cont.)

- ◆ Aggregation:
 $e' = \text{Aggr}[i:\mathcal{F}](e)$.
 $\phi_{\{e'\}}(\bar{x}, q)$ defined
 $\phi_{\{e\}}(\bar{x})$ and
 $(q = \text{Aggr}_F \bar{y}).(\phi_{\{e\}}(\bar{y}), y_i)$.

Theorem 2 – proof (cont.)

- ◆ Grouping:
 $e' = \text{Group}_m[\lambda S.e_1](e_2)$, where $e_1:u$ is over $SC \wedge \text{conjunction } \{S\}$, and $e_2:ts$ is over SC .
Let $x:t, y:s, z:u$. Then:
 $\text{phi}_{\{e'\}}(\bar{x}, \bar{z}) \wedge \text{defined } \exists \bar{y}$
 $\text{phi}_{\{e_2\}}(\bar{x}, \bar{y}) \wedge$
 $\text{phi}_{\{e_1\}}(\bar{z})[\text{phi}_{\{e_2\}}(\bar{x}, \bar{y})/S(\bar{v})]$

Theorem 3 (general)

Let e be a pure relational query in $\text{ALG_aggr}(\text{All}, \text{All})$, that is, an expression of type $b \dots b$ over a pure relational schema. Then e is local.

Translating L_{aggr} into L_c

- ◆ The logic L_c is simpler than L_{aggr} because it does not have aggregate terms.
- ◆ The price for the above will be that L_c has infinitary conjunctions and disjunctions.
- ◆ However, the translation ensures that for each infinite conjunction or disjunction, there is a uniform bound on the rank (Will be discussed later...) of formulae in it, and this property suffices to establish locality.

The logic L_c

- ◆ The structures are the same as for L_{aggr} .
- ◆ The only terms are variables (of any sort).
- ◆ Any constant c belongs Num is a term of the second sort.
- ◆ Atomic formulae are $R(x)$, where R belongs SC, and x is a tuple of terms of the appropriate sort, and x equals y , where x, y are terms of the same sort.
- ◆ Formulae are closed under Boolean connectives.

The logic L_C – cont.

- ◆ Formulae are infinitary connectives:
 - if $\phi_i, i \in I$, is a collection of formulae, then
 - $\bigvee_{i \in I} \phi_i$ and $\bigwedge_{i \in I} \phi_i$ are L_C formulae.
- ◆ They are also closed under first and second sort quantification.

The logic L_C – cont.

- ◆ For every i belongs N , there is a quantifier $\exists i$ that binds one first sort variable: that is, if $\phi(x, \bar{y})$ is a formula, then $\exists x \phi(x, \bar{y})$ is a formula whose free variables are \bar{y} .

The logic L_C – cont.

- ◆ The semantic will be:
 - $D \models \exists x \bar{\phi}(x, \bar{a})$ if there are i distinct elements $b_1 \dots b_i$ belongs A s.t.
 - $D \models \exists x \bar{\phi}(b_i, \bar{a})$, $1 \leq j \leq i$.
- ◆ That is, the existential quantifier is witnessed by at least i elements.
- ◆ Note: $\exists x \bar{\phi}$ is equivalent to $\exists x \bar{\phi}$.

The rank of a formula – L_C

- ◆ For L_C this is the quantifier rank, but the second sort quantification does not count.
- ◆ For each atomic ϕ , $\text{rk}(\phi) = 0$.
- ◆ For $\phi = \text{or}_i \psi$, $\text{rk}(\phi) = \sup_i \text{rk}(\psi)$, and likewise for and .
- ◆ $\text{rk}(\text{not } \phi) = \text{rk}(\phi)$.
- ◆ $\text{rk}(\text{exists } x \phi) = \text{rk}(\phi) + 1$ for x first sort;
- ◆ $\text{rk}(\text{exists } k \phi) = \text{rk}(\phi)$ for k second sort.

The rank of a formula – L_aggr

- ◆ For L_aggr the definition differs slightly.
- ◆ For a variable or a constant term, the rank is 0.
- ◆ The rank of an atomic formula is the maximum rank of a term in it.
- ◆ $\text{rk}(\phi_1 * \phi_2) = \max(\text{rk}(\phi_1), \text{rk}(\phi_2))$,
for $*$ belongs { \vee , \wedge }
- ◆ $\text{rk}(\neg \phi) = \text{rk}(\phi)$.

The rank of a formula – L_aggr

– cont.

- ◆ $\text{rk}(f(\tau_1 \dots \tau_n)) = \max_{\{1 \leq i \leq n\}} \text{rk}(\tau_i)$.
- ◆ $\text{rk}(\exists x \phi) = \text{rk}(\phi) + 1$ if x is first sort;
 $\text{rk}(\exists k \phi) = \text{rk}(\phi)$ if k is second sort.
- ◆ $\text{rk}(\text{Aggr}_{\mathcal{F}} \bar{y} . (\phi, \tau)) = \max(\text{rk}(\phi), \text{rk}(\tau)) + m$, where m is the number of first sort variables in y .

Translating L_{aggr} into L_C

Proposition 1: For every formulae $\phi(X)$ of $L_{\text{aggr}}(\text{All}, \text{All})$, there exists an equivalent formulae $\phi^0(X)$ of L_C s.t.

$$\text{rk}(\phi^0) \leq \text{rk}(\phi)$$

L_C is local

- ◆ Formulae of L_aggr have finite rank; hence they are translated into L_C formulae of finite rank.
- ◆ It can be shown by a simple induction that those formulae are local.
- ◆ The proof is based on Lemma1 – the Permutation Lemma.

The Permutation Lemma

Let D be first-sort, with $A = \text{adom}(D)$, and $r > 0$.

If a $\{3r+1\}$ -ary D exists, then there exists a permutation

$p: A \rightarrow A$ s.t. $\forall c \in D$ $p(c)$ belongs to A .

Putting everything together

- ◆ Let e be a pure relational expression of ALG_{agg}(ALL,ALL).
- ◆ By Theorem 2, it is expressible in $L_{agg}(ALL,ALL)$, and by Proposition 1, by a L_C formula of finite rank.
- ◆ Hence it is local.

SQL over ordered domains

- ◆ So far the only nonnumerical selection was of the form $\sigma_{\{i=j\}}$, testing equality of 2 attributes. We now extend the language to $ALG^{<}_{aggr}$ by allowing selections of the form $\sigma_{\{i<j\}}(e)$, where both i and j are of type b , and $<$ is some fixed linear ordering on the domain Dom .

Natural Numbers

- ◆ Let $\text{num} = \mathbb{N}$. We consider a version of ALG_aggr that has the most usual set of arithmetic and aggregate operators: namely, $+$, $<$ and constants for arithmetic, and the aggregate Σ . This suffices to express aggregates MIN, MAX, COUNT, TOTAL, but certainly not AVG, which, produces rational numbers.

Natural Numbers (Cont.)

- ◆ We shall use the notations:
 - SQL_N for $ALG_aggr(\{+, \cdot, <, 0, 1\}, \{\Sigma\})$,
and
 - $SQL_{N^<}$ for $ALG_aggr^<(\{+, \cdot, <, 0, 1\}, \{\Sigma\})$.
- ◆ We show how a well-known counting logic $FO(C)$ can be embedded into $SQL_{N^<}$.

Definition of FO(C)

- ◆ It is a two-sorted logic, with second sort being the sort of natural numbers. That is, a structure D is of the form $(\{a_1, \dots, a_n\}, (1, \dots, n), <, +, \dots, 1, n, R_1, \dots, R_t)$



Theorem 4.

Over ordered structures,
 $\text{FO}(C) \subseteq \text{SQL}_N^<$.

Reachability over ordered domains

- ◆ Notice that the reachability query, even over ordered domains of nodes, is order-independent; that is, the result does not depend on a particular ordering on the nodes, just on the graph structure.
- ◆ Order-independent queries in SQL_N and SQL_N^< are the same.



Proposition 3



- ◆ There exist order-independent non-local queries expressible in $\text{SQL}_N^{<}$. Thus, there are order-independent $\text{SQL}_N^{<}$ queries not expressible in SQL_N .



Conclusion



- ◆ Did SQL3 designers really have to introduce recursion, or is it expressible with what's already there?
- ◆ Our results show that they clearly had a good reason for adding a new construct , because:

Conclusion (Cont.)

1. Over unordered types, reachability queries cannot be expressed by the basic SQL *SELECT-FROM-WHERE-GROUPBY-HAVING* statements, in fact , all queries expressible by such statements are local.



Conclusion (Cont.)



2. Over ordered domains, with limited arithmetic, reachability queries are most likely inexpressible, but proving this is very hard. Adding more arithmetic operations might help, but only at the expense of encodings which are several thousand digits long – so the new construct is clearly justified.