

ECI 2015 Buenos Aires



# Fundamentos lógicos de bases de datos (Logical foundations of databases)

Diego Figueira

**Gabriele Puppis** 

CNRS LaBRI



- NP complexity class (SAT, 3COL, ....)
- Conjunctive Queries (correspondence with SQL and Relational Algebra)
- Homomorphisms and canonical structure
- Evaluation of CQ (NP-completeness)
- Containment, Equivalence, Minimisation of CQ (NP-completeness)
- Extension to functional dependencies (chased canonical structure)



On graphs: CQ  $\varphi$  is acyclic if  $G_\varphi$  is tree-like



underlying undirected graph is acyclic

**On graphs**: CQ  $\phi$  is **acyclic** if  $G_{\phi}$  is tree-like

On graphs: CQ  $\phi$  is acyclic if  $G_{\phi}$  is tree-like



On graphs: CQ  $\phi$  is acyclic if  $G_{\phi}$  is tree-like



On arbitrary structures: a CQ  $\phi$  is acyclic if it has a join tree

 $\phi(\bar{y}) = \exists \bar{z} . R_1(\bar{z}_1) \land ... \land R_m(\bar{z}_m)$ 

On graphs: CQ  $\varphi$  is acyclic if  $G_\varphi$  is tree-like



On arbitrary structures: a CQ  $\phi$  is acyclic if it has a join tree

$$\varphi(\bar{y}) = \exists \bar{z} . R_1(\bar{z}_1) \land ... \land R_m(\bar{z}_m)$$

- $\bullet$  nodes are the atoms  $R_i(\bar{z}_i)$
- $\bullet$  for every variable x of  $\varphi$  the set of  $\,R_i(\bar z_i)$  's with  $x\in \bar z_i$  forms a subtree of T

On graphs: CQ  $\varphi$  is acyclic if  $G_\varphi$  is tree-like

$$\phi(\mathbf{x},\mathbf{y}) = \exists z . E(\mathbf{x},z) \land E(z,t) \land E(y,z)$$



On arbitrary structures: a CQ  $\phi$  is acyclic if it has a join tree

$$\varphi(\bar{y}) = \exists \bar{z} . R_1(\bar{z}_1) \land ... \land R_m(\bar{z}_m)$$

If x occurs in two nodes, then it occurs in the path linking the two nodes.

- $\bullet$  nodes are the atoms  $R_i(\bar{z}_i)$
- for every variable x of  $\varphi$  the set of  $\,R_i(\bar z_i)$ 's with  $x\in \bar z_i$  forms a subtree of T

On graphs: CQ  $\varphi$  is acyclic if  $G_\varphi$  is tree-like

$$\phi(\mathbf{x},\mathbf{y}) = \exists z . E(\mathbf{x},z) \land E(z,t) \land E(y,z)$$

x Alternatively, if its canonical hyper-graph is α-acyclic.

On arbitrary structures: a CQ  $\phi$  is acyclic if it has a join tree

 $\varphi(\bar{y}) = \exists \bar{z} . R_1(\bar{z}_1) \land ... \land R_m(\bar{z}_m)$ 

- $\bullet$  nodes are the atoms  $R_i(\bar{z}_i)$
- for every variable x of  $\varphi$  the set of  $\,R_i(\bar z_i)$ 's with  $x\in \bar z_i$  forms a subtree of T

On graphs: C

A reduced hypergraph F = (V,E) is *a-acyclic* if for each  $U \subseteq V$ , if  $F|_U$  is connected and has more than one edge, then it has an **articulation set**.

 $\phi(\mathbf{x},\mathbf{y}) = \exists z . E(\mathbf{x},z) \land E(z,t) \land E(y,z)$ 

Alternatively, if its canonical hyper-graph is **α-acyclic**.

On arbitrary structures: a CQ  $\phi$  is acyclic if it has a join tree

 $\varphi(\bar{y}) = \exists \bar{z} . R_1(\bar{z}_1) \land ... \land R_m(\bar{z}_m)$ 

- $\bullet$  nodes are the atoms  $R_i(\bar{z}_i)$
- for every variable x of  $\varphi$  the set of  $\,R_i(\bar z_i)$ 's with  $x\in \bar z_i$  forms a subtree of T

 $\phi(\mathbf{x},\mathbf{y}) = \exists z . E(\mathbf{x},z) \land E(\mathbf{x},t) \land E(\mathbf{y},z)$ 









 $\varphi = \exists x,y,z,t . R(x,y,z) \land S(z,t) \land S(x,z) \land T(z) \land T(x)$ 











The evaluation problem for acyclic CQ sentences is in  $O(|\phi|.|D|)$ 

[Yannakakis]



The evaluation problem for acyclic CQ sentences is in  $O(|\phi|.|D|)$ 

The semi-join  $R \ltimes_{\{i_1=j_1,...,i_n=j_n\}} S = \{ (x_1,...,x_n) \in R \mid \text{there is } (y_1,...,y_m) \in S$ where  $x_{i_k} = y_{j_k}$  for all k} Note:  $R \ltimes_{\{i_1=j_1,...,i_n=j_n\}} S \subseteq R$ 

[Yannakakis]



- 1. Compute the join tree T for  $\varphi$
- 2. Populate the nodes of T with tuples from corresponding relations of D
- 3. For every leaf  $S(x_1,...,x_n)$  with parent  $R(y_1,...,y_m)$  replace tuples of parent with  $R \Join \{i=j | x_i = y_j\} S$ and delete the leaf  $S(x_1,...,x_n)$ .
- 4. Repeat until we are left with one node. If it contains a non-empty relation, then D satisfies  $\phi$ , otherwise it does not.





- 1. Compute the join tree T for  $\phi$
- 2. Populate the nodes of T with tuples from corresponding relations of D
- 3. For every leaf  $S(x_1,...,x_n)$  with parent  $R(y_1,...,y_m)$  replace tuples of parent with  $R \Join \{i=j | x_i = y_j\} S$ and delete the leaf  $S(x_1,...,x_n)$ .
- 4. Repeat until we are left with one node. If it contains a non-empty relation, then D satisfies  $\phi$ , otherwise it does not.



The evaluation problem for acyclic CQ sentences is in  $O(|\varphi|.|D|)$ 

[Yannakakis]

1. Compute the join tree T for  $\varphi$ 

remove all the tuples from the parent that do not match a tuple from the child

- 2. Populate the nodes of T with tuples from correspond
- 3. For every leaf  $S(x_1,...,x_n)$  with parent  $R(y_1,...,y_m)$  replace tuples of parent with  $R \Join \{i=j | x_i = y_j\} S$ and delete the leaf  $S(x_1,...,x_n)$ .
- 4. Repeat until we are left with one node. If it contains a non-empty relation, then D satisfies  $\phi$ , otherwise it does not.













```
\phi = \exists x, y, z, t . R(x, y, z) \land S(z, t) \land S(x, z) \land T(z) \land T(x)
R = \{(1,4,4), (4,1,4)\}
S = \{(4,5), (5,2), (4,4)\}
                                             S(z,t) : \{(4,5), (5,2), (4,4)\}
T = \{1, 2, 3, 4\}
                                             S(x,z) : \{(4,5), (5,2), (4,4)\}
```



#### How to compute a join tree?

GYO reducts [Graham, Yu, Ozsoyoglu]

An ear of a hypergraph (V,E) is a hyperedge e in E such that one of the following conditions holds:

(1) There is a witness e' in E, such that  $e' \neq e$  and each

vertex from e is either

(a) **only** in **e** or

(b) in **e'**; or

(2) e has no intersection with any other hyperedge.











Definition: The GYO **reduct** of a hyper-graph is the result of removing ears until no more ears are left.

Definition: The GYO **reduct** of a hyper-graph is the result of removing ears until no more ears are left.

Theorem: TFAE

- The GYO reduct of a hyper graph G is empty
- $\bullet$  A CQ  $\varphi$  having G as underlying canonical hyper-graph is acyclic
- The hyper graph G is  $\alpha$ -acyclic

Definition: The GYO **reduct** of a hyper-graph is the result of removing ears until no more ears are left.

Theorem: TFAE

- The GYO reduct of a hyper graph G is empty
- $\bullet$  A CQ  $\varphi$  having G as underlying canonical hyper-graph is acyclic
- The hyper graph G is  $\alpha$ -acyclic

We can test acyclicity by computing the GYO reduct!

#### Ears?










Acyclic!

#### How to compute a join tree?

#### How to compute a join tree?

**GYO algorithm** [Graham, Yu, Ozsoyoglu] Given the query  $Q = \{ R_1(X_1),...,R_n(X_n) \}$ Consider its canonical structure  $G_Q$ For  $R_i(X_i)$  an ear with witness  $R_j(Y_j)$ Put an edge between  $R_i(X_i)$  and  $R_j(X_j)$ , and remove  $R_i$  from Q. Repeat.

> **E.g.** R(x,y,z), S(x,y), T(x,x), R(x,x,y), T(y,y)

#### How to compute a join tree?

**GYO algorithm** [Graham, Yu, Ozsoyoglu] Given the query  $Q = \{ R_1(X_1),...,R_n(X_n) \}$ Consider its canonical structure  $G_Q$ For  $R_i(X_i)$  an ear with witness  $R_j(Y_j)$ Put an edge between  $R_i(X_i)$  and  $R_j(X_j)$ , and remove  $R_i$  from Q. Repeat.

> E.g. R(x,y,z), S(x,y), T(x,x), R(x,x,y), T(y,y) T(x,x) R(x,y,z) R(x,x,y) S(x,y)T(y,y)

#### How to compute a join tree?

E.g.  

$$R(x,y,z), S(x,y), T(x,x), R(x,x,y), T(y,y)$$
  
 $T(x,x)$   
 $R(x,y,z)$   
 $R(x,x,y)$   
 $S(x,y)$   
 $T(y,y)$ 

#### How to compute a join tree?

E.g.  

$$R(x,y,z), S(x,y), T(x,x), R(x,x,y), T(y,y)$$

$$T(x,x)$$

$$R(x,y,z) - R(x,x,y)$$

$$S(x,y) - T(y,y)$$

#### How to compute a join tree?

E.g.  

$$R(x,y,z), S(x,y), T(x,x), R(x,x,y), T(y,y)$$
  
 $T(x,x)$   
 $R(x,y,z)$   
 $R(x,y,z)$   
 $R(x,x,y)$ 

#### How to compute a join tree?



#### How to compute a join tree?

**GYO algorithm** [Graham, Yu, Ozsoyoglu] Given the query  $Q = \{ R_1(X_1),...,R_n(X_n) \}$ Consider its canonical structure  $G_Q$ For  $R_i(X_i)$  an ear with witness  $R_j(Y_j)$ Put an edge between  $R_i(X_i)$  and  $R_j(X_j)$ , and remove  $R_i$  from Q. Repeat.

Remove ears until you're left with only one!



E.g.  

$$R(x,y,z), S(x,y), T(x,x), R(x,x,y), T(y,y)$$
  
 $T(x,x)$   
 $R(x,y,z)$   
 $R(x,y,z)$   
 $R(x,x,y)$   
 $R(x,x,y)$ 



#### [Gottlob, Leone, Scarcello]

- Evaluation problem for boolean ACQ's is LOGCFL-complete
- $NL \subseteq LOGCFL \subseteq AC^1 \subseteq NC^2 \subseteq P$

the class of problems logspace-reducible to a context-free language

# Beyond acyclic CQ's

- (Hyper-)Treewidth = a measure of the cyclicity of (hyper-)graphs tw(Q) = tw(canonical hypergraph of Q) tw = 1 are forests, tw = 2 are graphs without K<sub>4</sub> as a minor, ...
- For fixed k,

computing whether Q has tw ≤k and calculating a tree decomposition can be done in **linear time** 

# Beyond acyclic CQ's

- (Hyper-)Treewidth = a measure of the cyclicity of (hyper-)graphs tw(Q) = tw(canonical hypergraph of Q) tw = 1 are forests, tw = 2 are graphs without K<sub>4</sub> as a minor, ...
- For fixed k,

computing whether Q has tw ≤k and calculating a tree decomposition can be done in **linear time** 

Bounded tree width queries = a class of CQ's so that for some k, every query has tw  $\leq k$ 

## Beyond acyclic CQ's

CQ's with bounded treewidth can be evaluated in PTIME

CQ's with bounded treewidth can be evaluated in PTIME

Containment of CQ's with bounded treewidth is in PTIME



 $\varphi \subseteq \psi \text{ iff } \psi(G_{\varphi}) \neq \emptyset$ 

Containment of CQ's with bounded treewidth is in PTIME

CQ's with bounded treewidth can be evaluated in PTIME

Containment of CQ's with bounded treewidth is in PTIME

A class *C* of CQ's can be evaluated in PTIME iff they have bounded tree width!

For graphs, assuming  $W[1] \neq FPT$ , and C a r.e. class of graphs. [Grohe, Schwentick, Segoufin]

# Querying with semi-joins

The semi-join  $R \ltimes_{\{i_1=j_1,...,i_n=j_n\}} S = \{ (x_1,...,x_n) \in R \mid \text{there is } (y_1,...,y_m) \in S$ where  $x_{i_k} = y_{j_k}$  for all k}

#### The semi-join algebra (SA): variant of RA with operations:

 $\ltimes$ ,  $\cup$ ,  $\pi$ ,  $\sigma$ ,  $\setminus$ , *dupcol* 

# Querying with semi-joins

The semi-join  $R \ltimes_{\{i_1=j_1,...,i_n=j_n\}} S = \{ (x_1,...,x_n) \in R \mid \text{there is } (y_1,...,y_m) \in S$ where  $x_{i_k} = y_{j_k}$  for all k}

### The semi-join algebra (SA): variant of RA with operations:

 $\ltimes$ ,  $\cup$ ,  $\pi$ ,  $\sigma$ ,  $\setminus$ , *dupcol* 

Output at most linear in the database. Further,

The evaluation problem for SA is in  $O(|\phi|.|D|)$ 

Logical characterisation: "stored-tuples guarded fragment of FO"

- every intermediate relation is **linear** in |D|
- we apply  $|\varphi|$  semi-joins

What if we allow intermediate relations to be **polynomial** in |D|?

Def.

### $FO^k$ = The fragment of FO restricted to k variable names





G



#### Def. FO<sup>k</sup> = The fragment of FO restricted to k variable names

 $\begin{aligned} \varphi(x) &= \text{``Every neighbour of } x \text{ has an outgoing path of length 2''} \\ &= \forall y. \left( E(x, y) \implies \exists z \exists w \left( E(y, z) \land E(z, w) \right) \right) \in FO^4 \\ &= \forall y. \left( E(x, y) \implies \exists x \left( E(y, x) \land \exists y E(x, y) \right) \right) \in FO^2 \end{aligned}$ 

G











The evaluation problem for  $FO^k$  is in PTIME (combined c.)

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Algorithm for a FO<sup>k</sup> formula  $\psi$  of **quantifier rank r**:

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...)))))$ "

Algorithm for a FO<sup>k</sup> formula  $\psi$  of quantifier rank r:

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...)))))$ "

Algorithm for a FO<sup>k</sup> formula  $\psi$  of quantifier rank r:

qr 1

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...)))))$ "

Algorithm for a FO<sup>k</sup> formula  $\psi$  of **quantifier rank r**:

qr 1

qr 2

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

→ Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...))))$ " qr 0 qr 1

Algorithm for a FO<sup>k</sup> formula  $\psi$  of quantifier rank r:

- 1. Evaluate qr=0 subformulas  $\alpha$  and output result in relations  $R_{0,\alpha}$
- 2. Evaluate qr=1 subformulas  $\beta$  based on  $R_{0,\alpha}$  and output in  $R_{1,\beta}$
- 3. Evaluate qr=2 subformulas  $\gamma$  based on  $R_{1,\beta}$  and output in  $R_{1,\gamma}$
- 4. . . .

qr2
## Bounded variable FO

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...)))))$ "

Algorithm for a FO<sup>k</sup> formula  $\psi$  of **quantifier rank r**:

- 1. Evaluate qr=0 subformulas  $\alpha$  and output result in relations  $R_{0,\alpha}$  $\rightsquigarrow |V|^k \cdot (|\alpha| \cdot |G|)^p$
- 2. Evaluate qr=1 subformulas  $\beta$  based on  $R_{0,\alpha}$  and output in  $R_{1,\beta}$
- 3. Evaluate qr=2 subformulas  $\gamma$  based on  $R_{1,\beta}$  and output in  $R_{1,\gamma}$
- 4. ... : r. ...

qr1

qr 2

## Bounded variable FO

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...)))))$ "

Algorithm for a FO<sup>k</sup> formula  $\psi$  of quantifier rank r:

1. Evaluate qr=0 subformulas  $\alpha$  and output result in relations  $R_{0,\alpha}$  $\rightsquigarrow |V|^k \cdot (|\alpha| \cdot |G|)^p$ 

2. Evaluate qr=1 subformulas  $\beta$  based on  $R_{0,\alpha}$  and output in  $R_{1,\beta}$  $\rightsquigarrow |V|^k \cdot (|\beta| \cdot (|G| + |R_1|))^p$ 

3. Evaluate qr=2 subformulas  $\gamma$  based on  $R_{1,\beta}$  and output in  $R_{1,\gamma}$ 



qr 1

qr 2

## Bounded variable FO

The evaluation problem for  $FO^k$  is in PTIME (combined c.)

Maximum number of nested quantifiers " $\exists x(... \forall y(... \exists x(... \exists z(...)))))$ "

Algorithm for a FO<sup>k</sup> formula  $\psi$  of quantifier rank r:

1. Evaluate qr=0 subformulas  $\alpha$  and output result in relations  $R_{0,\alpha}$  $\rightsquigarrow |V|^k \cdot (|\alpha| \cdot |G|)^p$ 

2. Evaluate qr=1 subformulas  $\beta$  based on  $R_{0,\alpha}$  and output in  $R_{1,\beta}$  $\rightsquigarrow |V|^k \cdot (|\beta| \cdot (|G| + |R_1|))^p$ 

3. Evaluate **qr=2** subformulas  $\gamma$  based on  $R_{1,\beta}$  and output in  $R_{1,\gamma} \leq |V|^k$  $\Rightarrow |V|^k \cdot (|\gamma| \cdot (|G| + |R_2|))^p \leq |V|^k$ 



qr 1

gr 2

Desirable:

• Given k and a FO query  $\phi$ , is  $\phi$  in FO<sup>k</sup>?  $\longrightarrow$  Undecidable (even w.o.  $\neg$ )

Desirable:

• Given k and a FO query  $\phi$ , is  $\phi$  in FO<sup>k</sup>?  $\longrightarrow$  Undecidable (even w.o.  $\neg$ )

• Given k and a CQ query  $\phi$ , is  $\phi$  in FO<sup>k</sup>?  $\longrightarrow$  NP-complete

Desirable:

• Given k and a FO query  $\phi$ , is  $\phi$  in FO<sup>k</sup>?  $\longrightarrow$  Undecidable (even w.o.  $\neg$ )

• Given k and a CQ query  $\phi$ , is  $\phi$  in FO<sup>k</sup>?  $\longrightarrow$  NP-complete

• Satisfiability for  $FO^k$ 

→ Undecidable if  $k \ge 3$  (Domino)

••• NEXPTIME-complete if k=2



#### Equivalence-RA

#### Equivalence-SQL

Equivalence-FO

Sat-FO

Domino

UNDECIDABLE

Eval-FO (combined)

QBF

**PSPACE** 

Eval-CQ (combined)

3COL

SAT

NP

Eval-FO<sup>k</sup> (combined)

**PTIME** 

Cont-CQ

Eval-FO (data)

LOGSPACE

## Descriptive complexity

What properties can be checked efficiently? E.g. 3COL can be tested in NP

Metatheorem "A property can be expressed in [insert some logic here] iff it can be checked in [some complexity class here]"

## Descriptive complexity

What properties can be checked efficiently? E.g. 3COL can be tested in NP

Metatheorem "A property can be expressed in [insert some logic here] iff it can be checked in [some complexity class here]"

 $\rightsquigarrow$  "A property is FO-definable iff it can be tested in AC<sup>0</sup>"

----- "A property is <code>3SO-definable</code> iff it can be tested in NP" [Fagin 73]

---> Open problem: which logic captures PTIME?

## Recursion



## Recursion

Can we enhance query languages with recursion ? E.g. express reachability properties

Datalog (semantics based on least fixpoint)
Ancestor(X,Y) :- Parent(X,Z), Ancestor(Z,Y)
Ancestor(X,X) :- .
?- Ancestor("Louis XIV",Y)

## Recursion

Can we enhance query languages with recursion ? E.g. express reachability properties

Datalog (semantics based on least fixpoint)
Ancestor(X,Y) :- Parent(X,Z), Ancestor(Z,Y)
Ancestor(X,X) :- .
?- Ancestor("Louis XIV",Y)

---> Incomparable with FO (has recursion, but is monotone)

---> Evaluation is in PTIME (for data complexity, but also for bounded arity)

## Semi-structured data

Tree-structured or graph-structures dbs in place of relational dbs.



## Semi-structured data

Tree-structured or graph-structures dbs in place of relational dbs.



Satisfiability for FO<sup>2</sup>[↓,~] is decidable [Bojanczyk, Muscholl, Schwentick, Segoufin 09]

## Incomplete information



# Incomplete information

How to correctly treat NULL values, missing tuples, noisy data ?



# Incomplete information

How to correctly treat NULL values, missing tuples, noisy data ?



#### Recap

- Relational Algebra = simple SQL = FO on active domain
- Evaluation, Satisfiability, Equivalence, Containment
- Data / Combined complexity
- Expressiveness of FO: EF games, 0-1 law, Rado structure, Locality
- Conjunctive Queries: Homomorphism lemma, Canonical structures
- Acyclic CQ
- FO<sup>k</sup>

• Abiteboul, Hull, Vianu, "Foundations of Databases", Addison-Wesley, 1995.

(freely available at <a href="http://webdam.inria.fr/Alice/">http://webdam.inria.fr/Alice/</a>)

• Libkin, "Elements of Finite Model Theory", Springer, 2004.

• Immerman, "Descriptive Complexity", Springer, 1999.