

# DATABASE THEORY

## Lecture 5: Conjunctive Queries

Markus Krötzsch

TU Dresden, 28 April 2016

### Review: FO Query Complexity

The evaluation of FO queries is

- PSPACE-complete for combined complexity
- PSPACE-complete for query complexity
- AC<sup>0</sup>-complete for data complexity

↪ PSPACE is rather high

↪ Are there relevant query languages that are simpler than that?

## Overview

1. Introduction | Relational data model
2. First-order queries
3. Complexity of query answering
4. Complexity of FO query answering
5. Conjunctive queries
6. Tree-like conjunctive queries
7. Query optimisation
8. Conjunctive Query Optimisation / First-Order Expressiveness
9. First-Order Expressiveness / Introduction to Datalog
10. Expressive Power and Complexity of Datalog
11. Optimisation and Evaluation of Datalog
12. Evaluation of Datalog (2)
13. Graph Databases and Path Queries
14. Outlook: database theory in practice

See course homepage [⇒ link] for more information and materials

### Conjunctive Queries

Idea: restrict FO queries to conjunctive, positive features

#### Definition

A **conjunctive query** (CQ) is an expression of the form

$$\exists y_1, \dots, y_m. A_1 \wedge \dots \wedge A_\ell$$

where each  $A_i$  is an atom of the form  $R(t_1, \dots, t_k)$ . In other words, a conjunctive query is an FO query that only uses conjunctions of atoms and (outer) existential quantifiers.

Example: “Find all lines that depart from an accessible stop” (as seen in earlier lectures)

$$\exists y_{\text{SID}}, y_{\text{Stop}}, y_{\text{To}}. \text{Stops}(y_{\text{SID}}, y_{\text{Stop}}, \text{"true"}) \wedge \text{Connect}(y_{\text{SID}}, y_{\text{To}}, x_{\text{Line}})$$

# Conjunctive Queries in Relational Calculus

The expressive power of CQs can also be captured in the relational calculus

## Definition

A **conjunctive query** (CQ) is a relational algebra expression that uses only the operations select  $\sigma_{n=m}$ , project  $\pi_{a_1, \dots, a_n}$ , join  $\bowtie$ , and renaming  $\delta_{a_1, \dots, a_n \rightarrow b_1, \dots, b_n}$ .

Renaming is only relevant in named perspective

→ CQs are also known as **SELECT-PROJECT-JOIN queries**

# Boolean Conjunctive Queries

A **Boolean conjunctive query** (BCQ) asks for a mapping from query variables to domain elements such that all atoms are true

**Example:** "Is there an accessible stop where some line departs?"

$\exists y_{SID}, y_{Stop}, y_{To}, y_{Line} \cdot \text{Stops}(y_{SID}, y_{Stop}, \text{"true"}) \wedge \text{Connect}(y_{SID}, y_{To}, y_{Line})$

Stops:

SID	Stop	Accessible
17	Hauptbahnhof	true
42	Helmholtzstr.	true
57	Stadtgutstr.	true
123	Gustav-Freytag-Str.	false
...	...	...

Connect:

From	To	Line
57	42	85
17	789	3
...	...	...

# Extensions of Conjunctive Queries

Two features are often added:

- **Equality:** CQs with equality can use atoms of the form  $t_1 \approx t_2$  (in relational calculus: table constants)
- **Unions:** unions of conjunctive queries are called UCQs (in this case the union is only allowed as outermost operator)

Both extensions truly increase expressive power (as shown in exercise)

**Features omitted on purpose:** negation and universal quantifiers  
→ the reason for this is query complexity (as we shall see)

# How Hard is it to Answer CQs?

If we know the variable mappings, it is easy to check:

- Checking if a single ground atom  $R(c_1, \dots, c_k)$  holds can be done in linear time
- Checking if a conjunction of ground atoms holds can be done in quadratic time

→ A candidate BCQ match can be verified in P

(There are  $n^m$  candidates:  $n$  size of domain;  $m$  number of query variables)

## Theorem

BCQ query answering is in NP for combined complexity (and also for query complexity).

→ Better than PSPACE (presumably)

## Can we do any better?

Not really. To see this, let's look at some other problems.

Consider two relational structures  $\mathcal{I}$  and  $\mathcal{J}$   
 (= database instances, interpretations, hypergraphs)

### Definition

A **homomorphism**  $h$  from  $\mathcal{I}$  to  $\mathcal{J}$  is a function  $h : \Delta^{\mathcal{I}} \rightarrow \Delta^{\mathcal{J}}$  such that, for all relation names  $R$ :

$$\text{if } \langle d_1, \dots, d_n \rangle \in R^{\mathcal{I}} \text{ then } \langle h(d_1), \dots, h(d_n) \rangle \in R^{\mathcal{J}}.$$

The **homomorphism problem** is the question if there is a homomorphism from  $\mathcal{I}$  to  $\mathcal{J}$ .

## BCQ Answering as Homomorphism Problem

The homomorphism problem can be reduced to BCQ answering:

- A relational structure  $\mathcal{I}$  gives rise to a CQ  $Q_{\mathcal{I}}$ :  
 replace domain elements by variables (one-to-one); add one query atom per relational tuple; existentially quantify all variables
- $\mathcal{I}$  has a homomorphism to  $\mathcal{J}$  if and only if  $\mathcal{J} \models Q_{\mathcal{I}}$

BCQ answering can be reduced to the homomorphism problem:

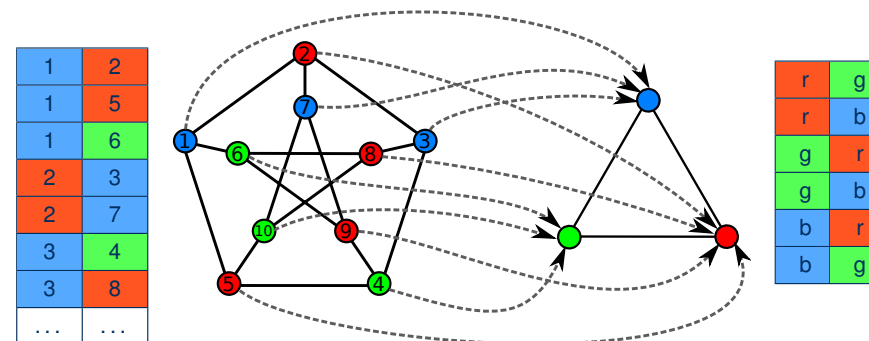
- Clear for BCQs that don't contain constants
- Eliminate query constants  $a$ : create new relation  $R_a = \{\langle a \rangle\}$ ;  
 replace  $a$  by a fresh variable  $x$  and add a query atom  $R_a(x)$

$\leadsto$  both problems are equivalent

## Example: Three-colouring as Homomorphism

$\mathcal{I}$ :

$\mathcal{J}$ :



3-colouring is NP-hard

$\leadsto$  the homomorphism problem is NP-hard

## Complexity of Conjunctive Query Answering

We showed that BCQ answering is in NP and that the homomorphism problem is NP-hard, therefore:

### Theorem

BCQ answering is

- NP-complete for combined complexity
- NP-complete for query complexity
- in  $AC^0$  for data complexity (inherited from FO queries)

# Constraint Satisfaction Problems

Another important problem equivalent to BCQ answering

## Definition

A **constraint satisfaction problem** (CSP) over a domain  $\Delta$  is given by a set of variables  $\{x_1, \dots, x_n\}$  and a set of constraints  $\{C_1, \dots, C_m\}$ , where each constraint  $C_i$  has the form  $\langle X_i, R_i \rangle$  with

- $X_i$  a list of variables from  $\{x_1, \dots, x_n\}$ ,
- $R_i$  a  $|X_i|$ -ary relation over  $\Delta$ .

A **solution** to the CSP is an assignment of variables to values from  $\Delta$  such that all constraints are satisfied (=all tuples occur in the respective relations).

↪ alternative notation for BCQ answering/homomorphism problem

# Equivalent Problems

Summing up, the following problems are equivalent:

- Answering a conjunctive query over a database instance
- Finding a homomorphism from a relational structure to another
- Solving a constraint satisfaction problem

Each of these problems is NP-complete

# CSP Example

A combinatorial crossword puzzle:

Domain:  $\Delta = \{A, \dots, Z\}$

Variables:  $x_1, \dots, x_{26}$

Constraints:

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	■	$x_6$
$x_7$	■	■	■	$x_8$	$x_9$	$x_{10}$
$x_{11}$	$x_{12}$	$x_{13}$	■	$x_{14}$	■	$x_{15}$
$x_{16}$	■	$x_{17}$	■	$x_{18}$	■	$x_{19}$
$x_{20}$	$x_{21}$	$x_{22}$	$x_{23}$	$x_{24}$	$x_{25}$	$x_{26}$

1 vertically:

H	E	A	R	T
H	O	N	E	Y
I	R	O	N	Y
L	O	G	I	C

1 horizontally:

H	A	P	P	Y
I	N	F	E	R
L	A	B	O	R
L	A	T	E	R

5 vertically:

R	A	D	I	O
R	E	T	R	O
Y	A	C	H	T
Y	E	R	B	A

...

# Towards Better Complexities

NP-complete problems are still intractable

↪ can we do better?

Problem: searching a match may require backtracking, eventually exploring all options

H	A	P	P	Y	■	■
O	■	■	■	A	■	■
N	E	W	■	C	■	■
E	■	A	■	H	■	■
Y	■	Y	■	T	■	■

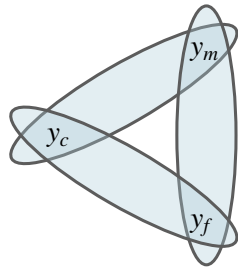
Intuition: life would be easier if we would not have to go back so much ...

↪ the problem is with the **cycles**

## Example: Cyclic CQs

“Is there a child whose parents are married with each other?”

$$\exists y_c, y_m, y_f. \text{mother}(y_c, y_m) \wedge \text{father}(y_c, y_f) \wedge \text{married}(y_m, y_f)$$

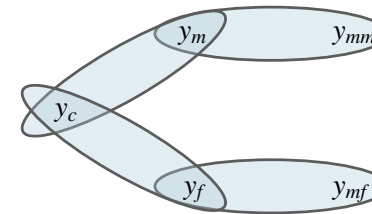


↪ cyclic query

## Example: Acyclic CQs

“Is there a child whose parents are married with someone?”

$$\exists y_c, y_m, y_f, y_{mm}, y_{mf}. \text{mother}(y_c, y_m) \wedge \text{father}(y_c, y_f) \wedge \text{married}(y_m, y_{mm}) \wedge \text{married}(y_{mf}, y_f)$$

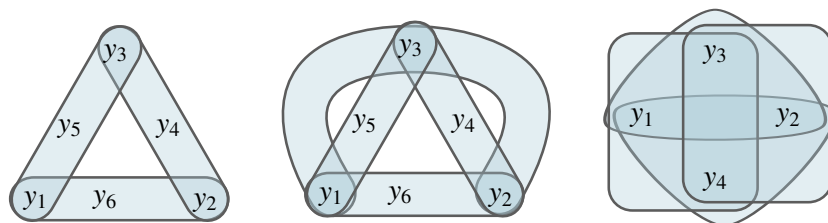


↪ acyclic query

## Defining Acyclic Queries

Queries in general are hypergraphs

↪ What does “acyclic” mean?



View hypergraphs as graphs to check acyclicity?

- **Primal graph:** same vertices; edges between each pair of vertices that occur together in a hyperedge
- **Incidence graph:** vertices and hyperedges as vertices, with edges to mark incidence (bipartite graph)

However: both graphs have cycles in almost all cases

## Acyclic Hypergraphs

**GYO-reduction** algorithm to check acyclicity:

(after Graham [1979] and Yu & Özsoyoğlu [1979])

Input: hypergraph  $H = \langle V, E \rangle$  (we don't need relation labels here)

Output: GYO-reduct of  $H$

Apply the following simplification rules as long as possible:

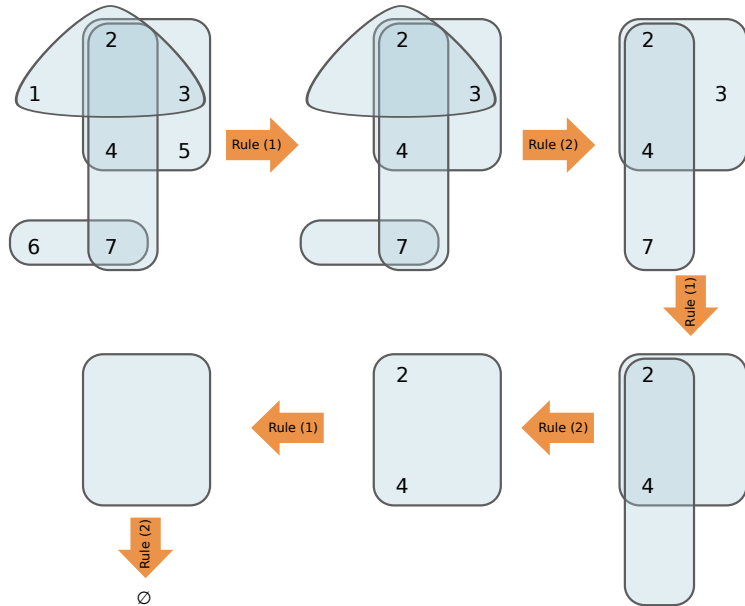
- (1) Delete all vertices that occur in at most one hyperedge
- (2) Delete all hyperedges that are empty or that are contained in other hyperedges

### Definition

A hypergraph is **acyclic** if its GYO-reduct is  $\langle \emptyset, \emptyset \rangle$ .

A CQ is **acyclic** if its associated hypergraph is.

## Example 1: GYO-Reduction

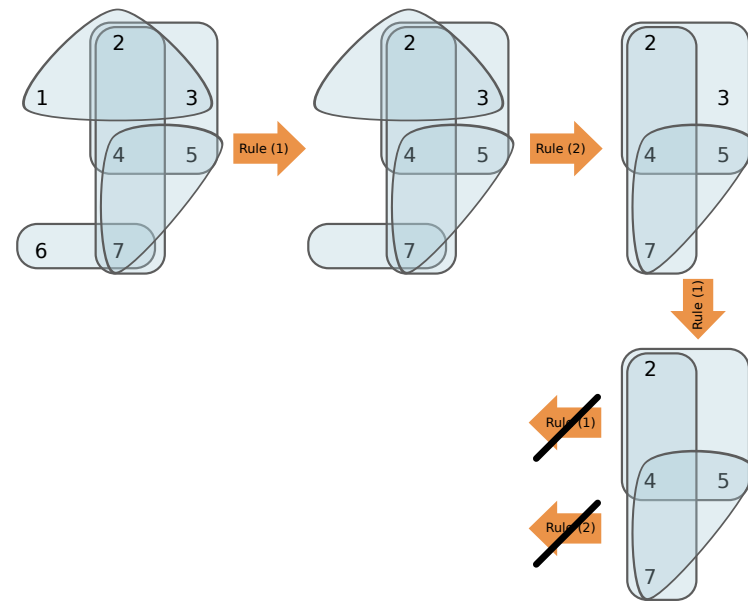


Markus Krötzsch, 28 April 2016

Database Theory

slide 21 of 31

## Example 2: GYO-Reduction



Markus Krötzsch, 28 April 2016

Database Theory

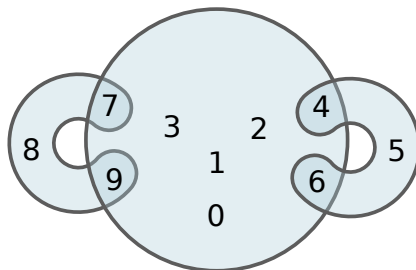
slide 22 of 31

## Alternative Version of GYO-Reduction

An **ear** of a hypergraph  $\langle V, E \rangle$  is a hyperedge  $e \in E$  that satisfies one of the following:

- (1) there is an edge  $e' \in E$  such that  $e \neq e'$  and every vertex of  $e$  is either only in  $e$  or also in  $e'$ , or
- (2)  $e$  has no intersection with any other hyperedge.

Example:



→ edges  $\langle 4, 5, 6 \rangle$  and  $\langle 7, 8, 9 \rangle$  are ears

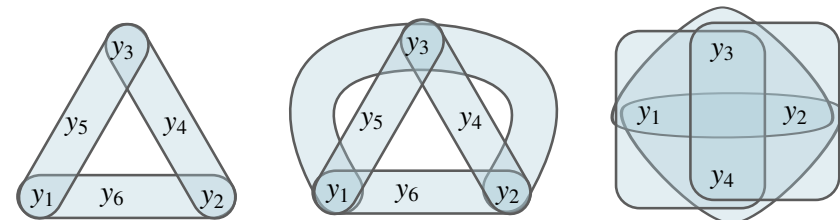
Markus Krötzsch, 28 April 2016

Database Theory

slide 23 of 31

## Examples

Any ears?



Markus Krötzsch, 28 April 2016

Database Theory

slide 24 of 31

# GYO'-Reduction

Input: hypergraph  $H = \langle V, E \rangle$

Output: GYO'-reduct of  $H$

Apply the following simplification rule as long as possible:

- Select an ear  $e$  of  $H$
- Delete  $e$
- Delete all vertices that only occurred in  $e$

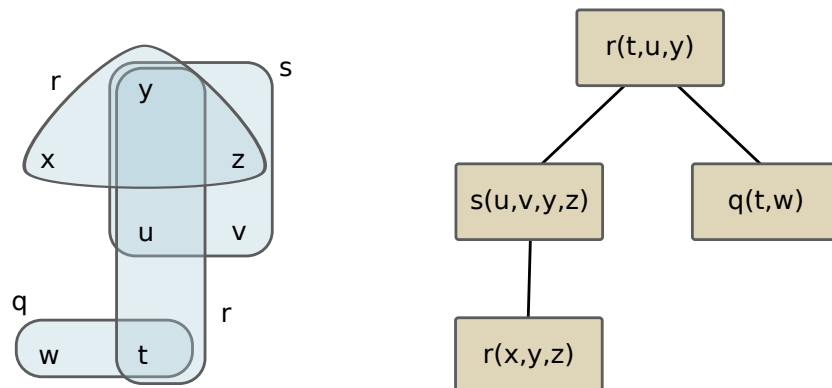
## Theorem

The GYO-reduct is  $\langle \emptyset, \emptyset \rangle$  if and only if the GYO'-reduct is  $\langle \emptyset, \emptyset \rangle$

↪ alternative characterization of acyclic hypergraphs

# Example: Join Tree

$$\exists x, y, z, t, u, v, w. (r(x, y, z) \wedge r(t, u, y) \wedge s(u, v, y, z) \wedge q(t, w))$$



# Join Trees

Both GYO algorithms can be implemented in linear time

Open question: what benefit does BCQ acyclicity give us?

Fact: if a BCQ is acyclic, then it has a join tree

## Definition

A **join tree** of a (B)CQ is an arrangement of its query atoms in a tree structure  $T$ , such that for each variable  $x$ , the atoms that refer to  $x$  are a connected subtree of  $T$ .

A (B)CQ that has a join tree is called a **tree query**.

# Processing Join Trees Efficiently

Join trees can be processed in polynomial time

Key ingredient: the semijoin operation

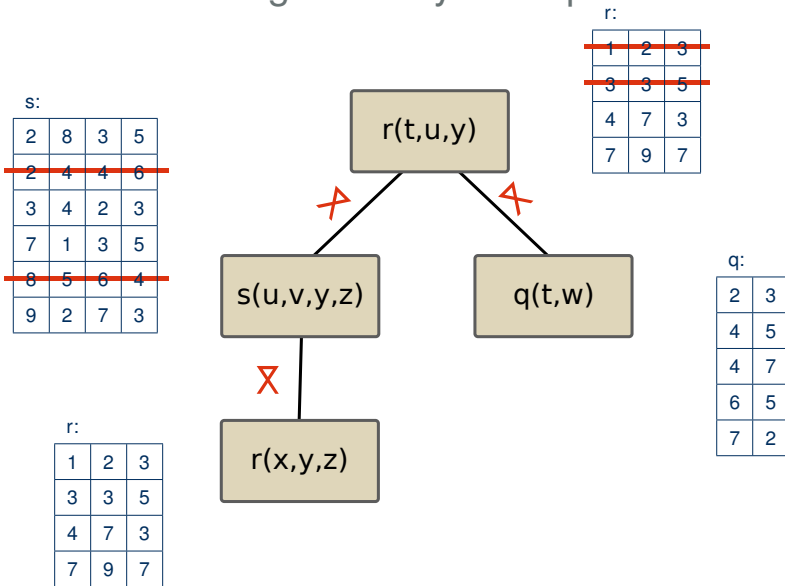
## Definition

Given two relations  $R[U]$  and  $S[V]$ , the **semijoin**  $R^J \bowtie S^J$  is defined as  $\pi_U(R^J \bowtie S^J)$ .

Join trees can now be processed by computing semijoins bottom-up

↪ Yannakakis' Algorithm

# Yannakakis' Algorithm by Example



# Yannakakis' Algorithm: Summary

Polynomial time procedure for answering BCQs

Does not immediately compute answers in the version given here  
 ~> modifications needed

Even tree queries can have exponentially many results,  
 but each can be computed (not just checked) in P  
 ~> **output-polynomial** computation of results

## Summary and Outlook

Conjunctive queries (CQs) are an important special case of FO queries

Boolean CQ answering, the homomorphism problem and constraint satisfaction problems are equivalent and NP-complete

CQ answering is simpler, namely in P, when CQs are tree queries

- Check acyclicity with GYO algorithm
- Evaluate query using Yannakakis' Algorithm

Open questions:

- Tree queries are rather special. Are there more general conditions for good queries?
- What about query optimisation?