

# Data Complexity of Answering Ontology-Mediated Queries with a Covering Axiom

Olga Gerasimova

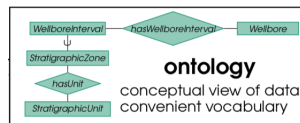
*Research advisor:* professor M.V. Zakharyashev

Department of Data Analysis and Artificial Intelligence  
Faculty of Computer Science  
National Research University Higher School of Economics

19 November, 2018

An **ontology** is an engineering artefact that

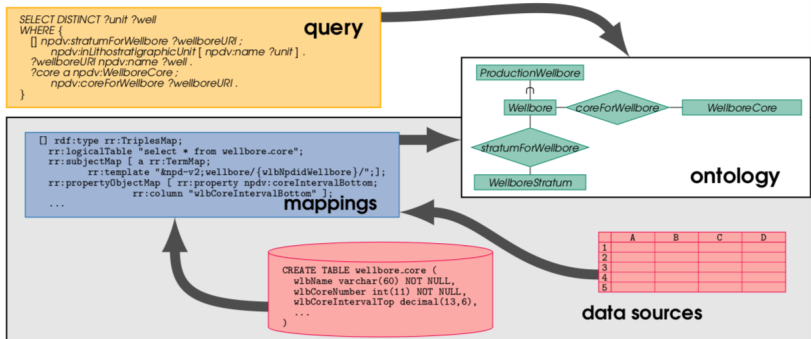
- is constituted by a **specific vocabulary** used to describe a certain domain of interest;
- is a set of **explicit assumptions (rules)** regarding the intended meaning of the vocabulary;
- is a **formal and machine manipulable** model of a domain of interest.



## Definition

**Ontology** – a formal report of *individuals*, *classes* and *properties* (or relations) that are considered in the framework of a domain of interest.

# Ontology-based data access (OBDA)



- 1 gives a high-level conceptual view of the data
- 2 provides the user with a convenient vocabulary for queries
- 3 allows the system to enrich incomplete data with background knowledge
- 4 supports queries to multiple and heterogeneous data sources

# Mapping

For example, we have the table

movie ID	<i>title</i>	year	kind
728	'Django Unchained'	2012	1
1257	'Game of Thrones'	2011	2
2543	'Blade Runner'	1982	1

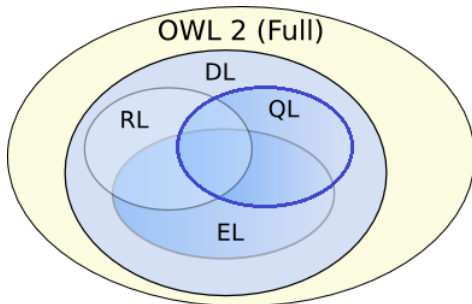
and the ontology vocabulary containing **classes** such as *Movie*, *Title*, *Year*, *Kind* and **relations** such as *has\_title*, *was\_released\_in* and *has\_kind*.

We need to translate the table into the data schemes related to the ontology vocabulary via mappings. Then, we obtain:

- *Movie*(728), *Movie*(1257), *Movie*(2543), *Title*(Django Unchained), *Title*(Game of Thrones), *Title*(Blade Runner), *Year*(2012), *Year*(2011), *Year*(1982), *Kind*(1), *Kind*(2);
- *has\_title*(728,Django Unchained), *was\_released\_in*(728,2012), *has\_kind*(728,1), ...

## Definition

**Description Logic (DL)** is an area of knowledge representation and reasoning in Artificial Intelligence and the Semantic Web that studies **logic-based formalisms** whose languages operate with **concepts** to represent classes of individuals in an application domain, and **roles** to represent binary relations between the individuals.



**OWL2QL** (OWL – Web Ontology Language) contains *individual names*  $a_i$ , *concept names*  $A_i$ , and *role names*  $P_i$  ( $i = 1; 2; \dots$ ).

$$\begin{aligned} R & ::= P_i \mid P_i^-, \\ B & ::= \perp \mid A_i \mid \exists R, \\ C & ::= B \mid \exists R.B \end{aligned}$$

**TBox**,  $T$ , – a finite set of (1) concept and role inclusions, (2) concept and role disjointness constraints of the forms

$$\begin{aligned} B \sqsubseteq C, \quad R_1 \sqsubseteq R_2 \\ B_1 \sqcap B_2 \sqsubseteq \perp, \quad R_1 \sqcap R_2 \sqsubseteq \perp \end{aligned}$$

**ABox**,  $A$ , – a finite set of assertions of the form  $A_k(a_i)$  and  $P_k(a_i; a_j)$  and *inequality constraints*  $a_i \neq a_j$  for  $i \neq j$ .

$T$  and  $A$  together constitute **the knowledge base**  $K = (T; A)$ .

- The user formulates a *query*  $q$  in the vocabulary of a given *ontology*  $T$ .

## The task of an OBDA system

To 'rewrite'  $q$  and  $T$  into a new query  $q_0$  in the vocabulary of the data such that, for any possible data  $A$  (in this vocabulary), the answers to  $q$  over  $(T; A)$  are precisely the same as the answers to  $q_0$  over  $A$ .

# Example

- (1)  $Dancer \sqsubseteq \exists isTrainedBy.Choreographer$       (5)  $LatinaCh \sqsubseteq Choreographer$   
(2)  $\exists participates.Competition \sqsubseteq Dancer$       (6)  $StandartCh \sqcap LatinaCh \sqsubseteq \perp$   
(3)  $Choreographer \sqsubseteq MasterOfSport$       (7)  $isTrainedBy^- \sqsubseteq trains$   
(4)  $StandartCh \sqsubseteq Choreographer$       (8)  $trains \sqsubseteq isTrainedBy^-$

$$q(x) = \exists y (MasterOfSport(y) \wedge trains(y, x))$$

$$q'(x) = \exists y [(MasterOfSport(y) \vee StandartCh(y) \vee LatinaCh(y) \vee Choreographer(y) \wedge (trains(y,x) \vee isTrainedBy(x, y)))] \vee Dancer(x) \vee \exists z participates(x, z)$$

$$\mathcal{A} = \{Dancer(a), Choreographer(b), trains(b, c), participates(d, e)\}$$

Answer to  $q(x)$  according to  $\mathcal{A}$ :  $\{\emptyset\}$

Answer to  $q'(x)$  according to  $\mathcal{A}$ :  $\{a, c, d\} \Leftrightarrow q(x)$  according to  $(\mathcal{T}; \mathcal{A})$



**Big problem 1:** how to determine the data complexity of answering a given ontology-mediated query  $(\mathcal{O}, \mathbf{q})$  with an ‘expressive’ ontology  $\mathcal{O}$  and a Boolean conjunctive query  $\mathbf{q}$ ?

That is, the complexity of deciding whether  $\mathcal{O}, \mathcal{D} \models \mathbf{q}$  for a given input data instance  $\mathcal{D}$

**Big problem 1':** how to determine whether  $(\mathcal{O}, \mathbf{q})$  is

- FO-rewritable, that is, whether there is an FO-sentence  $\mathbf{q}'$  such that, for any data  $\mathcal{D}$ , we have  $\mathcal{O}, \mathcal{D} \models \mathbf{q}$  iff  $\mathcal{D} \models \mathbf{q}'$  (AC<sup>0</sup>)
- linear-datalog-rewritable (in NL)
- datalog-rewritable (in PTime)

# Our little problem

## Covering axiom

$A \sqsubseteq T \sqcup F$  class  $A$  is covered by classes  $T$  and  $F$

**Example:**  $Animal \sqsubseteq Female \sqcup Male$

## Global covering and disjointness

$\top \sqsubseteq T \sqcup F$  everything is covered by  $T$  and  $F$

**Disjointness:**  $T \cap F \sqsubseteq \perp$        $\top \sqsubseteq Alive \sqcup Dead, Alive \cap Dead \sqsubseteq \perp$

### **Notation:**

$Cov_{\top} = \{T \sqsubseteq F \sqcup T\}$

$Cov_{\perp}^{\dagger} = \{T \sqsubseteq F \sqcup T, F \cap T \sqsubseteq \perp\}$

$Cov_A = \{A \sqsubseteq F \sqcup T\}$

$Cov_A^{\dagger} = \{A \sqsubseteq F \sqcup T, F \cap T \sqsubseteq \perp\}$

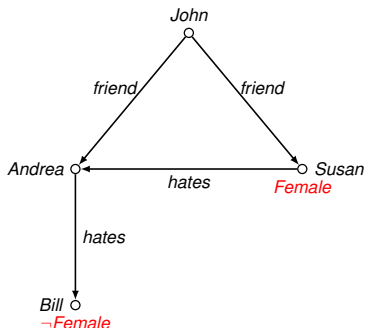
## Problem

Classify syntactically CQs  $q$  w.r.t. complexity of answering  $(Cov, q)$

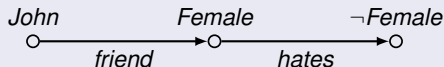
# Andrea example

Consider the ABox  $\mathcal{A}$ :

(John; Susan) : friend  
(John; Andrea) : friend  
(Susan; Andrea) : hates  
(Andrea; Bill) : hates  
Susan: Female  
Bill :  $\neg$ Female



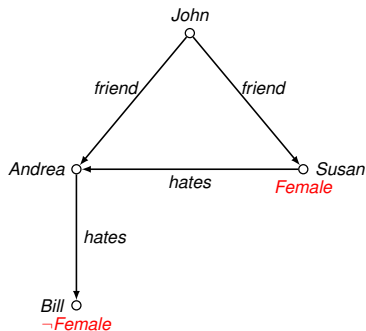
Does John have a female friend who hates a male (not female) person?



# Andrea example

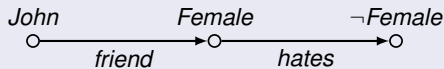
Consider the ABox  $\mathcal{A}$ :

(John; Susan) : friend  
(John; Andrea) : friend  
(Susan; Andrea) : hates  
(Andrea; Bill) : hates  
Susan: Female  
Bill :  $\neg$ Female



$\mathcal{T} = \{T \sqsubseteq \text{Female} \sqcup \text{Male}\}$

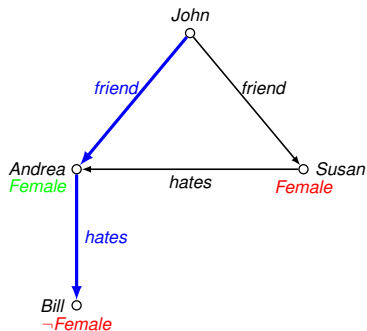
Does John have a female friend who hates a male (not female) person?



# Andrea example

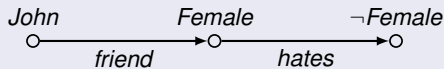
Consider the ABox  $\mathcal{A}$ :

(John; Susan) : friend  
(John; Andrea) : friend  
(Susan; Andrea) : hates  
(Andrea; Bill) : hates  
Susan: Female  
Bill :  $\neg$ Female



$\mathcal{T} = \{T \sqsubseteq \text{Female} \sqcup \text{Male}\}$

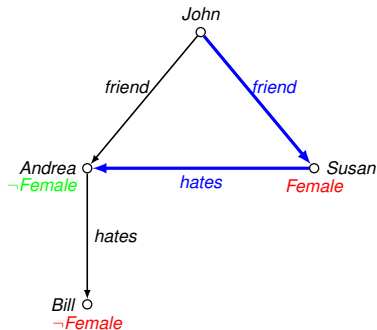
Does John have a female friend who hates a male (not female) person?



# Andrea example

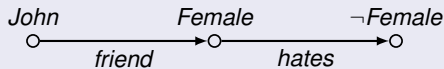
Consider the ABox  $\mathcal{A}$ :

(John; Susan) : friend  
(John; Andrea) : friend  
(Susan; Andrea) : hates  
(Andrea; Bill) : hates  
Susan: Female  
Bill :  $\neg$ Female



$\mathcal{T} = \{T \sqsubseteq \text{Female} \sqcup \text{Male}\}$

Does John have a female friend who hates a male (not female) person?



# Typical examples for $\text{Cov}_A = \{A \sqsubseteq F \sqcup T\}$

Complexity	CQ $q$	Explanation
$AC^0$	$F \circ \longrightarrow \circ$	if $q$ has only $F$ atoms but no $T$ , then the $F$ can be ignored
L	$F \circ \begin{array}{c} \longleftarrow \\ \longrightarrow \end{array} \circ T$	checks undirected reachability: $F \circ \longleftrightarrow \circ \longleftrightarrow \circ \longleftrightarrow \circ T$ the answer to $Q$ is 'yes'
NL	$F \circ \longrightarrow \circ T$	checks directed reachability: $F \circ \longrightarrow \circ \longrightarrow \circ \longrightarrow \circ T$ the answer to $Q$ is 'yes'
P	$T \circ \longrightarrow F \circ \longrightarrow T \circ$	evaluates monotone Boolean circuits $C$
coNP	$F \circ \longrightarrow F \circ \longrightarrow T \circ \longrightarrow T \circ$	checks CNF satisfiability

**Naïve idea:** classify  $q$  according to the number of occurrences of  $T, F$

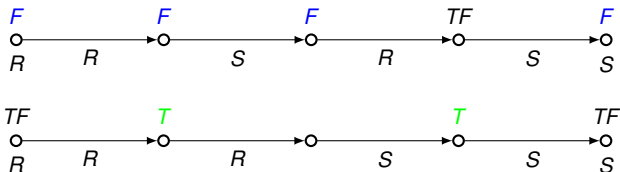
$q$  has NO **solitary occurrences** of  $F$  (or  $T$ )



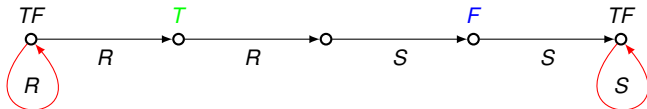
$Q = (\text{Cov}, q)$  is in AC<sup>0</sup>



$q$  is a rewriting of  $Q$



necessary & sufficient condition for **path CQs**, but not in general:



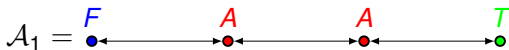


$q$  has **exactly one** atom  $F(x)$  and **exactly one** atom  $T(y)$ , for  $x \neq y$

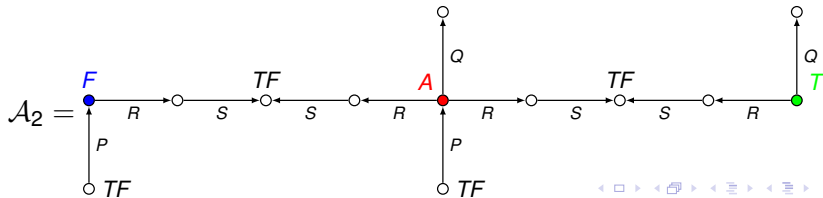


answering  $(\text{Cov}, q)$  is L-hard

- $q_1 = \text{O} \xrightarrow{F} \text{O} \xrightarrow{T} \text{O}$  is a *symmetric* query (L-complete)



- $q_2 = \text{O} \xrightarrow{P} \text{O} \xrightarrow{R} \text{O} \xrightarrow{S} \text{O} \xrightarrow{S} \text{O} \xrightarrow{R} \text{O} \xrightarrow{Q} \text{O}$   
(The middle three nodes are connected by orange arrows labeled R, S, S, R and are grouped under a bracket labeled "symmetric")

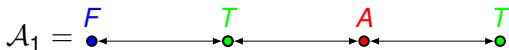


$q$  has **exactly one** atom  $F(x)$  and **exactly one** atom  $T(y)$ , for  $x \neq y$



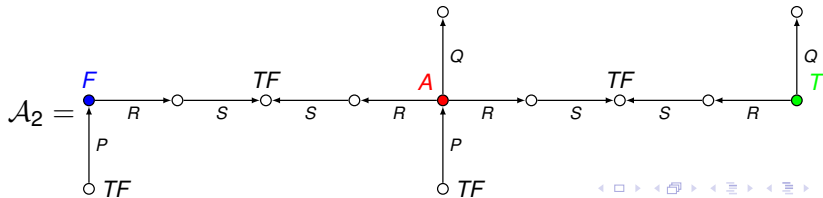
answering  $(\text{Cov}, q)$  is L-hard

•  $q_1 = \text{O} \xrightarrow{F} \text{O} \xrightarrow{T} \text{O}$  is a *symmetric* query (L-complete)



•  $q_2 = \text{O} \xrightarrow{P} \text{O} \xrightarrow{R} \text{O} \xrightarrow{S} \text{O} \xrightarrow{S} \text{O} \xrightarrow{R} \text{O} \xrightarrow{Q} \text{O}$

symmetric

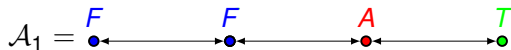


$q$  has **exactly one** atom  $F(x)$  and **exactly one** atom  $T(y)$ , for  $x \neq y$



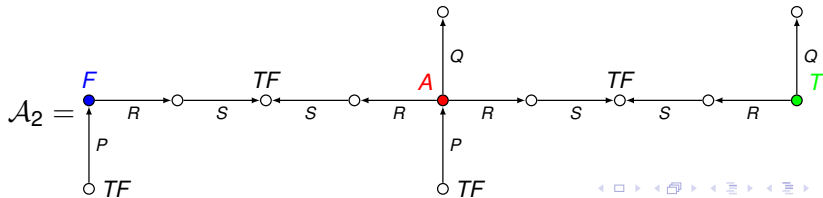
answering  $(\text{Cov}, q)$  is L-hard

•  $q_1 = \text{O} \xrightarrow{F} \text{O} \xrightarrow{T} \text{O}$  is a *symmetric* query (L-complete)



•  $q_2 = \text{O} \xrightarrow{P} \text{O} \xrightarrow{R} \text{O} \xrightarrow{S} \text{O} \xrightarrow{S} \text{O} \xrightarrow{R} \text{O} \xrightarrow{Q} \text{O}$

symmetric

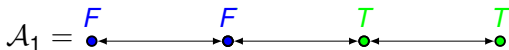


$q$  has **exactly one** atom  $F(x)$  and **exactly one** atom  $T(y)$ , for  $x \neq y$



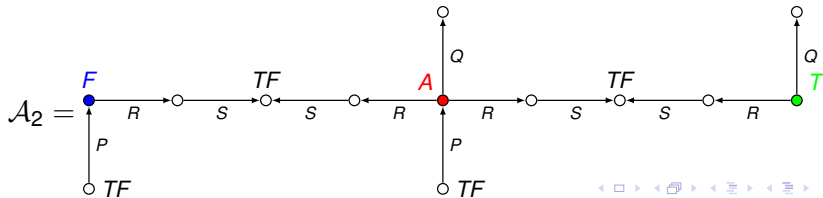
answering  $(\text{Cov}, q)$  is L-hard

•  $q_1 = \text{O} \xrightarrow{F} \text{O} \xrightarrow{T} \text{O}$  is a *symmetric* query (L-complete)



•  $q_2 = \text{O} \xrightarrow{P} \text{O} \xrightarrow{R} \text{O} \xrightarrow{S} \text{O} \xrightarrow{S} \text{O} \xrightarrow{R} \text{O} \xrightarrow{Q} \text{O}$


symmetric




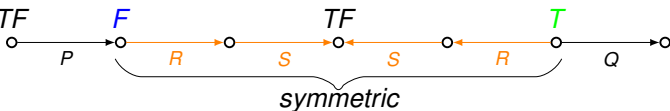
$q$  has **exactly one** atom  $F(x)$  and **exactly one** atom  $T(y)$ , for  $x \neq y$

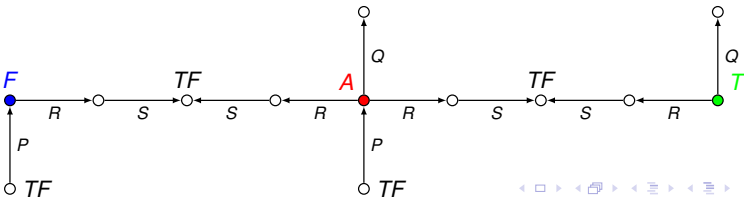


answering  $(\text{Cov}, q)$  is L-hard

•  $q_1 =$   is a *symmetric* query (L-complete)

$\mathcal{A}_1 =$  

•  $q_2 =$  

$\mathcal{A}_2 =$  

*path* CQs

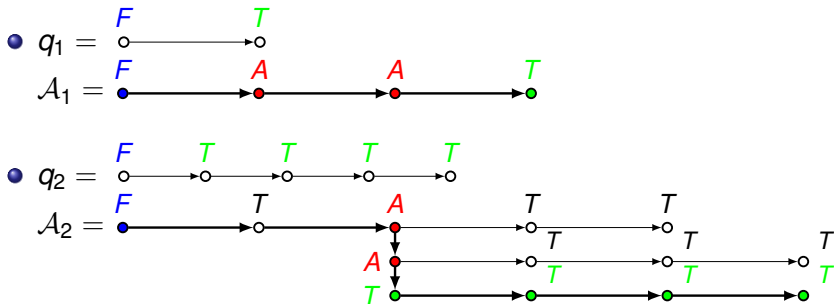
+

**at least one solitary  $F(x)$  and at least one solitary  $T(y)$**

$\Downarrow$

answering  $(\text{Cov}, \mathbf{q})$  is NL-hard

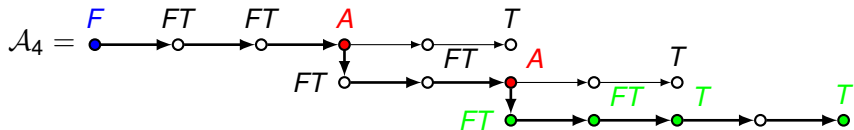
Examples of NL-complete queries:



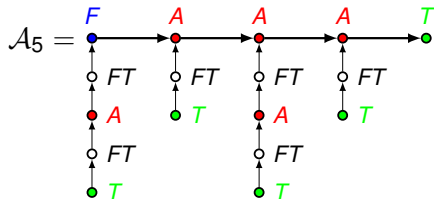
# NL-complete queries

•  $q_3 = \overset{F}{\circ} \rightarrow \circ \rightarrow \overset{T}{\circ} \rightarrow \circ \rightarrow \overset{T}{\circ} \rightarrow \circ \rightarrow \overset{T}{\circ}$  is similar to  $q_2$

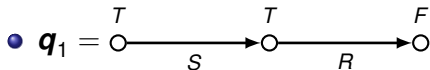
•  $q_4 = \overset{F}{\circ} \xrightarrow{FT} \circ \xrightarrow{FT} \circ \xrightarrow{T} \circ \rightarrow \circ \rightarrow \overset{T}{\circ}$



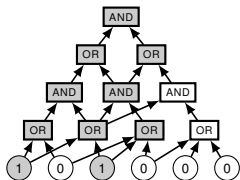
•  $q_5 = \overset{T}{\circ} \xrightarrow{FT} \circ \xrightarrow{F} \circ \xrightarrow{T} \circ$



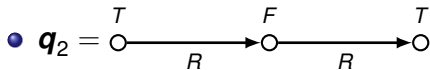
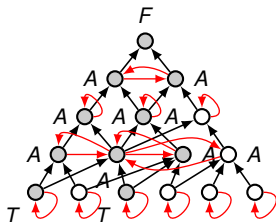
# P-hardness via monotone circuits



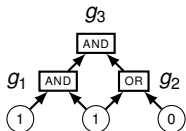
$\mathcal{C}_1$ :



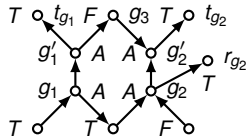
$\mathcal{A}_{\mathcal{C}_1}^\alpha$ :



$\mathcal{C}_2$ :



$\mathcal{A}_{\mathcal{C}_2}^\alpha$ :

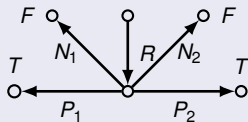


We can see that  $\mathcal{C}(\vec{\alpha}) = 1$  iff  $\text{Cov}, \mathcal{A}_{\mathcal{C}}^\alpha \models q$ .



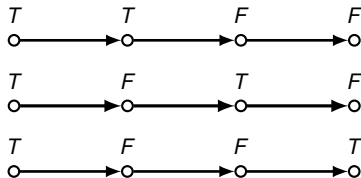
## Schaerf's (1993) example

$(\text{Cov}_T, \mathbf{q})$  with  $\mathbf{q}$  on the right is coNP-complete



## Conjecture

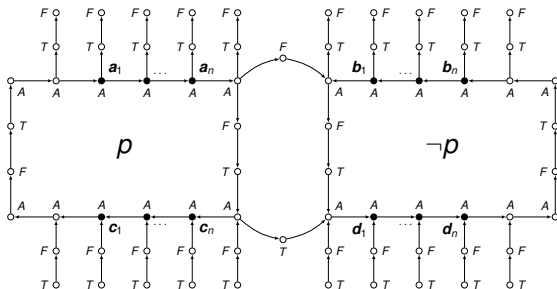
For (tree-shaped)  $\mathbf{q}$  with at least two solitary  $F$  and two solitary  $T$ ,  
 $(\text{Cov}_A, \mathbf{q})$  is coNP-hard



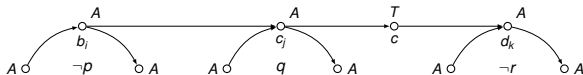
This is so for all 2-2 path-CQ of length 4

# Proving coNP-hardness by reduction to 3CNF

Gadget for  $p$  and  $\neg p$  in a 3CNF with  $n$  clauses:



Gadget for one clause  $c = (p \vee q \vee \neg r)$ :



The resulting ABox  $\mathcal{A}_\psi$  is such that  $\psi$  is satisfiable iff  $\text{Cov}_A, \mathcal{A}_\psi \not\models \mathbf{q}$

# Conclusion

- We have found quite a few syntactic and semantic sufficient and/or necessary conditions for OMQ (Cov; q) to be in this or that complexity class
- We are still aiming at a general complete classification, but the connections to other notoriously hard problems indicate that achieving this aim will be difficult

Thank you for attention!