

A data complexity and rewritability tetrachotomy of ontology-mediated queries with a covering axiom

O. Gerasimova¹ S. Kikot² M. Zakharyashev³
A. Kurucz⁴ V. Podolskii^{1,5}

¹National Research University Higher School of Economics, Moscow, Russia

²London Metropolitan University, U.K.

³Birkbeck, University of London, U.K.

⁴King's College, University of London, U.K.

⁵Steklov Mathematical Institute, Moscow, Russia

17th International Conference on Principles of Knowledge
Representation and Reasoning, Rhodes, 2020

Fact Set (a collection of *facts*, or ground atoms, a.k.a. the data)

$\mathcal{D} = \{Male(John), Parent(John, Sasha), Person(Sasha), Parent(Sasha, Ann), Female(Ann)\}$

Fact Set (a collection of *facts*, or ground atoms, a.k.a. the data)

$\mathcal{D} = \{Male(John), Parent(John, Sasha), Person(Sasha), Parent(Sasha, Ann), Female(Ann)\}$

Ontology (a set of sentences from some fragment of μ FO called *description logic*)

$\mathcal{O} = \{\forall x(Person(x) \rightarrow Male(x) \vee Female(x))\}$

Fact Set (a collection of *facts*, or ground atoms, a.k.a. the data)

$\mathcal{D} = \{Male(John), Parent(John, Sasha), Person(Sasha), Parent(Sasha, Ann), Female(Ann)\}$

Ontology (a set of sentences from some fragment of μ FO called *description logic*)

$\mathcal{O} = \{\forall x(Person(x) \rightarrow Male(x) \vee Female(x))\}$

Query (CQ, i.e. existentially quantified conjunction of atoms)

$q = \exists x \exists y Male(x) \wedge Parent(x, y) \wedge Female(y)$ (Is there a father and a daughter ?)

Fact Set (a collection of *facts*, or ground atoms, a.k.a. the data)

$\mathcal{D} = \{Male(John), Parent(John, Sasha), Person(Sasha), Parent(Sasha, Ann), Female(Ann)\}$

Ontology (a set of sentences from some fragment of μ FO called *description logic*)

$\mathcal{O} = \{\forall x(Person(x) \rightarrow Male(x) \vee Female(x))\}$

Query (CQ, i.e. existentially quantified conjunction of atoms)

$\mathbf{q} = \exists x \exists y Male(x) \wedge Parent(x, y) \wedge Female(y)$ (Is there a father and a daughter ?)

Semantics: $\mathcal{D}, \mathcal{O} \models \mathbf{q}$ if \mathbf{q} holds in all interpretations for \mathcal{D} and \mathcal{O}

$\mathcal{D}_1 = \mathcal{D} \cup \{Male(Sasha)\} \models \mathbf{q}$ under $x \mapsto Sasha$ and $y \mapsto Ann$

$\mathcal{D}_2 = \mathcal{D} \cup \{Female(Sasha)\} \models \mathbf{q}$ under $x \mapsto John$ and $y \mapsto Sasha$, thus $\mathcal{D}, \mathcal{O} \models \mathbf{q}$

Fact Set (a collection of *facts*, or ground atoms, a.k.a. the data)

$\mathcal{D} = \{Male(John), Parent(John, Sasha), Person(Sasha), Parent(Sasha, Ann), Female(Ann)\}$

Ontology (a set of sentences from some fragment of μ FO called *description logic*)

$\mathcal{O} = \{\forall x(Person(x) \rightarrow Male(x) \vee Female(x))\}$

Query (CQ, i.e. existentially quantified conjunction of atoms)

$\mathbf{q} = \exists x \exists y Male(x) \wedge Parent(x, y) \wedge Female(y)$ (Is there a father and a daughter ?)

Semantics: $\mathcal{D}, \mathcal{O} \models \mathbf{q}$ if \mathbf{q} holds in all interpretations for \mathcal{D} and \mathcal{O}

$\mathcal{D}_1 = \mathcal{D} \cup \{Male(Sasha)\} \models \mathbf{q}$ under $x \mapsto Sasha$ and $y \mapsto Ann$

$\mathcal{D}_2 = \mathcal{D} \cup \{Female(Sasha)\} \models \mathbf{q}$ under $x \mapsto John$ and $y \mapsto Sasha$, thus $\mathcal{D}, \mathcal{O} \models \mathbf{q}$

Ontology-mediated query (OMQ) $\mathcal{Q} = (\mathcal{O}, \mathbf{q})$ checks if

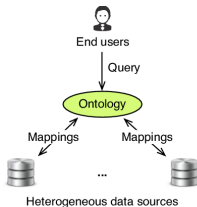
“there is a *Parent*-path over *Person* from *Male* to *Female*” (a rewriting of \mathcal{Q})

- can be expressed in linear Datalog;
- NL-complete in data complexity (other examples match AC^0 , L, P and coNP)

Data complexity of OMQ answering

Fundamental 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}



Fundamental 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 $\mathcal{D} \models \mathbf{q}'$ (AC⁰)
- linear-datalog-rewritable (in NL)
- datalog-rewritable (in PTime)

Boundedness of datalog programs (is recursion eliminable?)

Undecidable ([Vardi 1988](#)), even for sirups ([Marcinkowski 1996](#))
2ExpTime-complete for monadic programs ([Cosmadakis et al. 1988](#))

Linearisability of datalog programs

Transform a datalog program into an equivalent linear program
i.e., ≤ 1 IDB predicate in rule bodies ([Ullman, Afrati, Saraiya, etc.](#))

Markable disjunctive Datalog programs

A sufficient condition for plain Datalog rewritability of disjunctive Datalog programs.

([Grau, Nenov, Kaminsky 2016](#))

Constraint satisfaction problems

Classify CSP patterns according to their complexities ...

Constraint satisfaction problems

Classify CSP patterns according to their complexities ...

Reduction to CSPs

FO- and datalog-rewritability of OMQs with *SHIU* ontologies and atomic queries is decidable in NExpTime by a reduction to CSPs
([Bienvenu, ten Cate, Lutz, Wolter 2014](#))

Constraint satisfaction problems

Classify CSP patterns according to their complexities ...

Reduction to CSPs

FO- and datalog-rewritability of OMQs with *SHIU* ontologies and atomic queries is decidable in NExpTime by a reduction to CSPs
([Bienvenu, ten Cate, Lutz, Wolter 2014](#))

Trichotomy for \mathcal{EL} OMQs

AC⁰/NL/P trichotomy for the data complexity of answering OMQs with an \mathcal{EL} ontology and a CQ
Can be checked in ExpTime
([Lutz, Sabellek 2017-20](#))

Our approach

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:

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

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

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

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

Problem

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

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 Cov_A can be ignored
L	$F \circ \begin{array}{c} \longleftarrow \\ \longleftarrow \end{array} \circ \begin{array}{c} \longrightarrow \\ \longrightarrow \end{array} T$	checks undirected reachability: $F \circ \xrightarrow{A} \circ \xrightarrow{A} \circ \longrightarrow T$ the answer to Q is 'yes'
NL	$F \circ \longrightarrow \circ T$	checks directed reachability: $F \circ \xrightarrow{A} \circ \xrightarrow{A} \circ \longrightarrow T$ the answer to Q is 'yes'
P	$T \circ \longrightarrow \circ F \longrightarrow \circ T$	evaluates monotone Boolean circuits C
coNP	$F \circ \longrightarrow \circ F \longrightarrow \circ T \longrightarrow \circ T$	checks CNF satisfiability

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

How disjunctive predicates F and T may occur in q ?

$F(x) \in q$ is **solitary** if $T(x) \notin q$ $F(x), T(x) \in q$ is a **twin**

Theorem: AC^0 vs. L-hard (emphasise the danger of twins)

- If q has no solitary F then both (Cov_A^\perp, q) and (Cov_A, q) are in AC^0 , with q being an FO-rewriting
- Otherwise, for q without twins, (Cov_\top^\perp, q) , (Cov_\top, q) are L-hard
- (Cov_A^\perp, q) is in AC^0 iff q has no solitary F or T (nice criterion)
- deciding if (Cov_A, q) is in AC^0 is PSPACE-hard if q may have twins (and highly likely 2ExpTime-hard)

Syntactic tetrachotomy for path CQs and Cov_A^\perp

A CQ is a *path-CQ* if its binary atoms can be arranged in a directed path.

Theorem (main result)

For any $\mathcal{Q} = (\text{Cov}_A^\perp, \mathbf{q})$ with path-CQ \mathbf{q} , the following hold:

• if \mathbf{q} is a 0-CQ, then \mathcal{Q} is in AC^0 ;

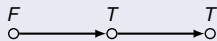
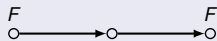
• if \mathbf{q} is a periodic 1-CQ,

then \mathcal{Q} is NL-complete;

• if \mathbf{q} is a non-periodic 1-CQ,

then \mathcal{Q} is P-complete;

• if \mathbf{q} is a 2-CQ, then \mathcal{Q} is coNP-complete.



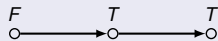
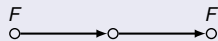
Syntactic tetrachotomy for path CQs and Cov_A^\perp

A CQ is a *path-CQ* if its binary atoms can be arranged in a directed path.

Theorem (main result)

For any $\mathcal{Q} = (\text{Cov}_A^\perp, \mathbf{q})$ with path-CQ \mathbf{q} , the following hold:

- if \mathbf{q} is a 0-CQ, then \mathcal{Q} is in AC^0 ;
- if \mathbf{q} is a periodic 1-CQ,
then \mathcal{Q} is NL-complete;
- if \mathbf{q} is a non-periodic 1-CQ,
then \mathcal{Q} is P-complete;
- if \mathbf{q} is a 2-CQ, then \mathcal{Q} is coNP-complete.



Questions ?

CQs without solitary F

$F(x) \in \mathbf{q}$ is **solitary** if $T(x) \notin \mathbf{q}$ $F(x), T(x) \in \mathbf{q}$ is a **twin**

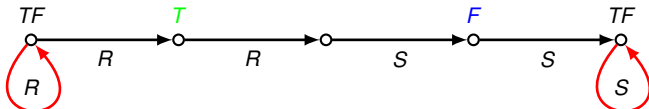
Theorem: AC^0 vs. L-hard

- If \mathbf{q} has no solitary F then both $(Cov_A^\perp, \mathbf{q})$ and (Cov_A, \mathbf{q}) are in AC^0 , with \mathbf{q} being an FO-rewriting
- Otherwise, for \mathbf{q} without twins, (Cov_\perp, \mathbf{q}) , (Cov_T, \mathbf{q}) are L-hard
- $(Cov_A^\perp, \mathbf{q})$ is in AC^0 iff \mathbf{q} has no solitary F or T (nice criterion)

Path CQs

(Cov_A, \mathbf{q}) with a path CQ \mathbf{q} is in AC^0 iff \mathbf{q} has no solitary F or T .
If path \mathbf{q} contains both solitary F and T , then (Cov_A, \mathbf{q}) is NL-hard

However, (Cov_A, \mathbf{q}) with the following (minimal) non-path \mathbf{q} is in AC^0



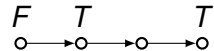
1-CQs (with one solitary F)


- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$

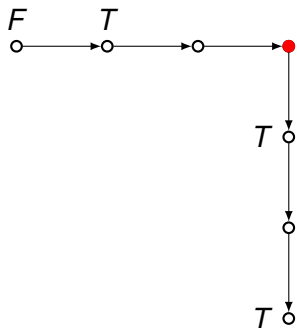




1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$

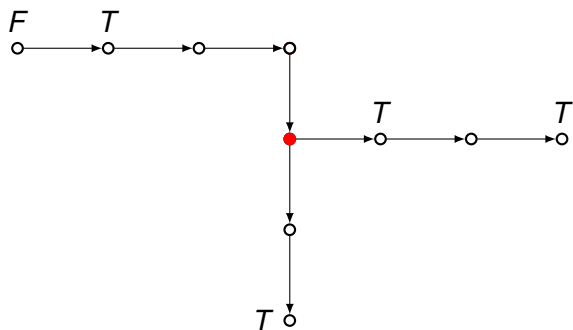


1. budding

1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$

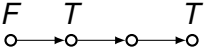


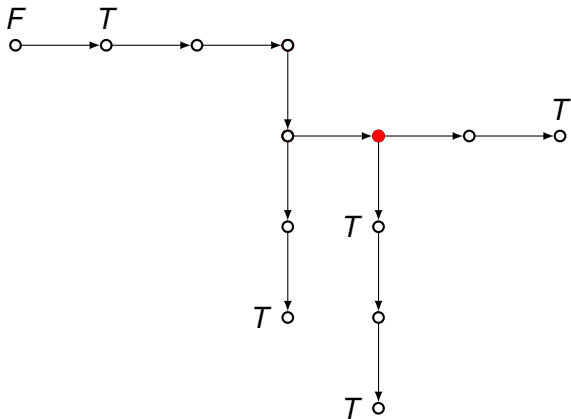
1. budding

2. budding

1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$ 



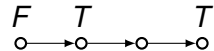
1. budding

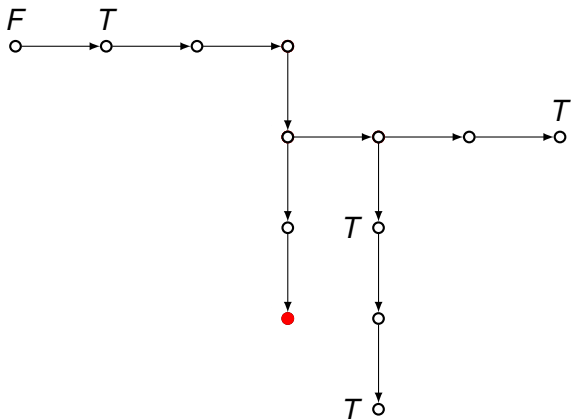
2. budding

3. budding

1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$ 



1. budding

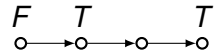
2. budding

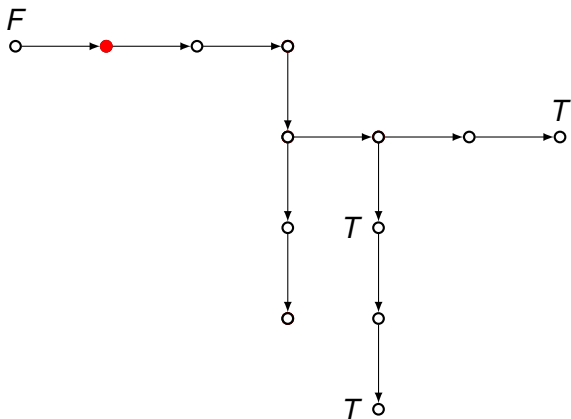
3. budding

4. pruning

1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$ 



1. budding

2. budding

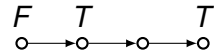
3. budding

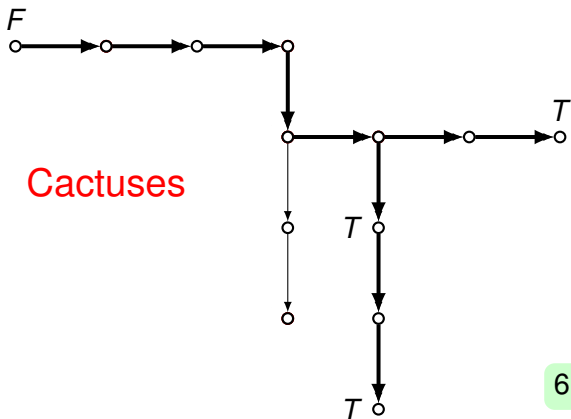
4. pruning

5. pruning

1-CQs (with one solitary F)

- Reducible to datalog, so can be done in P
- For tree-shaped q with root F , reducible to \mathcal{EL} atomic OMQs, so $AC^0/NL/P$ trichotomy decidable in EXPTIME by (Lutz, Sabellek)

Example: Understand models of (Cov_T, q) with $q =$ 



1. budding

2. budding

3. budding

4. pruning

5. pruning

6. branching degree 1

Cactuses (inspired by Cosmodakis et al. 1988, Lutz-Sabellek 2017-18)

For $\mathbf{Q} = (\text{Cov}_A, \mathbf{q})$ with 1-CQ \mathbf{q} , $\mathfrak{R}_{\mathbf{Q}}$ is the set of all cactuses for \mathbf{Q}
 $\mathfrak{R}_{\mathbf{Q}}^{\dagger}$ are minimal cactuses in $\mathfrak{R}_{\mathbf{Q}}$ (without proper sub-cactuses in $\mathfrak{R}_{\mathbf{Q}}$).

Cactus theorem for $\mathbf{Q} = (\text{Cov}_A, \mathbf{q})$ with 1-CQ \mathbf{q}

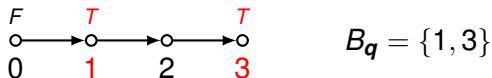
\mathbf{Q} is in AC^0 iff for every $\mathbf{C} \in \mathfrak{R}_{\mathbf{Q}}^{\dagger}$, there is a homomorphism $h: \mathbf{q} \rightarrow \mathbf{C}$;
 \mathbf{Q} is linear-datalog-rewritable (so in NL) if $\mathfrak{R}_{\mathbf{Q}}^{\dagger}$ is boundedly branching.

Proof: represent cactuses as words in a tree alphabet and construct
a tree automaton $\mathfrak{A}_{\mathbf{Q}}$ such that $\mathfrak{R}_{\mathbf{Q}}^{\dagger} \subseteq L(\mathfrak{A}_{\mathbf{Q}}) \subseteq \mathfrak{R}_{\mathbf{Q}}$.
For boundedly branching $\mathfrak{R}_{\mathbf{Q}}^{\dagger}$, $\mathfrak{A}_{\mathbf{Q}}$ is converted to a
linear-datalog-rewriting of \mathbf{Q}

Arithmetic progression theorem

Let \mathbf{q} be an F -path without twins with a single binary relation R .

Let $B_{\mathbf{q}}$ be the set of positions with T :



Theorem (NL/P dichotomy)

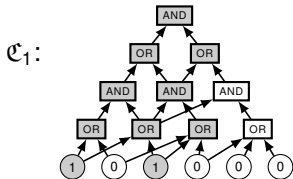
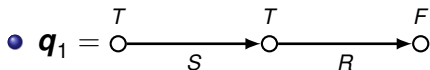
$\mathbf{Q} = (\text{Cov}_A, \mathbf{q})$ is in NL iff the set of integers $\{0\} \cup B_{\mathbf{q}}$ is an arithmetic progression (unless $NL = P$). Otherwise \mathbf{Q} is P-hard.

NB: For \mathbf{q} above: $(\text{Cov}_A, \mathbf{q})$ is P-hard as 0,1,3 is not a progression but $(\text{Cov}_T, \mathbf{q})$ is in NL (**A is essential!**)

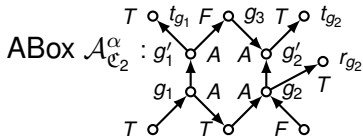
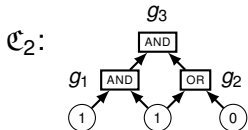
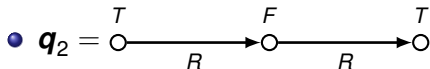
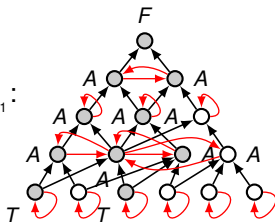
Question

How to generalise this theorem to Cov_T ?

P-hardness via monotone circuits



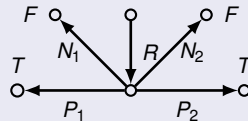
ABox $\mathcal{A}_{\mathcal{C}_1}^\alpha$:



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

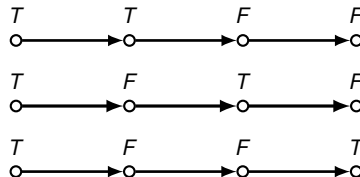
Schaerf's (1993) example

$(\text{Cov}_T, \mathbf{q})$ with \mathbf{q} on the right is coNP-complete
proof by reduction of 2+2 CNF



Conjecture

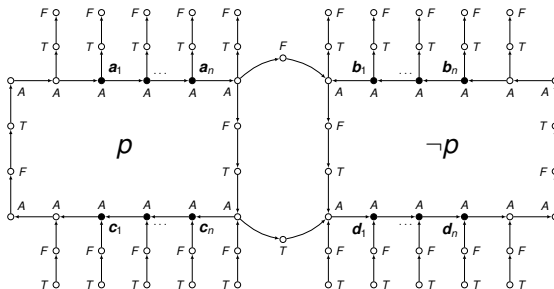
For (tree-shaped) \mathbf{q} with at least two solitary F and two solitary T ,
 $(\text{Cov}_A, \mathbf{q})$ is coNP-hard



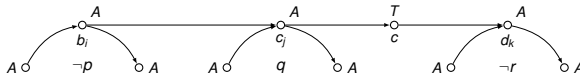
This is so for all 2-2 path-CQ of length 4

Proving coNP-hardness by reduction to 3CNF

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



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



The resulting ABox \mathcal{A}_ψ is s.t. 3CNF ψ is satisfiable iff $\text{Cov}_A, \mathcal{A}_\psi \not\models \mathbf{q}$