

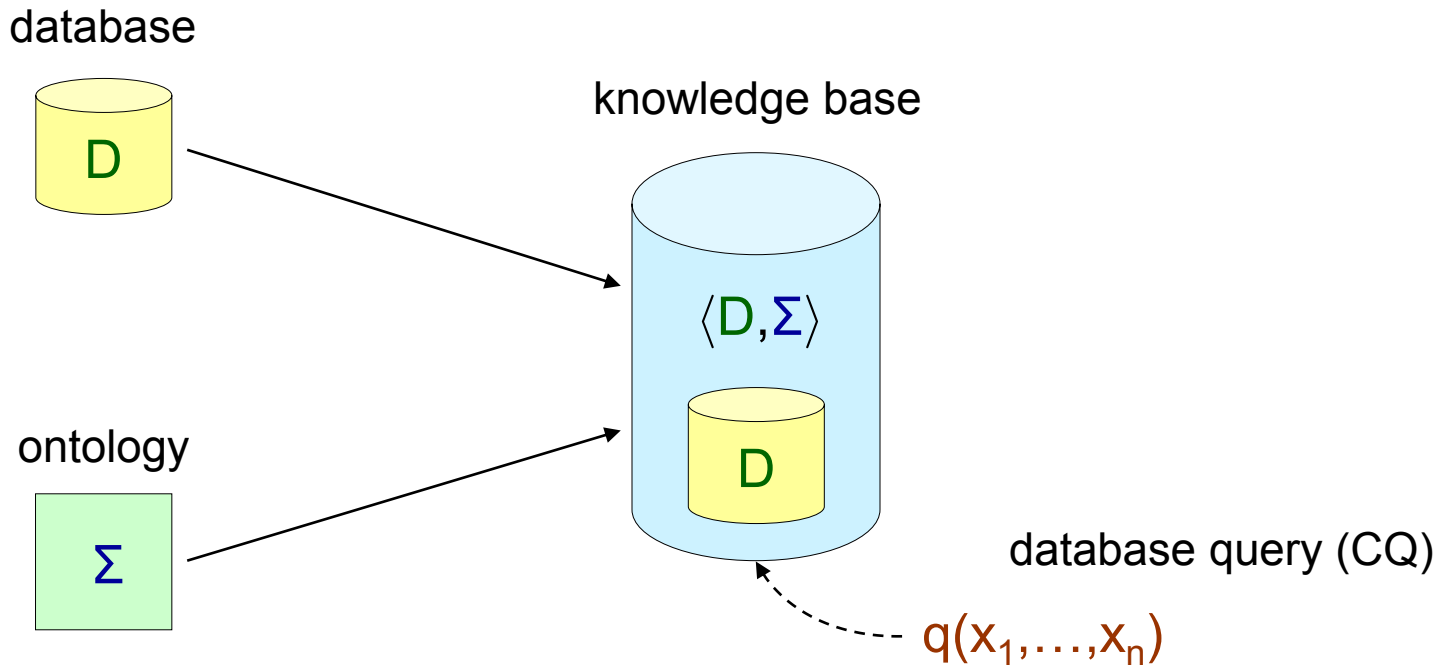
On the Succinctness of Query Rewriting for Datalog[±]

Shqiponja Ahmetaj¹ Andreas Pieris²

¹Institute of Logic and Computation, TU Wien, Austria

²School of Informatics, University of Edinburgh, UK

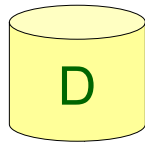
Ontology-Based Query Answering



$$\text{Certain-Answers}(q, D, \Sigma) = \{ (c_1, \dots, c_n) \in \text{dom}(D)^n \mid D \wedge \Sigma \models q(c_1, \dots, c_n) \}$$

Ontology-Mediated Queries

database

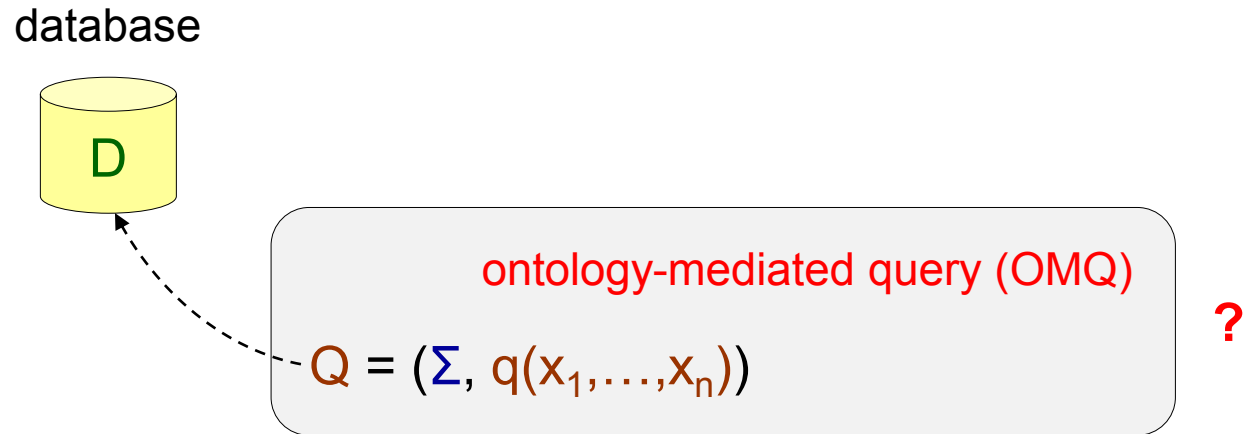


ontology-mediated query (OMQ)

$$Q = (\Sigma, q(x_1, \dots, x_n))$$

$$Q(D) = \text{Certain-Answers}(q, D, \Sigma)$$

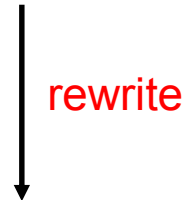
Scalable OMQ Evaluation



Exploit standard RDBMSs - efficient technology for answering queries

Query Rewriting

$$Q = (\Sigma, q(x_1, \dots, x_n))$$



$$Q_{\text{rew}}(x_1, \dots, x_n)$$

a query that can be executed by exploiting existing database technology

$$\text{for every database } D : Q(D) = Q_{\text{rew}}(D)$$

Query Rewriting: An Example

$\{ \forall x (\text{Person}(x) \rightarrow \exists y \text{HasFather}(x,y) \wedge \text{Person}(y)) \}$

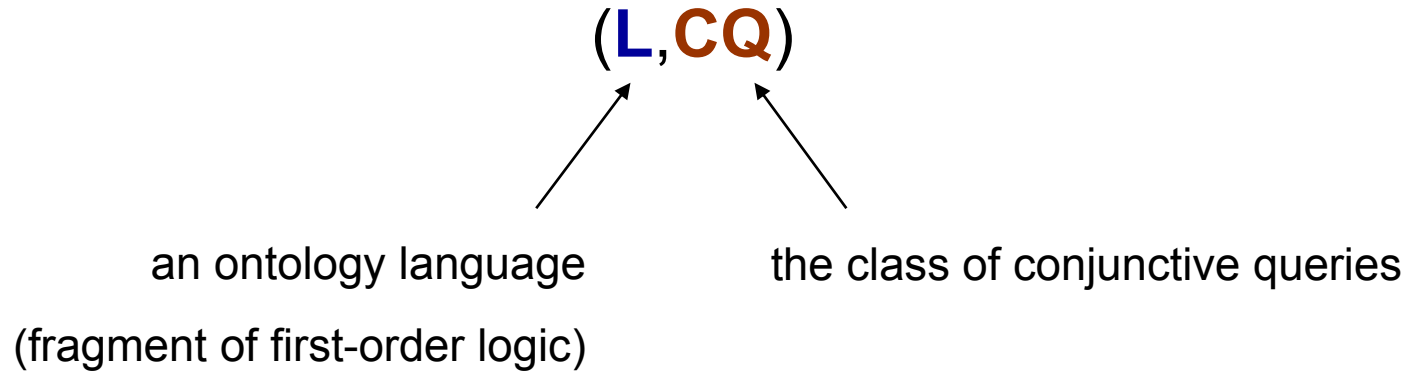
$\exists x (\text{Person}(x) \wedge \text{HasFather}(\text{John},x))$

$Q = (\Sigma, q)$

rewrite

$Q_{\text{rew}} = \exists x (\text{Person}(x) \wedge \text{HasFather}(\text{John},x)) \vee \text{Person}(\text{John})$

Query Rewritability



Definition: An OMQ language O is **QL-Rewritable** if every $Q \in O$ is QL-Rewritable

First-order (FO), $\exists FO^+$, Non-recursive Datalog (NDL), UCQ or Datalog

Query Rewritability: The Main Questions

1. **Can we isolate meaningful OMQ languages that are QL-Rewritable?**
2. **What is the price of QL-rewriting?**

...have been extensively studied for DL- and rule-based OMQ languages

Tuple-Generating Dependencies (TGDs)

(a.k.a. existential rules or Datalog[±] rules)

$$\forall \mathbf{x} \forall \mathbf{y} (\varphi(\mathbf{x}, \mathbf{y}) \rightarrow \exists \mathbf{z} \psi(\mathbf{x}, \mathbf{z}))$$

(**TGD**, **CQ**)

The Guarded Family

Weakly-Guarded

one body-atom contains all
the harmful \forall -variables



Guarded

one body-atom contains
all the \forall -variables



Linear

one body-atom

$$R(\mathbf{w}), \varphi(\mathbf{x}, \mathbf{y}) \rightarrow \exists \mathbf{z} \psi(\mathbf{x}, \mathbf{z}) - \mathbf{w} \subseteq \{\mathbf{x}, \mathbf{y}\}$$

[Cali, Gottlob & Kifer, KR 2008, J. Artif. Intell. Res. 2013]

$$R(\mathbf{x}, \mathbf{y}), \varphi(\mathbf{x}, \mathbf{y}) \rightarrow \exists \mathbf{z} \psi(\mathbf{x}, \mathbf{z})$$

[Cali, Gottlob & Kifer, KR 2008, J. Artif. Intell. Res. 2013]

$$R(\mathbf{x}, \mathbf{y}) \rightarrow \exists \mathbf{z} \psi(\mathbf{x}, \mathbf{z})$$

[Cali, Gottlob & Lukasiewicz, PODS 2009, J. Web Sem. 2012]

The Guarded Family

Theorem: It holds that

1. (**Linear**, **CQ**) is UCQ-Rewritable
2. (**Guarded**, **CQ**) is not FO-Rewritable, but is Datalog-Rewritable
3. (**Weakly-Guarded**, **CQ**) is not Datalog-Rewritable

(Guarded, CQ) is not FO-Rewritable

$$Q = (\{R(x,y), P(y) \rightarrow P(x)\}, P(c_n))$$

$D \supseteq \{P(c_1)\}$, and contains no other P-atom

Q_{rew} has to check for the existence of an R-path in D of **unbounded** length



compute the **transitive closure** of R - not possible via a first-order query

(**Weakly-Guarded**, **CQ**) is not Datalog-Rewritable

Evaluation of (**Weakly-Guarded**, **CQ**) queries is

EXPTIME-complete in data complexity

...in fact, (**Weakly-Guarded**^{stratified}, **CQ**) = EXPTIME, even w/o an order

[Gottlob, Rudolph & Šimkus, **PODS 2014**]

The Guarded Family

Theorem: It holds that

1. (**Linear**, **CQ**) is UCQ-Rewritable
2. (**Guarded**, **CQ**) is not FO-Rewritable, but is Datalog-Rewritable
3. (**Weakly-Guarded**, **CQ**) is not Datalog-Rewritable

The Guarded Family

Theorem: It holds that

1. (**Linear**, **CQ**) is UCQ-Rewritable
2. (**Guarded**, **CQ**) is not FO-Rewritable, but is Datalog-Rewritable
3. (**Weakly-Guarded**, **CQ**) is not Datalog-Rewritable

(Linear, CQ) is UCQ-Rewritable

Via a resolution-based algorithm - XRewrite

ALGORITHM 1: The algorithm XRewrite

Input: a CQ q over a schema \mathcal{R} and a set Σ of TGDs over \mathcal{R}

Output: the perfect rewriting of q w.r.t. Σ

```
 $i := 0;$   
 $Q_{\text{REW}} := \{(q, r, u)\};$   
repeat  
   $Q_{\text{TEMP}} := Q_{\text{REW}};$   
  foreach  $\langle q, x, u \rangle \in Q_{\text{TEMP}}$ , where  $x \in \{r, f\}$  do  
    foreach  $\sigma \in \Sigma$  do applicability condition for TGDs  
      /* rewriting step  
      foreach  $S \subseteq \text{body}(q)$  such that  $\sigma$  is applicable to  $S$  do  
         $i := i + 1;$   
         $q' := \gamma_{S, \sigma^i}(q[S/\text{body}(\sigma^i)]);$   
        if there is no  $\langle q'', r, \star \rangle \in Q_{\text{REW}}$  such that  $q' \simeq q''$  then  
           $Q_{\text{REW}} := Q_{\text{REW}} \cup \{(q', r, u)\};$   
        end  
      end  
    end apply useful reduction steps, but only useful ones  
    /* factorization step  
    foreach  $S \subseteq \text{body}(q)$  which is factorizable w.r.t.  $\sigma$  do  
       $q' := \gamma_S(q);$   
      if there is no  $\langle q'', \star, \star \rangle \in Q_{\text{REW}}$  such that  $q' \simeq q''$  then  
         $Q_{\text{REW}} := Q_{\text{REW}} \cup \{(q', f, u)\};$   
      end  
    end  
  end  
  /* query  $q$  is now explored  
   $Q_{\text{REW}} := (Q_{\text{REW}} \setminus \{(q, x, u)\}) \cup \{(q, x, e)\};$   
end  
until  $Q_{\text{TEMP}} = Q_{\text{REW}};$   
 $Q_{\text{FIN}} := \{q \mid \langle q, r, e \rangle \in Q_{\text{REW}}\};$   
return  $Q_{\text{FIN}}$ 
```

(Linear, CQ) is UCQ-Rewritable

Via a resolution-based algorithm - XRewrite

Given an OMQ $Q = (\Sigma, q)$ from (Linear, CQ)

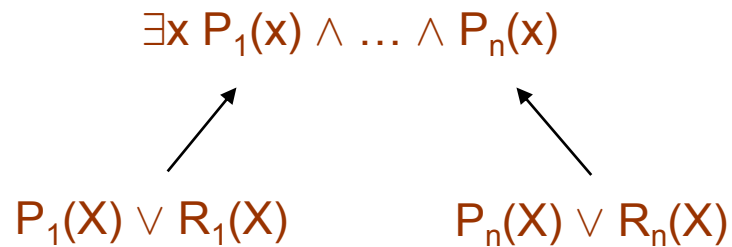
1. The height of $XRewrite(Q)$ is at most $|q|$
2. The size of $XRewrite(Q)$ is at most $\#pred(\Sigma)^{|q|} \cdot (arity(\Sigma) \cdot |q|)^{arity(\Sigma) \cdot |q|}$

worst-case optimal

Lower Bound for (**Linear**, **CQ**)

$$\Sigma = \{R_i(x) \rightarrow P_i(x)\}_{i \in \{1, \dots, n\}}$$

$$q = \exists x (P_1(x) \wedge \dots \wedge P_n(x))$$



\Rightarrow we need to consider 2^n disjuncts

Target More Succinct Query Languages

Theorem: For (**Linear**, **CQ**) there is

- No \exists FO⁺/NDL-rewriting of polynomial size
- No FO-rewriting of polynomial size (unless the PH collapses)

Proof: Via succinctness of monotone Boolean circuits

NOTE: The above proof exploits databases with a single domain element

Two Domain Elements

$$Q = (\Sigma, q(x_1, \dots, x_n))$$



rewrite in polynomial time

$$Q_{\text{rew}}(x_1, \dots, x_n)$$

for every database $D_{01} : Q(D_{01}) = Q_{\text{rew}}(D_{01})$

$\supseteq \{ \text{Zero}(0), \text{One}(1) \}$

Polynomial Rewritings

... assuming two domain elements

Theorem: A (**Linear**, **CQ**) query can be rewritten in polynomial time as:

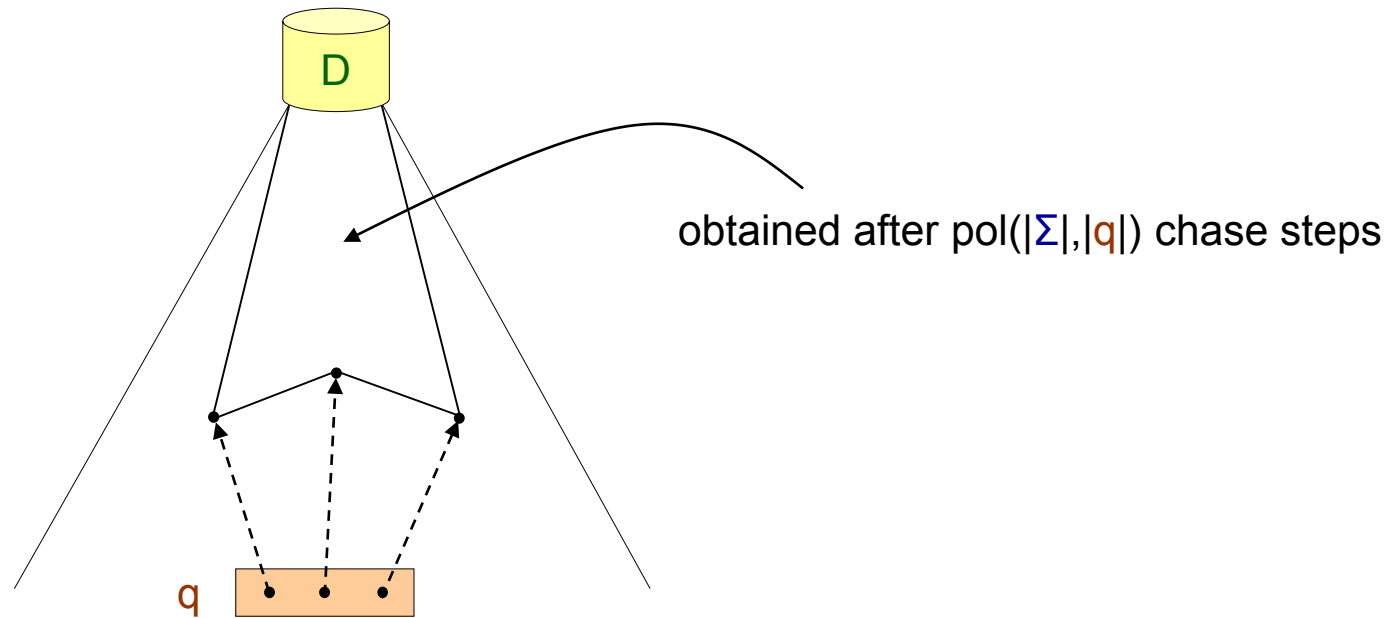
- An $\exists\text{FO}^+$ /NDL query for bounded arity predicates
- An FO query for arbitrary signatures

Proof:

- Bounded arity signatures - via the polynomial witness property
- Arbitrary signatures - via proof generators

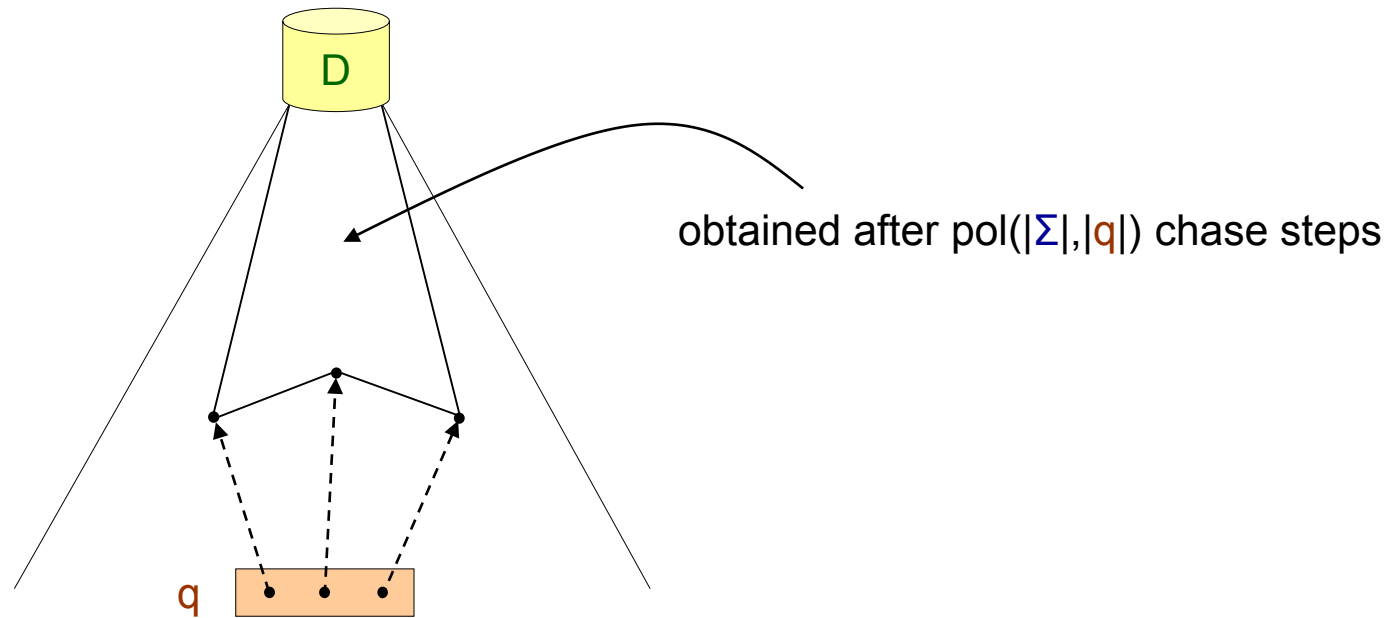
Polynomial Witness Property (PWP)

Definition: $(\mathbf{L}, \mathbf{CQ})$ enjoys the PWP if: there exists a polynomial $\text{pol}(\cdot)$ such that
for every $Q = (\Sigma, q(\mathbf{x})) \in (\mathbf{L}, \mathbf{CQ})$, database D , and $\mathbf{t} \in \text{dom}(D)^{|\mathbf{x}|}$
 $\mathbf{t} \in Q(D) \Rightarrow q(\mathbf{t})$ can be entailed after $\text{pol}(|\Sigma|, |q|)$ chase steps



Polynomial Witness Property (PWP)

Theorem: PWP \Rightarrow \exists FO⁺/NDL-rewritings constructible in polynomial time,
focusing on databases with at least two constants



Witnesses and Linearity

$$0 = \{ \}$$

$$q = \exists z \exists o \text{Number}(o, \dots, o, z, o)$$

$$\{ \text{Number}(x_1, \dots, x_{n-i}, \underbrace{z, o, \dots, o, z, o}_{i-1}) \rightarrow \text{Number}(x_1, \dots, x_{n-i}, \underbrace{o, z, \dots, z, z, o}_{i-1}) \}_{i \in \{1, \dots, n\}}$$

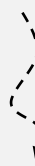
0 Number(0, ..., 0, 0, 0, 1)



1 Number(0, ..., 0, 1, 0, 1)

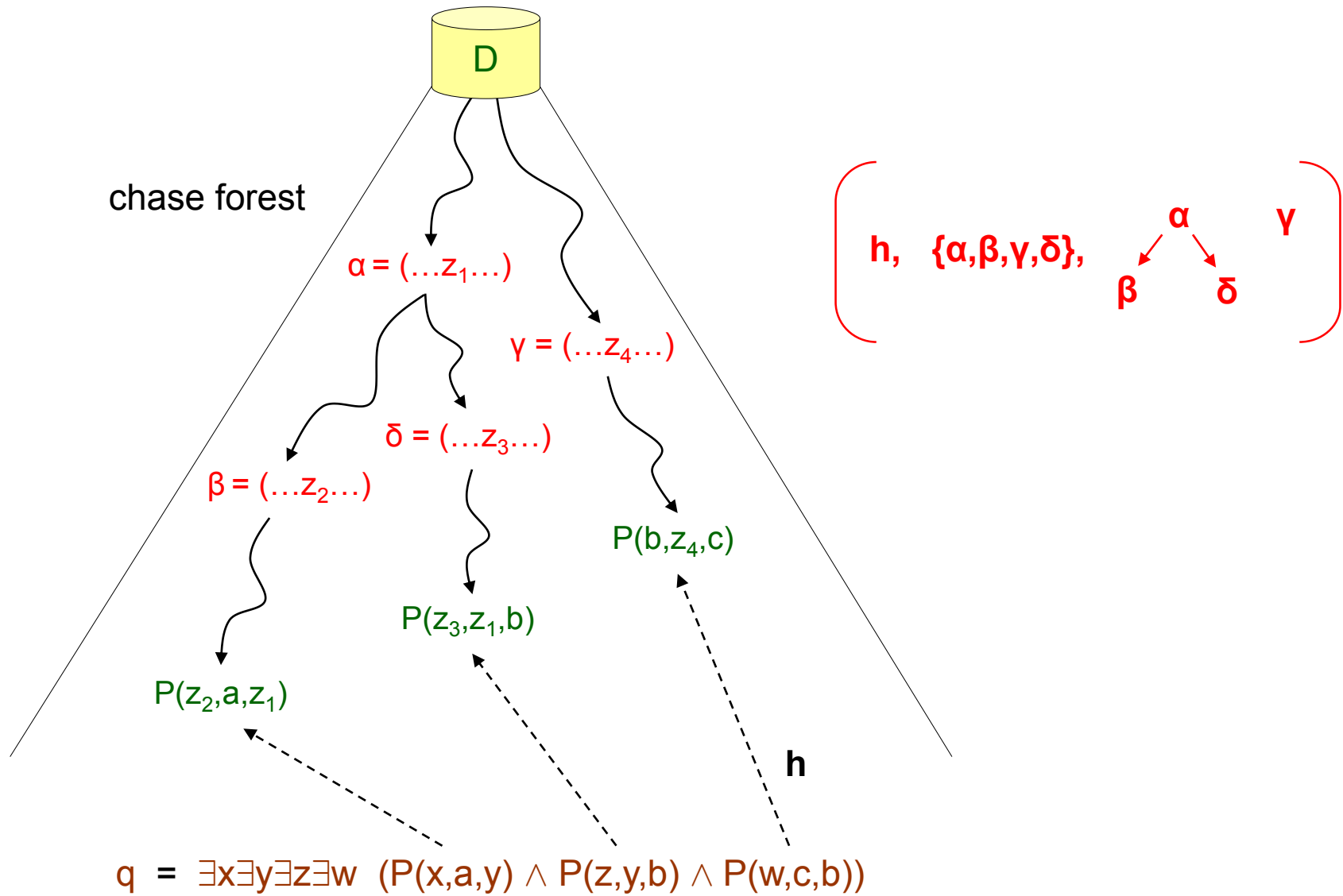


2 Number(0, ..., 1, 0, 0, 1)

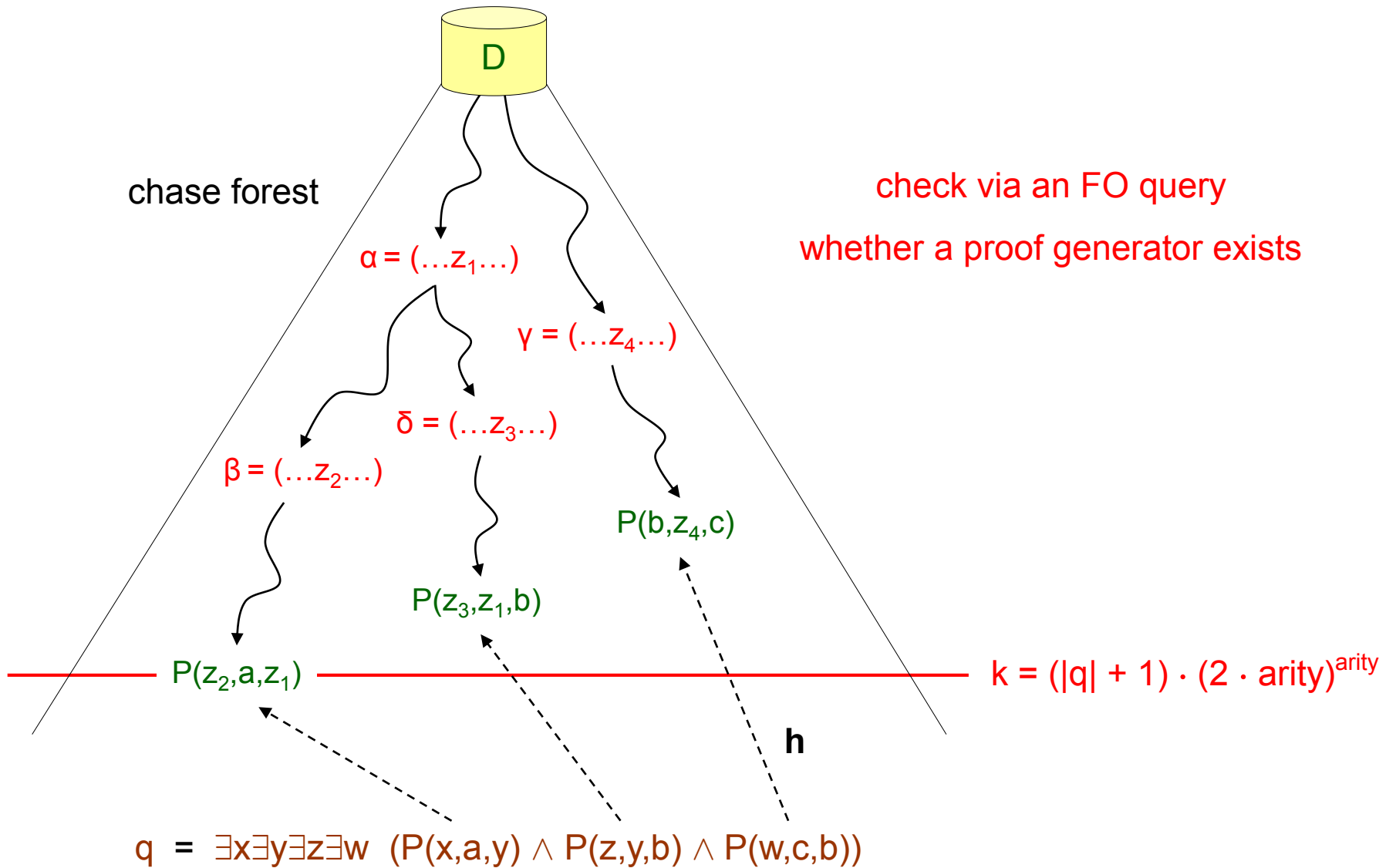


2^n Number(1, ..., 1, 1, 0, 1)

Proof Generator



Proof Generator



Polynomial Rewritings

... assuming two domain elements

Theorem: A (**Linear**, **CQ**) query can be rewritten in polynomial time as:

- An $\exists\text{FO}^+$ /NDL query for bounded arity predicates
- An FO query for arbitrary signatures

Proof:

- Bounded arity signatures - via the polynomial witness property
- Arbitrary signatures - via proof generators

The Guarded Family

Theorem: It holds that

1. (Linear, CQ) is UCQ-Rewritable
2. (Guarded, CQ) is not FO-Rewritable, but is Datalog-Rewritable
3. (Weakly-Guarded, CQ) is not Datalog-Rewritable

(Guarded, CQ) is Datalog-Rewritable

Via inference rules - inspired by DLs

$$\frac{\alpha \rightarrow \beta \wedge A}{\alpha \rightarrow A} \quad A \text{ has no existential variables}$$

$$\frac{\alpha \rightarrow \beta \quad \gamma_1 \wedge \gamma_2 \rightarrow \delta}{\alpha \wedge h(\gamma_1) \rightarrow \beta \wedge h(\delta)} \quad \begin{array}{l} \gamma_1 \wedge \gamma_2 \rightarrow \delta \text{ is a Datalog rule,} \\ h \text{ is a homomorphism from} \\ \gamma_2 \text{ to } \beta \text{ with } \text{vars}(h(\gamma_1)) \subseteq \\ \text{vars}(\alpha). \end{array}$$

$$\frac{\alpha \rightarrow \beta}{g(\alpha) \rightarrow g(\beta)} \quad g : \text{vars}(\alpha) \rightarrow \text{vars}(\alpha)$$

(Guarded, CQ) is Datalog-Rewritable

Via inference rules - inspired by DLs

Given an OMQ $Q = (\Sigma, q)$ from (Guarded, CQ),

the size of the Datalog rewriting is at most $2^{(\#\text{pred}(\Sigma) \cdot \#\text{body-vars}(\Sigma)^{\text{arity}(\Sigma)})}$

worst-case optimal?

Polynomial Rewritings

... assuming two domain elements

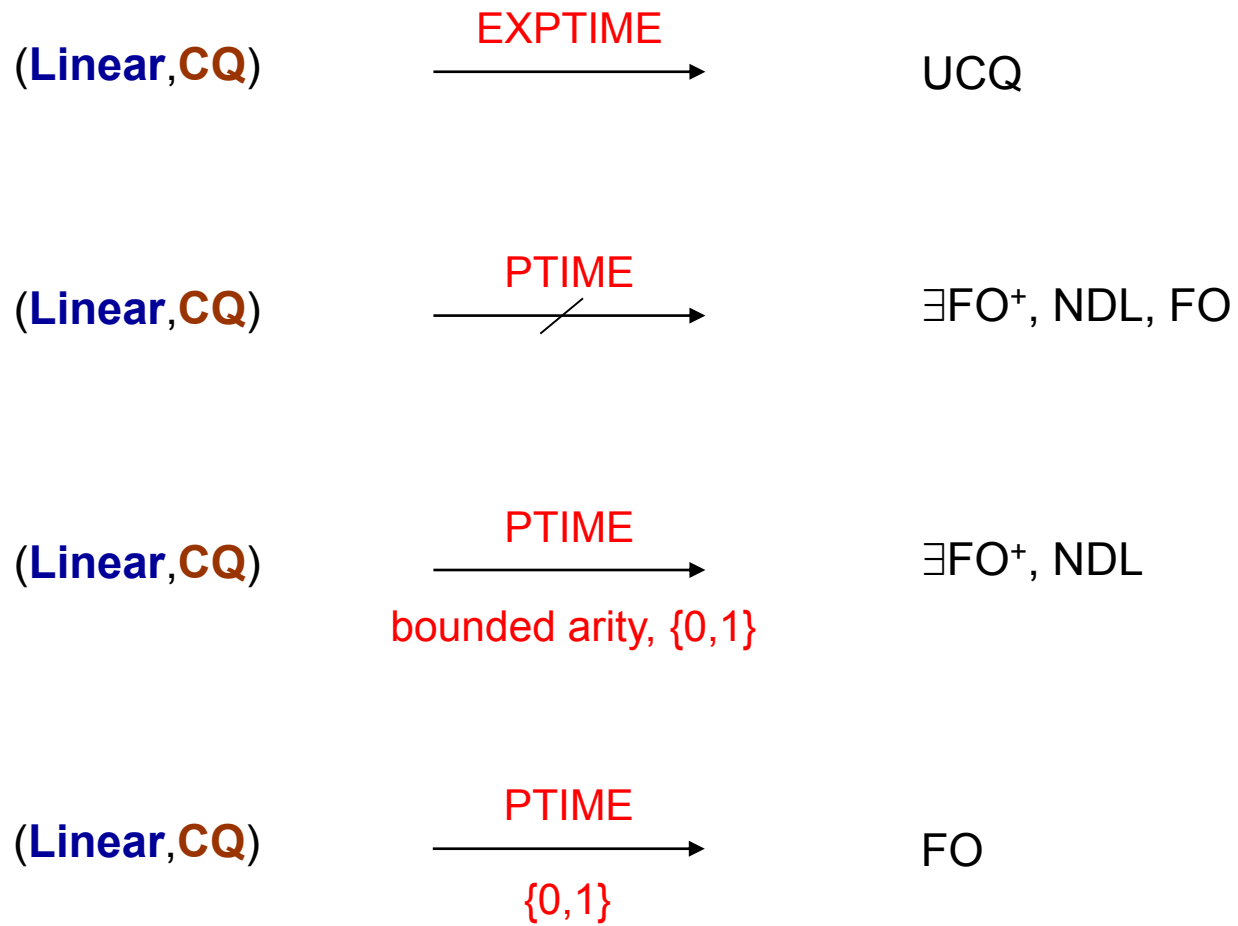
Theorem: A (**Guarded**, **Full CQ**) or (**Guarded**, **Acyclic CQ**) query over bounded arity predicates can be rewritten in polynomial time as a Datalog query

Proof: Via types

- Build all possible types
- Mark “bad” types
- From marked types to Datalog rules that capture ground consequences

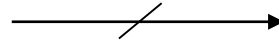
Tuesday, March 27, ICDT4: Logic and Dependencies

Wrap Up



Wrap Up

(Guarded, CQ)



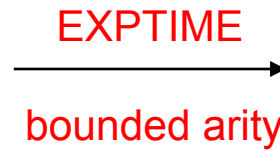
FO

(Guarded, CQ)



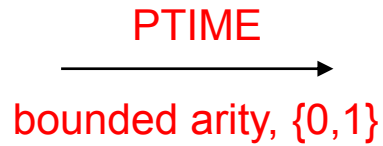
Datalog

(Guarded, CQ)



Datalog

(Guarded, FCQ/ACQ)



Datalog

Open Problems

(**Linear**, **CQ**) $\xrightarrow[\text{bounded arity, } \{0,1\}]{\text{PTIME}}$ UCQ

(**Linear**, **CQ**) $\xrightarrow[\{0,1\}]{\text{PTIME}}$ NDJ

(**Guarded**, **CQ**) $\xrightarrow[\text{bounded arity, } \{0,1\}]{\text{PTIME}}$ Datalog

Thank You!