

REWRITING CONJUNCTIVE QUERIES UNDER DL CONSTRAINTS

Héctor Pérez-Urbina, Boris Motik, and Ian Horrocks

Computing Laboratory
University of Oxford

International Workshop on
Logic in Databases
May 2008





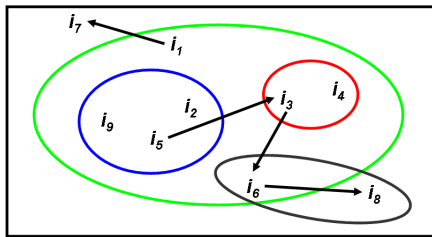
OUTLINE

- Introduction
 - Description Logics (DLs)
 - DLs and DBs
- Answering queries via query rewriting
- Query rewriting for \mathcal{ELHI}
- Complexity results
- Conclusion
- Future work



THE DESCRIPTION LOGIC WORLD

- Things (**individuals**), relationships between things (**roles**) and sets of things (**concepts**)



A DL **Knowledge Base** (KB) is a tuple $\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$

- **TBox** \approx Conceptual **schema**
- **ABox** \approx (Partial) Database **instance**

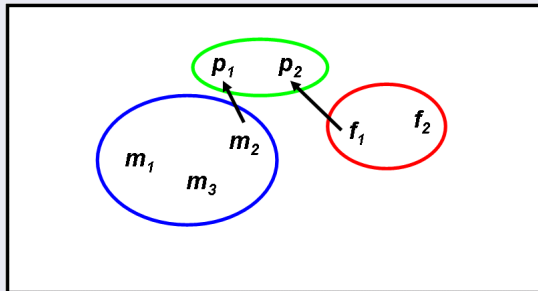


DL KNOWLEDGE BASES

EXAMPLE

$\mathcal{T} = \{\text{Parent} \sqsubseteq \exists \text{hasChild}, \text{Mother} \sqsubseteq \text{Parent}, \text{Father} \sqsubseteq \text{Parent}\}$

$\mathcal{A} = \{\text{Parent}(p_1), \text{Parent}(p_2), \text{Father}(f_1), \text{Father}(f_2), \text{Mother}(m_1),$
 $\text{Mother}(m_2), \text{Mother}(m_3), \text{hasChild}(m_2, p_1), \text{hasChild}(f_1, p_2)\}$



- Parent
- Mother
- Father

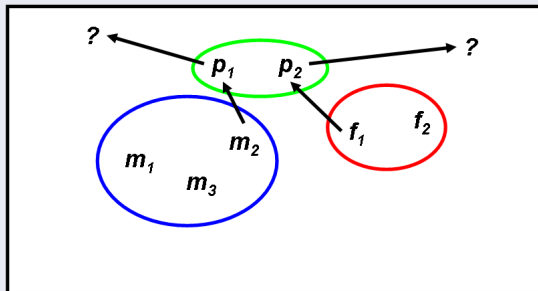


DL KNOWLEDGE BASES

EXAMPLE

$\mathcal{T} = \{\text{Parent} \sqsubseteq \exists \text{hasChild}, \text{Mother} \sqsubseteq \text{Parent}, \text{Father} \sqsubseteq \text{Parent}\}$

$\mathcal{A} = \{\text{Parent}(p_1), \text{Parent}(p_2), \text{Father}(f_1), \text{Father}(f_2), \text{Mother}(m_1),$
 $\text{Mother}(m_2), \text{Mother}(m_3), \text{hasChild}(m_2, p_1), \text{hasChild}(f_1, p_2)\}$



- Parent
- Mother
- Father

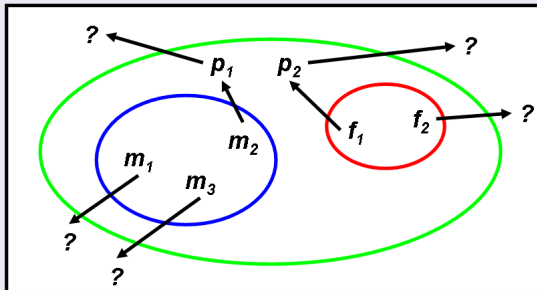


DL KNOWLEDGE BASES

EXAMPLE

$\mathcal{T} = \{\text{Parent} \sqsubseteq \exists \text{hasChild}, \text{Mother} \sqsubseteq \text{Parent}, \text{Father} \sqsubseteq \text{Parent}\}$

$\mathcal{A} = \{\text{Parent}(p_1), \text{Parent}(p_2), \text{Father}(f_1), \text{Father}(f_2), \text{Mother}(m_1),$
 $\text{Mother}(m_2), \text{Mother}(m_3), \text{hasChild}(m_2, p_1), \text{hasChild}(f_1, p_2)\}$



- Parent
- Mother
- Father



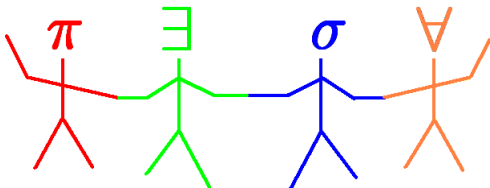
DESCRIPTION LOGIC APPLICATIONS

- DLs underpin the W3C standard **OWL** ontology languages
- Conceptual modeling
 - Systematized Nomenclature of Medicine (**Snomed**)
 - National Cancer Institute (**NCI**)
 - Generalized Architecture for Languages, Encyclopaedias and Nomenclatures in medicine (**Galen**)
- **Semantic Web**
- Information **Integration**
- Data **exchange**
- Data **warehousing**



DESCRIPTION LOGICS AND DATABASES

- TBox \approx Conceptual **schema** (set of dependencies)
- ABox \approx (Partial) Database **instance**
- Query answering (QA) over DL KBs \equiv QA over **incomplete** DBs
 - First-order entailment!
 - $\text{ans}(\phi, \mathcal{K}) = \{\vec{a} \mid \mathcal{K} \models \phi[\vec{a}_1/x_1, \dots, \vec{a}_n/x_n]\}$





QUERYING A KNOWLEDGE BASE

EXAMPLE

$$\mathcal{T} = \{\exists\text{hasParent.Human} \sqsubseteq \text{Human}, \text{Child} \sqsubseteq \exists\text{hasParent.Human}, \\ \text{Father} \sqsubseteq \text{Human}, \text{hasFather} \sqsubseteq \text{hasParent}\}$$
$$\mathcal{A} = \{\text{Father}(f_1), \text{Father}(f_2), \text{Child}(c_1), \text{hasFather}(s_1, f_2)\}$$

- $Q = \text{Human}(x)$
- $\mathcal{T} \models$
 - $\text{Father} \sqsubseteq \text{Human}$
 - $\text{Child} \sqsubseteq \text{Human}$
 - $\exists\text{hasFather.Human} \sqsubseteq \text{Human}$
- $\text{ans}(Q, \mathcal{K}) = \{\langle f_1 \rangle, \langle f_2 \rangle, \langle c_1 \rangle, \langle s_1 \rangle\}$
- Consider the TBox **first!**



QUERY ANSWERING VIA QUERY REWRITING

QUERY REWRITING OVER DL TBOXES

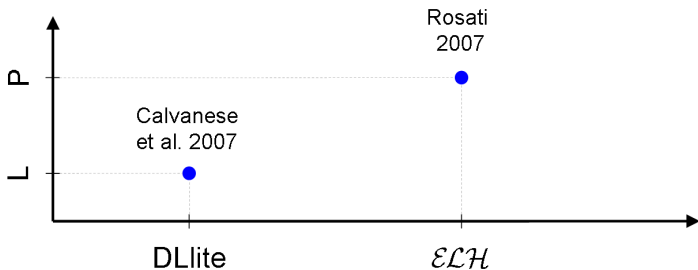


s.t. $\text{ans}(Q, \langle \mathcal{T}, \mathcal{A} \rangle) = \text{ans}(Q', \langle \emptyset, \mathcal{A} \rangle)$, for **all** \mathcal{A}

- Compile the **knowledge** of \mathcal{T} into Q
 - $Q = \text{Human}(x)$
 - $Q' = \text{Human}(x) \vee \text{Child}(x) \vee \text{Father}(x) \vee \dots$
- **Avoid** repeating work with **different** ABoxes



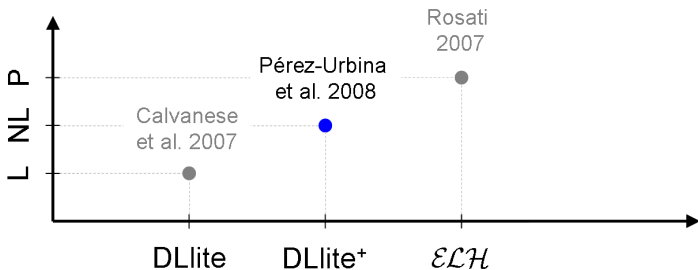
RELATED WORK



- Type of Rewriting
 - DL-Lite \mapsto Union of CQs
 - \mathcal{ELH} \mapsto Datalog program



RELATED WORK

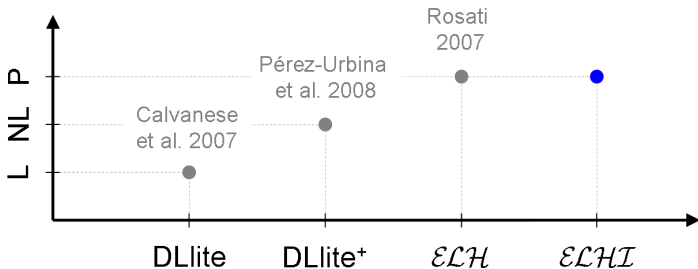


■ Type of Rewriting

- DL-Lite \mapsto Union of CQs
- DL-Lite⁺ \mapsto Linear datalog program
- \mathcal{ELH} \mapsto Datalog program



OUR GOAL



- Generalization and extension of existing approaches
- Worst-case optimal
 - DL-Lite \mapsto Union of CQs
 - DL-Lite⁺ \mapsto Linear datalog program
 - $\mathcal{ELH}, \mathcal{ELHI}$ \mapsto Datalog program



QUERYING \mathcal{ELHI} KBS

\mathcal{ELHI}

- $B \rightarrow A \mid \exists R \mid \exists R.A \mid B_1 \sqcap B_2$
- $R \rightarrow P \mid P^-$
- **TBox** assertions: $B_1 \sqsubseteq B_2$, and $R_1 \sqsubseteq R_2$
- **ABox** assertions: $A(a)$, and $P(a, b)$

QUERIES

- $Q = \langle Q_P, P \rangle$
- P is a set of Horn clauses



USING RESOLUTION

- Goal: Worst-case **optimal** query rewriting algorithm for \mathcal{ELHI} TBoxes
- Idea: Use **known** approaches
 - Motik and Kazakov: **resolution!**

EXAMPLE

$$\mathcal{T} = \{\text{Human} \sqsubseteq \exists \text{hasParent.Human}\}$$
$$\exists(\mathcal{T}) = \{\text{hasParent}(x, f(x)) \leftarrow \text{Human}(x), \text{Human}(f(x)) \leftarrow \text{Human}(x)\}$$

- Functional terms + cycles = **trouble!**



OUR ALGORITHM IN A NUTSHELL

Given a **conjunctive** query $Q = \langle Q_P, P \rangle$ and a *ELHI* TBox \mathcal{T}

- 1. Derive **enough** consequences from Q and \mathcal{T}
- 2. Get rid of **functional** symbols
- 3. Optimize by **unfolding**



1. SATURATION

- Resolution with **free selection** calculus \mathcal{R}^{DL}
- **Saturate** $\Xi(\mathcal{T})$ and P to derive **new** consequences

EXAMPLE

$$Q = \langle Q_P, \{Q_P(x) \leftarrow \text{Human}(x)\} \rangle$$

$$\mathcal{T} = \{ \exists \text{hasParent}^- . \text{Human} \sqsubseteq \text{Human}, \text{Parent} \sqsubseteq \exists \text{hasParent}^- . \text{Human} \}$$

$$(1) \quad \text{Human}(x) \leftarrow \underline{\text{hasParent}(y, x)} \wedge \text{Human}(y)$$

$$(2) \quad \underline{\text{hasParent}(g(x), x)} \leftarrow \text{Parent}(x)$$

$$(3) \quad \underline{\text{Human}(g(x))} \leftarrow \text{Parent}(x)$$

$$(4) \quad \underline{Q_P(x)} \leftarrow \text{Human}(x)$$

$$(5) \quad \underline{\text{Human}(x) \leftarrow \text{Parent}(x) \wedge \text{Human}(g(x))} \quad (1, 2)$$

$$(6) \quad \text{Human}(x) \leftarrow \underline{\text{Parent}(x)} \quad (3, 5)$$

...

- Main challenge: ensure **termination!**



2. ELIMINATION OF FUNCTIONAL TERMS

$\text{ff}(Q, \mathcal{T})$: all **function-free** clauses in the saturation of $\Xi(\mathcal{T}) \cup P$

EXAMPLE CONTD.

$$\text{ff}(Q, \mathcal{T}) = \{ \text{Human}(x) \leftarrow \text{hasParent}(y, x) \wedge \text{Human}(y), \\ \text{Human}(x) \leftarrow \text{Parent}(x), Q_P(x) \leftarrow \text{Human}(x) \}$$

LEMMA

Given a **conjunctive** query Q and a \mathcal{ELHI} TBox \mathcal{T} ,
 $\text{ans}(Q, \langle \mathcal{T}, \mathcal{A} \rangle) = \text{ans}(\text{ff}(Q, \mathcal{T}), \langle \emptyset, \mathcal{A} \rangle)$, for **all** \mathcal{A}

- $\text{ff}(Q, \mathcal{T})$ is a **rewriting**!
- but is it **optimal**?



WHAT HAPPENS WITH DL-LITE⁺?

EXAMPLE

$$Q = \langle Q_P, \{Q_P(x) \leftarrow \text{Human}(x)\} \rangle$$

$$\mathcal{T} = \{ \exists \text{hasParent}. \text{Human} \sqsubseteq \text{Human}, \text{hasFather} \sqsubseteq \text{hasParent} \}$$

$$\text{ff}(Q, \mathcal{T}) = \{ \text{Human}(x) \leftarrow \text{hasParent}(y, x) \wedge \text{Human}(y), \\ \text{hasParent}(x, y) \leftarrow \text{hasFather}(x, y), Q_P(x) \leftarrow \text{Human}(x) \}$$

- We expect a **linear** datalog program!

FROM NONLINEAR TO LINEAR DATALOG PROGRAMS

$$(1) \text{ Human}(x) \leftarrow \text{hasParent}(y, x) \wedge \text{Human}(y)$$

$$(2) \text{ hasParent}(x, y) \leftarrow \text{hasFather}(x, y)$$

$$(1) \text{ Human}(x) \leftarrow \text{hasParent}(y, x) \wedge \text{Human}(y)$$

$$(2) \text{ Human}(x) \leftarrow \text{hasFather}(y, x) \wedge \text{Human}(y)$$



3. UNFOLDING

$\text{rew}(Q, T)$: **unfold** all clauses of the forms

- $A(x) \leftarrow B(x)$,
- $A(x) \leftarrow R(x, y)$,
- $A(x) \leftarrow R(y, x)$,
- $S(x, y) \leftarrow R(x, y)$, and
- $S(x, y) \leftarrow R(y, x)$

in $\text{ff}(Q, T)$, and **get rid** of them

EXAMPLE CONTD.

$$\text{rew}(Q, T) = \{ \text{Human}(x) \leftarrow \text{hasParent}(y, x) \wedge \text{Human}(y), \\ \text{Human}(x) \leftarrow \text{hasFather}(y, x) \wedge \text{Human}(y), \\ Q_P(x) \leftarrow \text{Human}(x) \}$$



OUR ALGORITHM IN A NUTSHELL

PHASES

Given a **conjunctive** query Q and a *ELHI* TBox \mathcal{T}

- 1. **Saturate** to get consequences from Q and \mathcal{T} using \mathcal{R}^{DL}
- 2. Get rid of **functional** symbols: $\text{ff}(Q, \mathcal{T})$
- 3. Optimize by **unfolding**: $\text{rew}(Q, \mathcal{T})$

THEOREM

Given a **conjunctive** query Q and a *ELHI* TBox \mathcal{T} ,
 $\text{ans}(Q, \langle \mathcal{T}, \mathcal{A} \rangle) = \text{ans}(\text{rew}(Q, \mathcal{T}), \langle \emptyset, \mathcal{A} \rangle)$, for **all** \mathcal{A}



PROPERTIES OF THE REWRITING

LEMMA

If \mathcal{T} is in \mathcal{ELHI} or \mathcal{ELH} , then $\text{rew}(Q, \mathcal{T})$ is a **datalog program**

LEMMA

Let $\text{rew}(Q, \mathcal{T}) = \langle Q_P, P' \rangle$. If \mathcal{T} is in DL-Lite^+ , then $P' = U_Q \cup U_C$, such that $\langle Q_P, U_Q \rangle$ is a **union of conjunctive queries** and $\langle Q_P, U_C \rangle$ is a **linear datalog query**

LEMMA

If \mathcal{T} is in DL-Lite , then $\text{rew}(Q, \mathcal{T})$ is a **union of conjunctive queries**



COMPLEXITY ANALYSIS

THEOREM

For a **conjunctive query** Q and a \mathcal{ELHI} KB $\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$, deciding whether $\vec{a} \in \text{ans}(Q, \mathcal{K})$ is **P-complete** data complexity

COROLLARY

Given a **conjunctive query** Q and a KB $\mathcal{K} = \langle \mathcal{T}, \mathcal{A} \rangle$, if \mathcal{T} is in \mathcal{ELH} , DL-Lite^+ , or DL-Lite , we can compute the answers to $\text{rew}(Q, \mathcal{T})$ over \mathcal{A} in **P**, **NL**, or **L** respectively w.r.t. data complexity

- $\text{rew}(Q, \mathcal{T})$ is **optimal!**



CONCLUSION

NEW RESULT

QA over \mathcal{ELHI} KBs is **P-complete** data complexity

CONTRIBUTION

Conjunctive query **rewriting** algorithm for \mathcal{ELHI} KBs

- Worst-case **optimal** for sublanguages of \mathcal{ELHI} for which QA is **P-complete**, **NL-complete**, or in **L** data complexity
- **Generalization** and **extension** of existing approaches
- Straightforward use of **DB technology** for query evaluation



CURRENT AND FUTURE WORK

CURRENT WORK

Nominals: \mathcal{ELHIO}

- $\{Pope\}$

FUTURE WORK

Prototype system in an [Information Integration](#) setting

- ComparaGRID



THANKS!

