

# QUONTO: ontology-based data access and integration using relational technology

Giuseppe De Giacomo



SAPIENZA  
UNIVERSITÀ DI ROMA

Semantic Days 2009

# Motivation: ontologies and data

- The best current standard DL reasoning systems can deal with moderately large ABoxes.  $\sim 10^4$  individuals (*and this is a big achievement of the last years!*)
- But data of interests in typical information systems are much **larger**  $\sim 10^6 - 10^9$  individuals
- The best technology to deal with large amounts of data are **relational databases**.

## Question:

How can we use ontologies together with large amounts of data?

## Answer:

**Yes**, by using the system **QUONTO**.

# Which query language?

Which query language to use?

Two extreme cases:

- 1 **Just classes and properties** of the ontology  $\leadsto$  instance checking
  - Ontology languages are tailored for capturing intensional relationships.
  - They are quite **poor as query languages**:  
Cannot refer to same object via multiple navigation paths in the ontology, i.e., allow only for a limited form of JOIN, namely chaining.
- 2 **Full SQL** (or equivalently, first-order logic)
  - Problem: in the presence of incomplete information, query answering becomes **undecidable** (FOL validity).

A good compromise are (unions of) **conjunctive queries**.

A **conjunctive query (CQ)** is a first-order query of the form

$$q(\vec{x}) \leftarrow \exists \vec{y}. R_1(\vec{x}, \vec{y}) \wedge \dots \wedge R_k(\vec{x}, \vec{y})$$

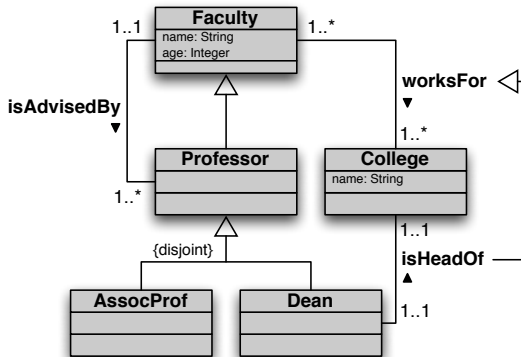
where each  $R_i(\vec{x}, \vec{y})$  is an atom using (some of) the free variables  $\vec{x}$ , the existentially quantified variables  $\vec{y}$ , and possibly constants.

*Note:*

- CQs contain no disjunction, no negation, no universal quantification.
- Correspond to SQL/relational algebra **select-project-join (SPJ) queries** – the most frequently asked queries.
- They can form the core of **SPARQL** queries.

# Example of conjunctive query

Professor	$\sqsubseteq$	Faculty
AssocProf	$\sqsubseteq$	Professor
Dean	$\sqsubseteq$	Professor
AssocProf	$\sqsubseteq$	$\neg$ Dean
Faculty	$\sqsubseteq$	$\exists$ age
$\exists$ age $^{-}$	$\sqsubseteq$	Integer
$\exists$ worksFor	$\sqsubseteq$	Faculty
$\exists$ worksFor $^{-}$	$\sqsubseteq$	College
Faculty	$\sqsubseteq$	$\exists$ worksFor
College	$\sqsubseteq$	$\exists$ worksFor $^{-}$
	$\vdots$	



$q(nf, nd, av) \leftarrow \exists f, c, d.$

$worksFor(f, c) \wedge isHeadOf(d, c) \wedge name(f, nf) \wedge name(d, nd) \wedge$   
 $age(f, av) \wedge age(d, av)$

## Query answering: certain answers to a query

Let  $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$  be an ontology,  $\mathcal{I}$  an interpretation for  $\mathcal{O}$ , and  $q(\vec{x}) \leftarrow \exists \vec{y}. conj(\vec{x}, \vec{y})$  a CQ.

Def.: The **answer** to  $q(\vec{x})$  over  $\mathcal{I}$ , denoted  $q^{\mathcal{I}}$

... is the set of **tuples  $\vec{c}$  of constants of  $\mathcal{A}$**  such that the formula  $\exists \vec{y}. conj(\vec{c}, \vec{y})$  evaluates to true in  $\mathcal{I}$ .

We are interested in finding those answers that hold in all models of an ontology.

Def.: The **certain answers** to  $q(\vec{x})$  over  $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$ , denoted  $cert(q, \mathcal{O})$

... are the **tuples  $\vec{c}$  of constants of  $\mathcal{A}$**  such that  $\vec{c} \in q^{\mathcal{I}}$ , for **every model  $\mathcal{I}$**  of  $\mathcal{O}$ .

# Complexity of query answering in ontologies

Studied extensively for (unions of) CQs and various ontology languages:

	Combined complexity	Data complexity
Plain databases	NP-complete	in LOGSPACE <sup>(2)</sup>
OWL 2 (and less)	2EXPTIME-complete	coNP-hard <sup>(1)</sup>

(1) Already for a TBox with a single disjunction!

(2) This is what we need to scale with the data.

## Question

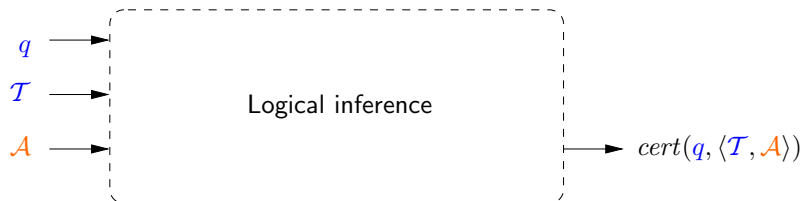
- Can we find interesting DLs for which the query answering problem can be solved efficiently (i.e., in LOGSPACE)?
- Can we leverage relational database technology for query answering?

## Answer

**Yes, but we need new foundations!**

No more tableaux coming from logic, but **chase** coming from databases as main tool for reasoning!

# Inference in query answering

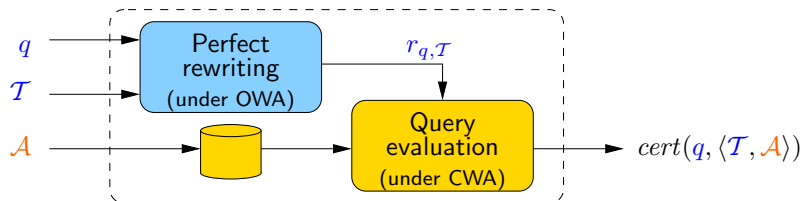


To be able to deal with data efficiently, we need to separate the contribution of  $\mathcal{A}$  from the contribution of  $q$  and  $\mathcal{T}$ .

→ Query answering by **query rewriting**.



# Query rewriting



Query answering can **always** be thought as done in two phases:

- 1 Perfect rewriting: from  $q$  and  $\mathcal{T}$  generate a new query  $r_{q,\mathcal{T}}$ .
- 2 Query evaluation: evaluate  $r_{q,\mathcal{T}}$  over the ABox  $\mathcal{A}$  seen as a complete database.  
 $\rightsquigarrow$  Produces  $cert(q, \langle \mathcal{T}, \mathcal{A} \rangle)$ .

Note: The “always” holds if we pose no restriction on the language in which to express the rewriting  $r_{q,\mathcal{T}}$ .

The expressiveness of the ontology language affects the **query language into which we are able to rewrite CQs**:

- When we can rewrite into FOL/SQL.  
↪ Query evaluation can be done in SQL, i.e., via an RDBMS (Note: FOL is in LOGSPACE).
- When we can rewrite into an NLOGSPACE-hard language.  
↪ Query evaluation requires (at least) linear recursion.
- When we can rewrite into a PTIME-hard language.  
↪ Query evaluation requires full recursion (e.g., Datalog).
- When we can rewrite into a CONP-hard language.  
↪ Query evaluation requires (at least) power of Disjunctive Datalog.

# The QUONTO description logic: $DL-Lite_{\mathcal{A}}$

- QUONTO is based on  $DL-Lite_{\mathcal{A}}$ .
- $DL-Lite_{\mathcal{A}}$  is carefully designed to provide robust foundations for Ontology-Based Data Access: Query answering for UCQ is:
  - NP-complete in query complexity – as relational DBs
  - PTIME in the size of the TBox
  - LOGSPACE in size of ABox (data complexity) – as relational DBs
  - queries can be rewritten into FOL/SQL – allows delegating reasoning on data to a RDMBS!
- Inference based on (inverted) chase and not on tableaux.

ISA between classes	$A_1 \sqsubseteq A_2$
Disjointness between classes	$A_1 \sqsubseteq \neg A_2$
Domain and range of properties	$\exists P \sqsubseteq A_1 \quad \exists P^- \sqsubseteq A_2$
Mandatory participation ( <i>min card</i> = 1)	$A_1 \sqsubseteq \exists P \quad A_2 \sqsubseteq \exists P^-$
Functionality of relations ( <i>max card</i> = 1)	<b>(funct</b> $P$ ) <b>(funct</b> $P^-$ )
ISA between properties	$Q_1 \sqsubseteq Q_2$
Disjointness between properties	$Q_1 \sqsubseteq \neg Q_2$

- Captures all the basic constructs of **UML Class Diagrams** and of the **ER Model** ...
- ... *except covering constraints* in generalizations. – if we add them, query answering becomes **CONP-hard** in data complexity
- A substantial fragment, chosen as one **one of** the three candidate **OWL 2 Profiles**: **OWL 2 QL**.
- Extends (the DL compatible part of) the ontology language **RDFS**.

# Beyond $DL\text{-Lite}_A$ : results on data complexity

	lhs	rhs	funct.	Prop. incl.	Data complexity of query answering
0	$DL\text{-Lite}_A$		$\sqrt{*}$	$\sqrt{*}$	in LOGSPACE
1	$A \mid \exists P.A$	$A$	—	—	NLOGSPACE-hard
2	$A$	$A \mid \forall P.A$	—	—	NLOGSPACE-hard
3	$A$	$A \mid \exists P.A$	$\checkmark$	—	NLOGSPACE-hard
4	$A \mid \exists P.A \mid A_1 \sqcap A_2$	$A$	—	—	PTIME-hard
5	$A \mid A_1 \sqcap A_2$	$A \mid \forall P.A$	—	—	PTIME-hard
6	$A \mid A_1 \sqcap A_2$	$A \mid \exists P.A$	$\checkmark$	—	PTIME-hard
7	$A \mid \exists P.A \mid \exists P^-.A$	$A \mid \exists P$	—	—	PTIME-hard
8	$A \mid \exists P \mid \exists P^-$	$A \mid \exists P \mid \exists P^-$	$\checkmark$	$\checkmark$	PTIME-hard
9	$A \mid \neg A$	$A$	—	—	CONP-hard
10	$A$	$A \mid A_1 \sqcup A_2$	—	—	CONP-hard
11	$A \mid \forall P.A$	$A$	—	—	CONP-hard

## Notes:

- \* with the “proviso” of not specializing functional properties.
- NLOGSPACE and PTIME hardness holds already for instance checking.
- For CONP-hardness in line 10, a TBox with a single assertion  $A_L \sqsubseteq A_T \sqcup A_F$  suffices!  $\rightsquigarrow$  No hope of including covering constraints.

# Example

TBox:  $\text{Professor} \sqsubseteq \exists \text{teaches}$   
 $\exists \text{teaches}^- \sqsubseteq \text{Course}$

Query:  $q(x) \leftarrow \text{teaches}(x, y), \text{Course}(y)$

Perfect Reformulation:  $q(x) \leftarrow \text{teaches}(x, y), \text{Course}(y)$   
 $q(x) \leftarrow \text{teaches}(x, y), \text{teaches}(-, y)$   
 $q(x) \leftarrow \text{teaches}(x, -)$   
 $q(x) \leftarrow \text{Professor}(x)$

ABox:  $\text{teaches}(\text{john}, \text{kdbd})$   
 $\text{Professor}(\text{mary})$

It is easy to see that  $\text{Eval}(\text{SQL}(r_{q, \mathcal{T}}), \text{DB}(\mathcal{A}))$  in this case produces as answer  $\{\text{john}, \text{mary}\}$ .

- Includes support for:
  - *DL-Lite<sub>A</sub>*
  - Identification path constrains
  - Denial constrains
  - Epistemic constrains
  - Union of conjunctive queries – *expressed in Datalog or SPARQL*
  - Epistemic queries –*expressed in SparSQL*
- Reasoning services are highly optimized
- Can be used with internal and external DBMS (include drivers for Oracle, DB2, IBM Information Integrator, SQL Server, MySQL, etc.)
- Implemented in Java – *API are available for selected projects upon request*
- Several wrapped versions publicly available at:  
<http://www.dis.uniroma1.it/~quonto/>     (*or just google “quonto”*)



# QUONTO wrapped versions

<http://www.dis.uniroma1.it/~quonto/>

## DIG Server wrapper + ODBA Protégé plugin

by Mariano Rodriguez-Muro, Univ. Bolzano

The screenshot displays the Protégé 3.3.1 ODBA plugin interface. The main window is titled "univ-bench-NewMap Protégé 3.3.1" and shows a "DATASOURCE MANAGER" tab. On the left, the "DATASOURCE BROWSER" shows a project named "univ-bench-NewMap" with a data source named "Sapienza". Below this, the "For datasource:" section provides configuration details for "Sapienza":

- Type: RDEMS
- Mapping Type: ODBAMappings
- Source ID: univ-bench-NewMap
- JDBC URL: jdbc:mysql://localhost/
- Database Name: stud
- Database Username: quonto
- Database Password: quonto
- JDBC Driver: com.mysql.jdbc.Driver

The main editor area shows a list of SQL queries for various classes:

- Dean**
  - Query: `SELECT university_dean.CODE AS DEANCODE, university.CODE AS UNIVCODE FROM university_dean, university WHERE university_dean.UNIVCODE = university.CODE`
- WorksFor**
  - Query: `SELECT professor.CODE AS PROFCODE, university.CODE AS UNIVCODE FROM professor, university WHERE professor.UNIVCODE = university.CODE`
- teaching**
  - Query: `SELECT course_assignment.PROFESSOR_CODE, course_assignment.COURSE_CODE FROM course_assignment`
- Ass\_Professor**
  - Query: `SELECT professor_data.CODE FROM professor_data`
- University**
  - Query: `SELECT CODE, CITY, DESCRIPTION FROM university`
- hasAlumnus**
- Course**
- Student**

# QUONTO wrapped versions

<http://www.dis.uniroma1.it/~quonto/>

**Qtoolkit**: simple graphical interface, only standard ABoxes (no connection to external DBs)



(BasicVersion) REL. ALPHA 0.137

---

QToolkit is a simple textual interface to the QuOnto system.  
QToolkit makes use of H2 as an embedded RDBMS.

© Dipartimento di Informatica e Sistemistica  
SAPIENZA Università di Roma

---

QToolkit is developed by Emma Di Pasquale, Marco Ruzzi, Claudio Corona, and Domenico Fabio Savo

**ROWLkit**: first implementation of the OWL 2 QL Profile



REL. ALPHA 1.0

ROWLKit is a very simple GUI toolkit for OWL 2 QL ontologies. The main features are:

- language check: parses an OWL file and verifies if such file is expressed in the OWL 2 QL fragment
- intensional reasoning: ontology classification and basic classes (and properties) subsumption and satisfiability
- query answering: evaluates union of conjunctive queries expressed in SPARQL

ROWLKit makes use of H2 as an embedded RDBMS.

© Dipartimento di Informatica e Sistemistica  
SAPIENZA Università di Roma

ROWLKit is developed by Marco Ruzzi, Claudio Corona, and Domenico Fabio Savo

- Ontology-based data access and data integration is **ready for prime time** – see *Calvanese&De Giacomo's tutorial*
- **QUONTO** provides serious proof of concept of this.
- We are successfully applying QUONTO in various **full-fledged case studies** – see *Diego Calvanese's talk*
- We are currently **looking for projects** where to apply such technology further!
- This technology is **ready to become a product!**

People involved in this work:

- Sapienza Università di Roma
  - Claudio Corona
  - Domenico Lembo
  - Maurizio Lenzerini
  - Antonella Poggi
  - Riccardo Rosati
  - Marco Ruzzi
  - Domenico Fabio Savo
- Libera Università di Bolzano
  - Diego Calvanese
  - Mariano Rodriguez Muro
- Several past master and PhD students (thanks!)