

Reasoning for Ontology Engineering and Usage

Part 4: Integrating Data into Ontologies

Diego Calvanese¹, Giuseppe De Giacomo²,
Demo by Mariano Rodríguez-Muro¹

¹ Free University of Bozen-Bolzano,



FREE UNIVERSITÄT BOZEN
LIBERA UNIVERSITÀ DI BOLZANO
FREE UNIVERSITY OF BOZEN - BOLZANO

² SAPIENZA Università di Roma



SAPIENZA
UNIVERSITÀ DI ROMA



Tutorial at 7th Int. Semantic Web Conference (ISWC 2008)
Karlsruhe – Oct. 27, 2008

- 1 Introduction
- 2 Querying data through ontologies
- 3 *DL-Lite_A*: an ontology language for accessing data
- 4 Ontology-based data integration
- 5 References



- 1 Introduction
- 2 Querying data through ontologies
- 3 *DL-Lite_A*: an ontology language for accessing data
- 4 Ontology-based data integration
- 5 References



- We have seen that current DL reasoning systems can deal with relatively large ABoxes. $\leadsto 10^4$ individuals
- This is small, if compared to data found in various contexts: biological data, scientific data, enterprise data, ...
 $\leadsto 10^5 - 10^9$ 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?



Challenges when integrating data into ontologies

Deal with well-known tradeoff between **expressive power** of the ontology language and **complexity** of dealing with (i.e., performing inference over) ontologies in that language.

Requirements come from the specific setting:

- We have to fully take into account the ontology.
~> **inference**
- We have to deal very large amounts of data.
~> **relational databases**
- We want flexibility in querying the data.
~> **expressive query language**
- We want to keep the data in the sources, and not move it around.
~> **map** data sources to the ontology (cf. [Data Integration](#))



Questions addressed in this part of the tutorial

- 1 Which is the “right” **query language**?
- 2 Which is the “right” **ontology language**?
- 3 How can we bridge the **semantic mismatch** between the ontology and the data sources?
- 4 How can **tools for ontology-based data access and integration** fully take into account all these issues?



- 1 Introduction
- 2 Querying data through ontologies
- 3 *DL-Lite_A*: an ontology language for accessing data
- 4 Ontology-based data integration
- 5 References



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).



Conjunctive queries (CQs)

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.

We will also use the simpler Datalog notation:

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

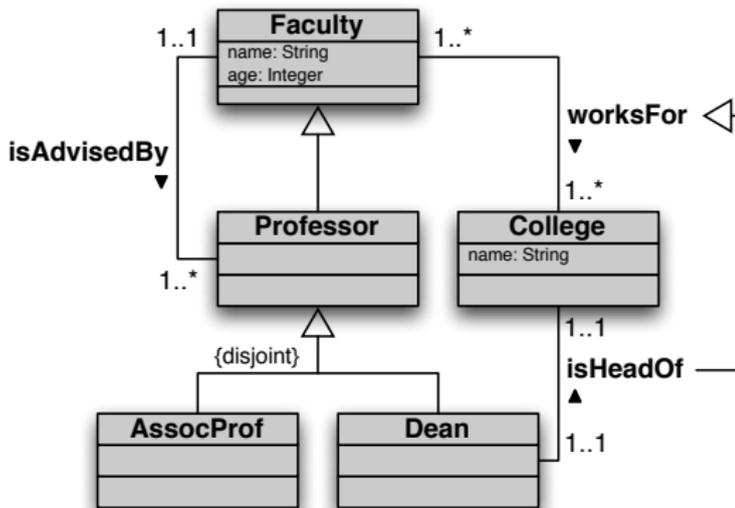
Note:

- 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 also be written as **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	



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

$worksFor(f, c) \wedge isHeadOf(d, c) \wedge name(f, nf) \wedge name(d, nd) \wedge$
 $age(f, af) \wedge age(d, ad) \wedge af = ad$



Conjunctive queries and SQL – Example

Relational alphabet:

`worksFor(fac, coll), isHeadOf(dean, coll), name(p, n), age(p, a)`

Query: return name, age, and name of dean of all faculty that have the same age as their dean.

Expressed in SQL:

```
SELECT NF.name, AF.age, ND.name
FROM worksFor W, isHeadOf H, name NF, name ND, age AF, age AD
WHERE W.fac = NF.p AND W.fac = AF.p AND
      H.dean = ND.p AND H.dean = AD.p AND
      W.coll = H.coll AND AF.a = AD.a
```

Expressed as a CQ:

$$q(nf, af, nd) \leftarrow \text{worksFor}(f1, c1), \text{isHeadOf}(d1, c2), \\ \text{name}(f2, nf), \text{name}(d2, nd), \text{age}(f3, af), \text{age}(d3, ad), \\ f1 = f2, f1 = f3, d1 = d2, d1 = d3, c1 = c2, af = ad$$


There are fundamentally different assumptions when addressing query answering in different settings:

- **traditional database assumption**
- **knowledge representation assumption**

Note: for the moment we assume to deal with an ordinary ABox, which however may be very large and thus is stored in a database.



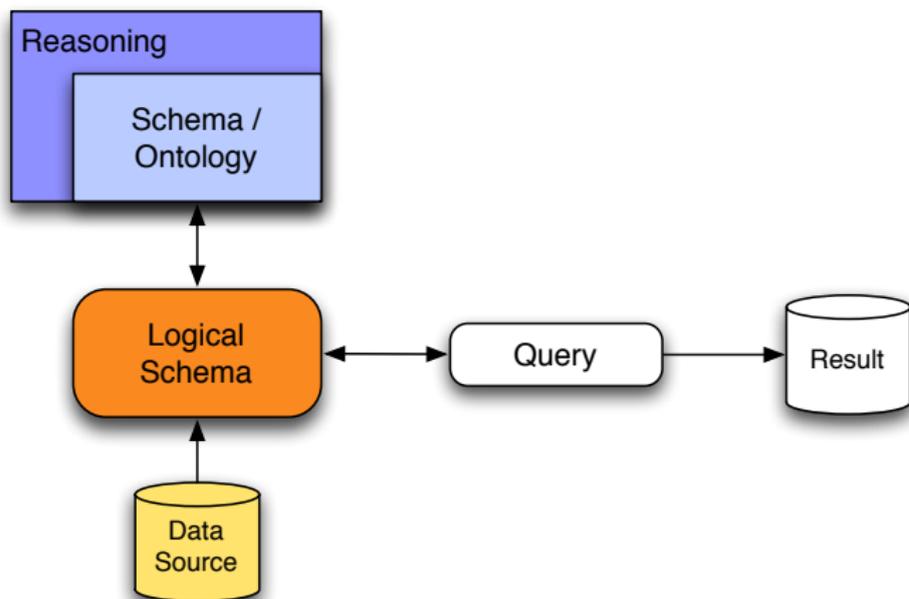
Query answering under the database assumption

- Data are completely specified (CWA), and typically large.
- Schema/intensional information used in the design phase.
- At **runtime**, the data is assumed to satisfy the schema, and therefore the **schema is not used**.
- Queries allow for complex navigation paths in the data (cf. SQL).

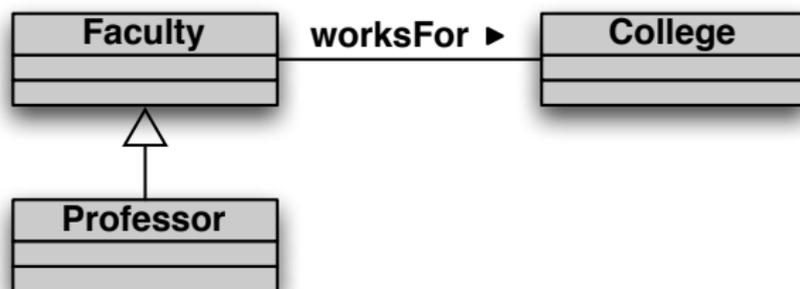
↪ Query answering amounts to **query evaluation**, which is computationally easy.



Query answering under the database assumption (cont'd)



Query answering under the database assumption – Example



For each class/property we have a (complete) table in the database.

DB: Faculty = { john, mary, nick }

Professor = { john, nick }

College = { collA, collB }

worksFor = { (john,collA), (mary,collB) }

Query: $q(x) \leftarrow \exists c. \text{Professor}(x), \text{College}(c), \text{worksFor}(x, c)$

Answer: { john }



Query answering under the KR assumption

- An ontology imposes constraints on the data.
- Actual data may be incomplete or inconsistent w.r.t. such constraints.
- The system has to take into account the constraints during query answering, and overcome incompleteness or inconsistency.

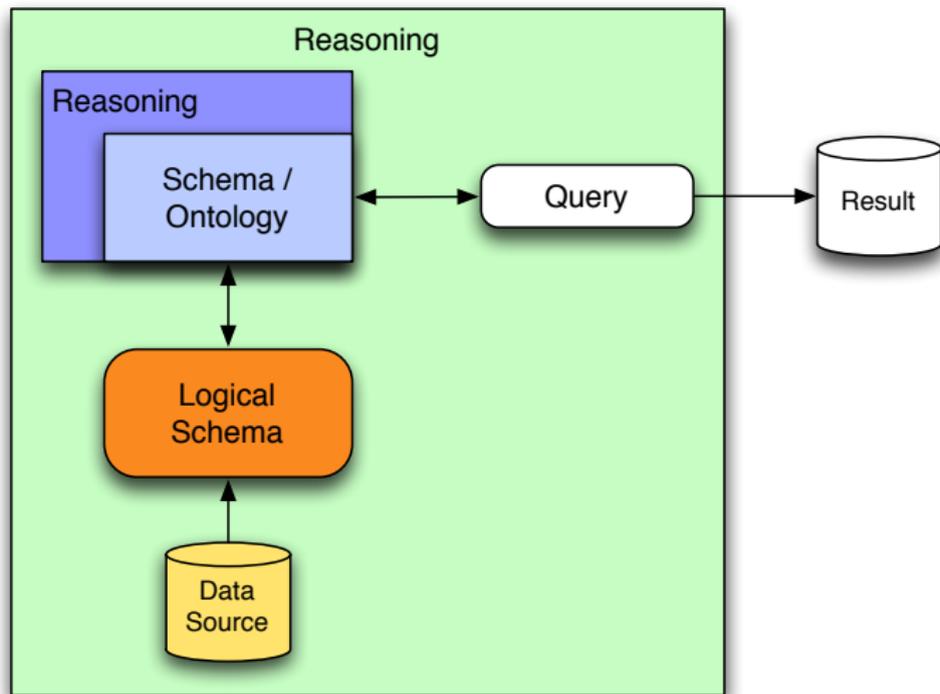
↪ Query answering amounts to **logical inference**, which is computationally more costly.

Note:

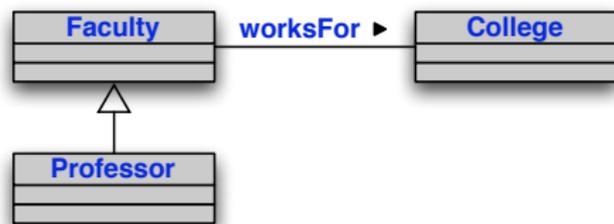
- Size of the data is not considered critical (comparable to the size of the intensional information).
- Queries are typically simple, i.e., atomic (a class name), and query answering amounts to instance checking.



Query answering under the KR assumption (cont'd)



Query answering under the KR assumption – Example



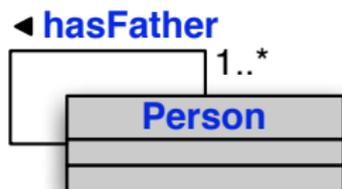
The tables in the database may be **incompletely specified**, or even missing for some classes/properties.

DB: Professor \supseteq { john, nick }
College \supseteq { collA, collB }
worksFor \supseteq { (john,collA), (mary,collB) }

Query: $q(x) \leftarrow \text{Faculty}(x)$

Answer: { john, nick, mary }

Query answering under the KR assumption – Example 2



Each person has a father, who is a person.

DB: $\text{Person} \supseteq \{ \text{john, nick, toni} \}$
 $\text{hasFather} \supseteq \{ (\text{john,nick}), (\text{nick,toni}) \}$

Queries: $q_1(x, y) \leftarrow \text{hasFather}(x, y)$

$q_2(x) \leftarrow \exists y. \text{hasFather}(x, y)$

$q_3(x) \leftarrow \exists y_1, y_2, y_3. \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2), \text{hasFather}(y_2, y_3)$

$q_4(x, y_3) \leftarrow \exists y_1, y_2. \text{hasFather}(x, y_1), \text{hasFather}(y_1, y_2), \text{hasFather}(y_2, y_3)$

Answers: to q_1 : $\{ (\text{john,nick}), (\text{nick,toni}) \}$

to q_2 : $\{ \text{john, nick, toni} \}$

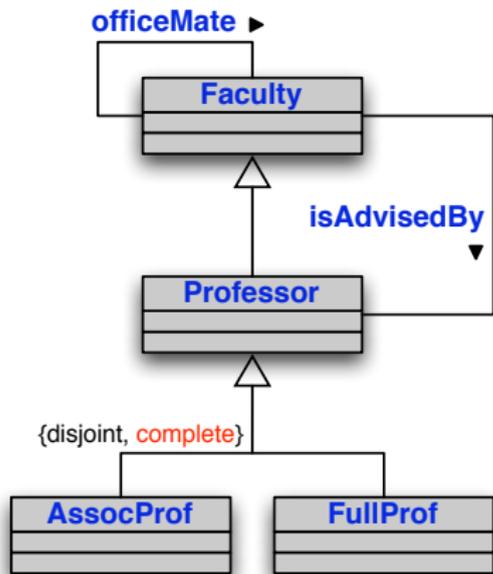
to q_3 : $\{ \text{john, nick, toni} \}$

to q_4 : $\{ \}$

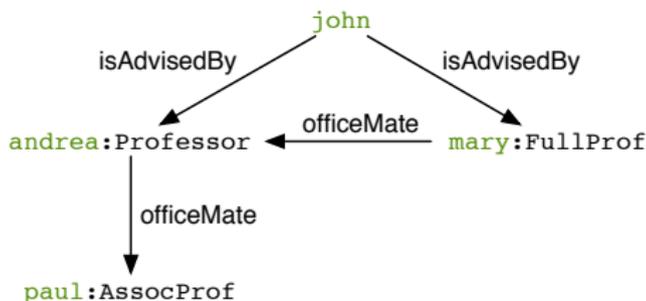


QA under the KR assumption – Andrea's Example

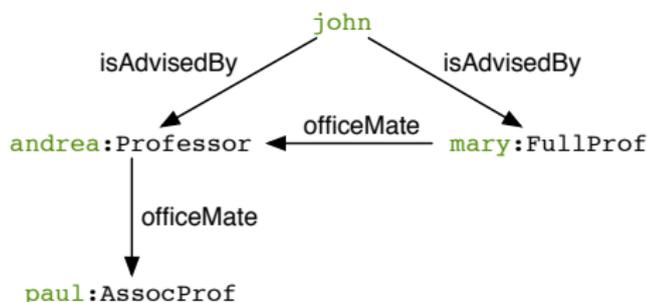
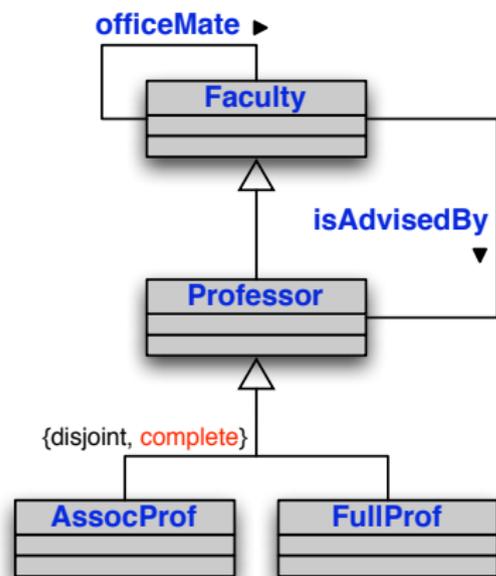
FullProf \equiv AssocProf \sqcup FullProf



- Faculty \supseteq { andrea, nick, mary, john }
- Professor \supseteq { andrea, nick, mary }
- AssocProf \supseteq { nick }
- FullProf \supseteq { mary }
- isAdvisedBy \supseteq { (john, andrea), (john, mary) }
- officeMate \supseteq { (mary, andrea), (andrea, nick) }



QA under the KR assumption – Andrea's Example (cont'd)



$q() \leftarrow \exists y, z.$
 $\text{isAdvisedBy}(\text{john}, y), \text{FullProf}(y),$
 $\text{officeMate}(y, z), \text{AssocProf}(z)$

Answer: **yes!**

To determine this answer, we need to resort to **reasoning by cases**.

We have to face the difficulties of both assumptions:

- The actual **data** is stored in external information sources (i.e., databases), and thus its size is typically **very large**.
- The ontology introduces **incompleteness** of information, and we have to do logical inference, rather than query evaluation.
- We want to take into account at **runtime** the **constraints** expressed in the ontology.
- We want to answer **complex database-like queries**.
- We may have to deal with multiple information sources, and thus face also the problems that are typical of data integration.



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} .

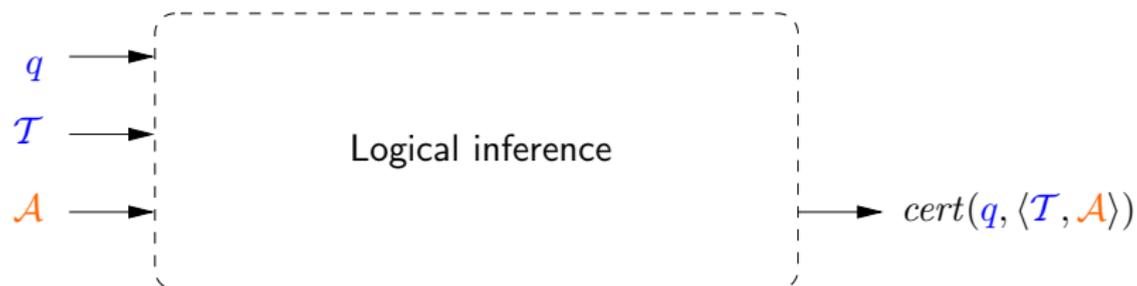
We are interested in finding those answers that hold in all models of an ontology.

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} .



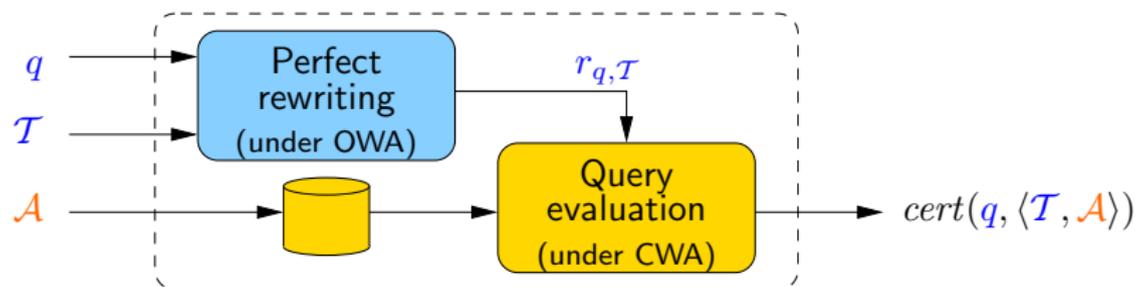
Inference in query answering



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



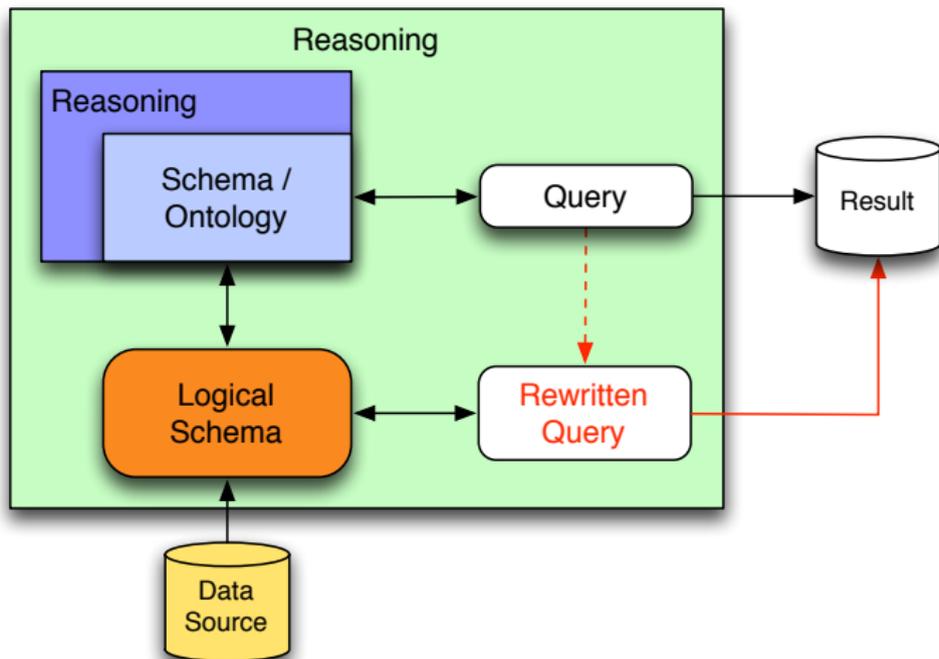
Query answering can **always** be thought as done in two phases:

- 1 **Perfect rewriting**: produce from q and the TBox \mathcal{T} a new query $r_{q,\mathcal{T}}$ (called the perfect rewriting of q w.r.t. \mathcal{T}).
- 2 **Query evaluation**: evaluate $r_{q,\mathcal{T}}$ over the ABox \mathcal{A} seen as a complete database (and without considering the TBox \mathcal{T}).
 \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 $r_{q,\mathcal{T}}$.



Query rewriting (cont'd)



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**.
 \rightsquigarrow 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.
 \rightsquigarrow Query evaluation requires (at least) **linear recursion**.
- When we can rewrite into a **PTIME-hard** language.
 \rightsquigarrow Query evaluation requires full recursion (e.g., **Datalog**).
- When we can rewrite into a **CONP-hard** language.
 \rightsquigarrow Query evaluation requires (at least) power of **Disjunctive Datalog**.



Complexity of query answering in DLs

Problem of rewriting is related to **complexity of query answering**.

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

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

- (1) Already for a TBox with a single disjunction (see Andrea's example).
(2) This is what we need to scale with the data.

Questions

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



- 1 Introduction
- 2 Querying data through ontologies
- 3 *DL-Lite_A*: an ontology language for accessing data
- 4 Ontology-based data integration
- 5 References



The *DL-Lite* family

- A family of DLs optimized according to the tradeoff between expressive power and **complexity** of query answering, with emphasis on **data**.
- Carefully designed to have nice computational properties for answering UCQs (i.e., computing certain answers):
 - The same complexity as relational databases.
 - In fact, query answering can be delegated to a relational DB engine.
 - The DLs of the *DL-Lite* family are essentially the maximally expressive ontology languages enjoying these nice computational properties.
- We present *DL-Lite_A*, an expressive member of the *DL-Lite* family.

DL-Lite_A provides robust foundations for Ontology-Based Data Access.



TBox assertions:

- Class inclusion assertions: $B \sqsubseteq C$, with:

$$\begin{array}{l} B \longrightarrow A \mid \exists Q \\ C \longrightarrow C \mid \neg C \end{array}$$

- Property inclusion assertions: $Q \sqsubseteq R$, with:

$$\begin{array}{l} Q \longrightarrow P \mid P^- \\ R \longrightarrow Q \mid \neg Q \end{array}$$

- Functionality assertions: (**funct** Q)
- **Proviso:** functional properties cannot be specialized.

ABox assertions: $A(c)$, $P(c_1, c_2)$, with c_1, c_2 constants

Note: DL-Lite_A distinguishes also between object and data properties (ignored here).



Semantics of the $DL-Lite_A$ assertions

Assertion	Syntax	Example	Semantics
class incl.	$B \sqsubseteq C$	Father $\sqsubseteq \exists$ child	$B^I \subseteq C^I$
o-prop. incl.	$Q \sqsubseteq R$	father \sqsubseteq anc	$Q^I \subseteq R^I$
v.dom. incl.	$E \sqsubseteq F$	$\rho(\text{age}) \sqsubseteq \text{xsd:int}$	$E^I \subseteq F^I$
d-prop. incl.	$U \sqsubseteq V$	offPhone \sqsubseteq phone	$U^I \subseteq V^I$
o-prop. funct.	(funct Q)	(funct father)	$\forall o, o', o''. (o, o') \in Q^I \wedge$ $(o, o'') \in Q^I \rightarrow o' = o''$
d-prop. funct.	(funct U)	(funct ssn)	$\forall o, v, v'. (o, v) \in U^I \wedge$ $(o, v') \in U^I \rightarrow v = v'$
mem. asser.	$A(c)$	Father(bob)	$c^I \in A^I$
mem. asser.	$P(c_1, c_2)$	child(bob, ann)	$(c_1^I, c_2^I) \in P^I$
mem. asser.	$U(c, d)$	phone(bob, '2345')	$(c^I, \text{val}(d)) \in U^I$

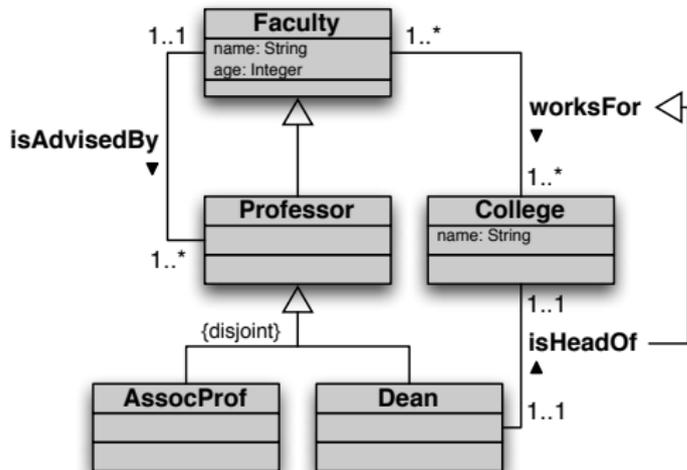


Capturing basic ontology constructs in $DL-Lite_A$

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 ($min\ card = 1$)	$A_1 \sqsubseteq \exists P \quad A_2 \sqsubseteq \exists P^-$
Functionality of relations ($max\ card = 1$)	($funct\ P$) ($funct\ P^-$)
ISA between properties	$Q_1 \sqsubseteq Q_2$
Disjointness between properties	$Q_1 \sqsubseteq \neg Q_2$



Example



Professor \sqsubseteq Faculty
 AssocProf \sqsubseteq Professor
 Dean \sqsubseteq Professor
 AssocProf \sqsubseteq \neg Dean

Faculty \sqsubseteq \exists age
 \exists age $^-$ \sqsubseteq xsd:int
 (func \sqsubseteq age)

\exists worksFor \sqsubseteq Faculty
 \exists worksFor $^-$ \sqsubseteq College
 Faculty \sqsubseteq \exists worksFor
 College \sqsubseteq \exists worksFor $^-$

\exists isHeadOf \sqsubseteq Dean
 \exists isHeadOf $^-$ \sqsubseteq College
 Dean \sqsubseteq \exists isHeadOf
 College \sqsubseteq \exists isHeadOf $^-$
 isHeadOf \sqsubseteq worksFor

(func \sqsubseteq isHeadOf)
 (func \sqsubseteq isHeadOf $^-$)

Note: $DL\text{-}Lite_{\mathcal{A}}$ cannot capture completeness of a hierarchy. This would require **disjunction** (i.e., **OR**).



- Captures all the basic constructs of **UML Class Diagrams** and of the **ER Model** ...
- ... **except covering constraints** in generalizations.
- Is **one of** the three candidate **OWL 2 Profiles**.
- Extends (the DL fragment of) the ontology language **RDFS**.
- Is completely symmetric w.r.t. **direct and inverse properties**.
- Does **not** enjoy the **finite model property**, i.e., reasoning and query answering differ depending on whether we consider or not also infinite models.



Query answering in $DL-Lite_A$

Based on **query reformulation**: given an (U)CQ and an ontology:

- 1 **Compute its perfect rewriting**, which turns out to be a UCQ.
- 2 **Evaluate the perfect rewriting** on the ABox seen as a DB.

To **compute the perfect rewriting**, starting from the original (U)CQ, iteratively get a CQ to be processed and either:

- **expand** positive inclusions & **simplify** redundant atoms, or
- **unify** atoms in the CQ to obtain a more specific CQ to be further expanded.

Each result of the above steps is added to the queries to be processed.

Note: negative inclusions and functionalities play a role in ontology satisfiability, but not in query answering.



Query answering in $DL-Lite_{\mathcal{A}}$ – Example

TBox: $\text{Professor} \sqsubseteq \exists \text{worksFor}$
 $\exists \text{worksFor}^{-} \sqsubseteq \text{College}$

Query: $q(x) \leftarrow \text{worksFor}(x, y), \text{College}(y)$

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

ABox: $\text{worksFor}(\text{john}, \text{collA})$ $\text{Professor}(\text{john})$
 $\text{worksFor}(\text{mary}, \text{collB})$ $\text{Professor}(\text{nick})$

Evaluating the last two queries over the ABox (seen as a DB) produces as answer $\{\text{john}, \text{nick}, \text{mary}\}$.



Complexity of reasoning in $DL-Lite_A$

Ontology satisfiability and all classical DL reasoning tasks are:

- Efficiently tractable in the size of $TBox$ (i.e., $PTime$).
- Very efficiently tractable in the size of the $ABox$ (i.e., $LOGSPACE$).

In fact, reasoning can be done by constructing suitable FOL/SQL queries and evaluating them over the $ABox$ (**FOL-rewritability**).

Query answering for CQs and UCQs is:

- $PTime$ in the size of $TBox$.
- $LOGSPACE$ in the size of the $ABox$.
- Exponential in the size of the **query** (**NP-complete**).
Bad? ... not really, this is exactly as in relational DBs.

Can we go beyond $DL-Lite_A$?

No! By adding essentially any additional constructor we lose these nice computational properties.

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

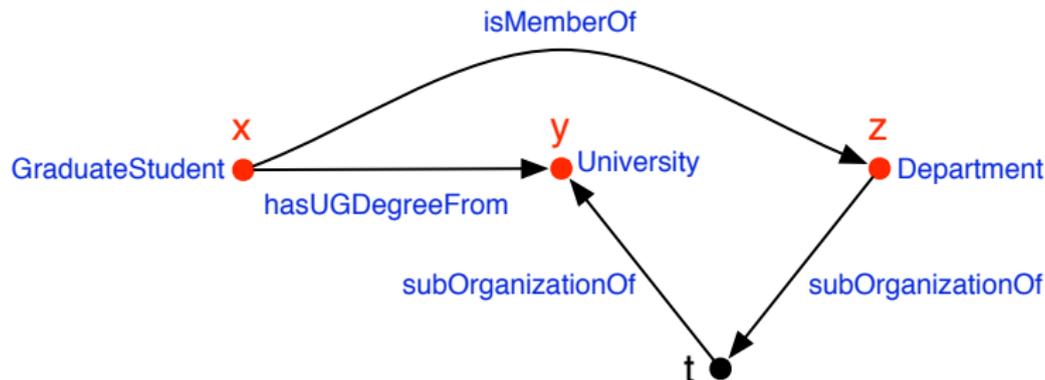
Notes:

- * 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 of query

$q(x, y, z) \leftarrow$ GraduateStudent(x), University(y), Department(z),
hasUndergraduateDegreeFrom(x, y), isMemberOf(x, z),
subOrganizationOf(z, t), subOrganizationOf(t, y)



SELECT ?X ?Y ?Z WHERE

?X rdf:type 'GraduateStudent' . ?Y rdf:type 'University' .

?Z rdf:type 'Department' .

?X :hasUndergraduateDegreeFrom ?Y . ?X :isMemberOf ?Z .

?Z subOrganizationOf ?T . ?T subOrganizationOf ?Y



- 1 Introduction
- 2 Querying data through ontologies
- 3 *DL-Lite_A*: an ontology language for accessing data
- 4 **Ontology-based data integration**
- 5 References



As said, $DL-Lite_{\mathcal{A}}$ can be seen as a fragment of OWL 2 specifically designed to deal with **large amounts of data**:

- It allows for **query answering** (and checking ontology satisfiability) in **LOGSPACE** wrt the size of the data (i.e., the ABox).
- Reasoning with data in $DL-Lite_{\mathcal{A}}$ can be **delegated to a relational DBMS**.
- $DL-Lite_{\mathcal{A}}$ captures all constructs in typical conceptual models, such as **UML Class Diagrams** and **ER**.



Next we look at **semantical interoperation** between an **ontology-based system** and **external systems**, such as traditional information systems.

- In doing this, we consider for concreteness an actual scenario, where semantic interoperation is the core issue ...
 \rightsquigarrow ... **data integration** ...
- ... and we show the notable advantages that ontologies can bring to this scenario.



Data integration is the problem of providing unified and transparent access to a set of autonomous and heterogeneous sources.

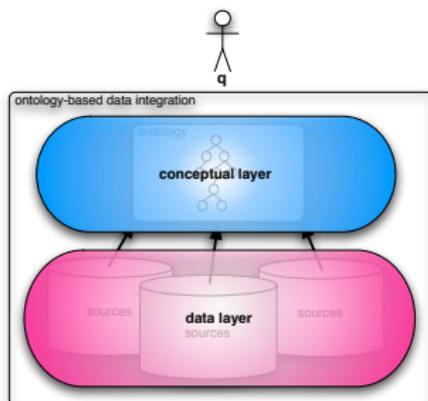
From [Bernstein & Haas, CACM Sept. 2008]:

- Large enterprises spend a great deal of time and money on information integration (e.g., 40% of information-technology shops' budget).
- Market for data integration software estimated to grow from \$2.5 billion in 2007 to \$3.8 billion in 2012 (+8.7% per year)
[IDC. Worldwide Data Integration and Access Software 2008-2012 Forecast. Doc No. 211636 (Apr. 2008)]
- Data integration is a large and growing part of science, engineering, and biomedical computing.



Ontology-based data integration: conceptual layer & data layer

Ontology-based data integration is based on the idea of decoupling information access from data storage.

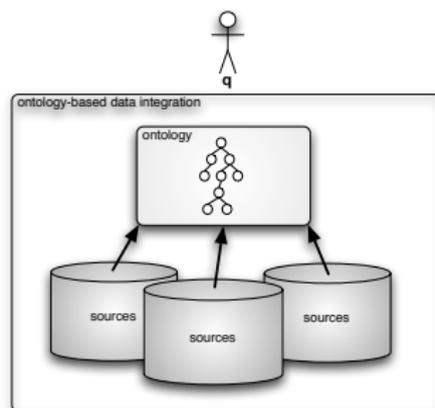


Clients access only the **conceptual layer** ... while the **data layer**, hidden to clients, manages the data.

~> Technological concerns (and changes) on the managed data become fully transparent to the clients.



Ontology-based data integration: architecture



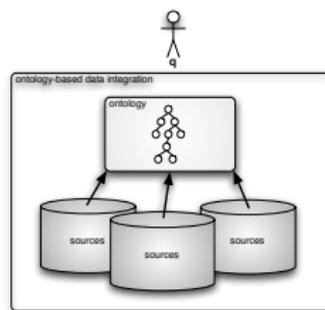
Based on three main components:

- **Ontology**, used as the conceptual layer to give clients a unified conceptual “global view” of the data.
- **Data sources**, these are external, independent, heterogeneous, multiple information systems.
- **Mappings**, which semantically link data at the sources with the ontology (*key issue!*)



Ontology-based data integration: the conceptual layer

The ontology is used as the conceptual layer, to give clients a unified conceptual global view of the data.

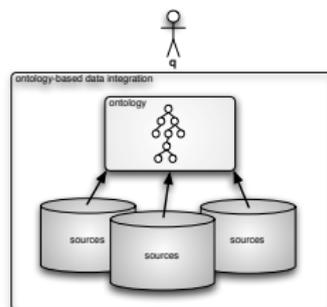


Note: in standard information systems, UML Class Diagram or ER is used at **design time**, ...
... here we use ontologies at **runtime**!



Ontology-based data integration: the sources

Data sources are external, independent, heterogeneous, multiple information systems.



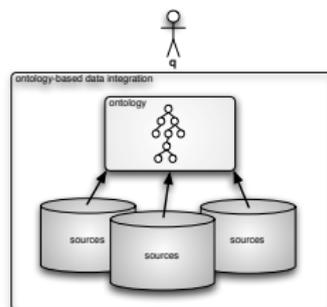
By now we have industrial solutions for:

- Distributed database systems
- Distributed query optimization
- Tools for source wrapping
- Systems for database federation, e.g., IBM Information Integrator
but notice no open-source federated databases yet!



Ontology-based data integration: the sources

Data sources are external, independent, heterogeneous, multiple information systems.



Based on these industrial solutions we can:

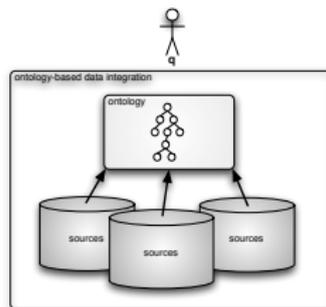
- 1 Wrap the sources and see all of them as relational databases.
- 2 Use federated database tools to see the multiple sources as a single one.

~> We can see the sources as a single (remote) relational database.



Ontology-based data integration: mappings

Mappings semantically link data at the sources with the ontology.



Scientific literature on data integration in databases has shown that ...

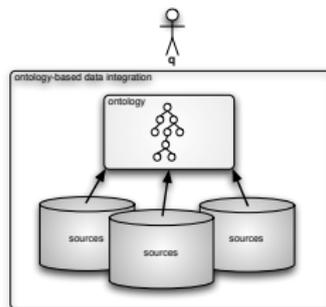
... generally we cannot simply **map** single relations to single elements of the global view (the ontology) ...

... we need to rely on **queries!**



Ontology-based data integration: mappings

Mappings semantically link data at the sources with the ontology.



Several general forms of mappings based on queries have been considered:

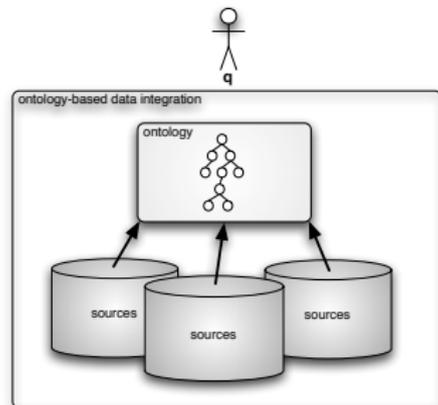
- GAV: map a query over the source to an element in the global view
– *most used form of mappings*
- LAV: map a relation in the source to a query over the global view
– *mathematically elegant, but practically useless (data in the sources are not clean enough!)*
- GLAV: map a query over the sources to a query over the global view
– *the most general form of mappings*

This is a key issue (more on this later).



Ontology-based data integration: incomplete information

It is assumed, even in standard data integration, that the information that the global view has on the data is incomplete!



Important

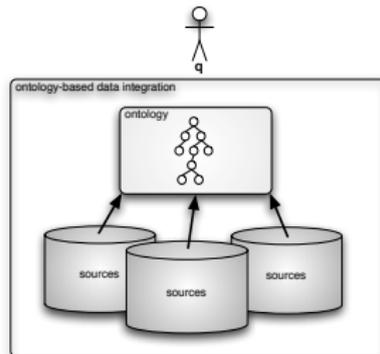
Ontologies are logical theories \leadsto they are perfectly suited to deal with **incomplete information!**



- Query answering amounts to compute **certain answers**, given the global view, the mapping and the data at the sources ...
- ... but query answering may be costly in ontologies (even without mapping and sources).



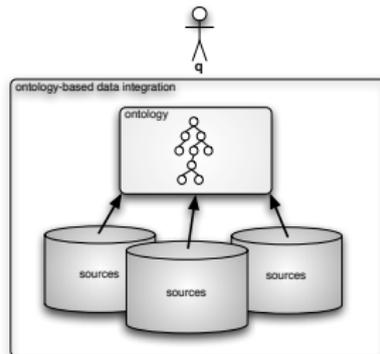
Ontology-based data integration: the $DL-Lite_A$ solution



- We require the data sources to be **wrapped** and presented as relational sources. \rightsquigarrow *“standard technology”*
- We make use of a **data federation tool**, such as IBM Information Integrator, to present the yet to be (semantically) integrated sources as a single relational database. \rightsquigarrow *“standard technology”*
- We make use of the **$DL-Lite_A$** technology presented above for the conceptual view on the data, to **exploit effectiveness of query answering**. \rightsquigarrow *“new technology”*



Ontology-based data integration: the *DL-Lite_A* solution



Are we done? Not yet!

- The (federated) source database is **external** and **independent** from the conceptual view (the ontology).
- **Mappings** relate information in the sources to the ontology. \rightsquigarrow sort of virtual ABox

We use GAV (global-as-view) mappings: the result of an (arbitrary) SQL query on the source database is considered a (partial) extension of a concept/role.

- Moreover, we exploit the distinction between **objects** and **values** in *DL-Lite_A* to deal with the notorious **impedance mismatch problem!**



Impedance mismatch problem

The impedance mismatch problem

- In **relational databases**, information is represented in forms of tuples of **values**.
- In **ontologies** (or more generally object-oriented systems or conceptual models), information is represented using both **objects** and values ...
 - ... with objects playing the main role, ...
 - ... and values a subsidiary role as fillers of object's attributes.

~> *How do we reconcile these views?*

Solution: We need **constructors** to create objects of the ontology out of tuples of values in the database.

Note: from a formal point of view, such constructors can be simply Skolem functions!



Impedance mismatch: the technical solution

Let Γ_V be the alphabet of constants (values) appearing in the sources. We introduce a new alphabet Λ of function symbols, where each function symbol has an associated arity, specifying the number of arguments it accepts.

We inductively define the set $\tau(\Lambda, \Gamma_V)$ of all (Skolem) **terms** of the form $f(d_1, \dots, d_n)$ such that

- $f \in \Lambda$,
- the arity of f is $n > 0$, and
- $d_1, \dots, d_n \in \Gamma_V$.

We use $\tau(\Lambda, \Gamma_V)$ to denote the instances of concepts in the ontology. The unique name assumption is now enforced on such a set.

↪ No confusion between the values stored in the database and the terms denoting objects.



An **ontology with mappings** is characterized by a triple

$\mathcal{O}_m = \langle \mathcal{T}, \mathcal{S}, \mathcal{M} \rangle$ such that:

- \mathcal{T} is a TBox;
- \mathcal{S} is a (federated) relational database representing the sources;
- \mathcal{M} is a set of **mapping assertions**, each one of the form*

$$\Phi(\vec{x}) \rightsquigarrow \Psi(\mathbf{f}(\vec{x}), \vec{x})$$

where

- $\Phi(\vec{x})$ is an arbitrary SQL query over \mathcal{S} , returning attributes \vec{x}
- $\Psi(\mathbf{f}(\vec{x}), \vec{x})$ is a conjunctive query over \mathcal{T} **without non-distinguished variables**, whose variables, possibly occurring in terms, i.e., $\mathbf{f}(\vec{x})$, are from \vec{x} .

* Note: this is a form of GAV mapping



Example

Let \mathcal{S} be the database constituted by a set of relations with the signature:

$D_1[\text{SSN}, \text{PROJ}, \text{D}]$, $D_2[\text{SSN}, \text{NAME}]$, $D_3[\text{CODE}, \text{NAME}]$, $D_4[\text{CODE}, \text{SSN}]$

- Relation D_1 stores tuples (s, p, d) , where s and p are strings and d is a date, such that s is the social security number of a temporary employee, p is the name of the project s/he works for (different projects have different names), and d is the ending date of the employment.
- Relation D_2 stores tuples (s, n) of strings consisting of the social security number s of an employee and her/his name n .
- Relation D_3 stores tuples (c, n) of strings consisting of the code c of a manager and her/his name n .
- Relation D_4 relates managers' code with their social security number.



Example (cont'd)

Consider the ontology with mappings $\mathcal{O}_m = \langle \mathcal{T}, \mathcal{M}, DB \rangle$ such that \mathcal{T} is

TempEmp \sqsubseteq Employee

Manager \sqsubseteq Employee

Employee \sqsubseteq Person

Employee $\sqsubseteq \exists \text{worksFor}$

$\exists \text{worksFor}^- \sqsubseteq \text{Project}$

Person $\sqsubseteq \exists \text{persName}$

(**funct** *persName*)

Project $\sqsubseteq \exists \text{projName}$

(**funct** *projName*)

TempEmp $\sqsubseteq \exists \text{until}$

$\exists \text{until} \sqsubseteq \exists \text{worksFor}$

(**funct** *until*)

Manager $\sqsubseteq \neg \exists \text{until}$

and \mathcal{M} is defined by using $\Lambda = \{\mathbf{pers}, \mathbf{proj}, \mathbf{mgr}\}$, all of which are function symbols of arity 1.



Example (cont'd)

Mapping assertions \mathcal{M} :

- M_{m_1} : SELECT SSN, PROJ, D
FROM D_1 \rightsquigarrow TempEmp(**pers**(SSN)),
worksFor(**pers**(SSN), **proj**(PROJ)),
projName(**proj**(PROJ), PROJ),
until(**pers**(SSN), D)
- M_{m_2} : SELECT SSN, NAME
FROM D_2 \rightsquigarrow Employee(**pers**(SSN)),
persName(**pers**(SSN), NAME)
- M_{m_3} : SELECT SSN, NAME
FROM D_3, D_4 \rightsquigarrow Manager(**pers**(SSN)),
persName(**pers**(SSN), NAME)
WHERE $D_3.CODE = D_4.CODE$
- M_{m_4} : SELECT CODE, NAME
FROM D_3 \rightsquigarrow Manager(**mgr**(CODE)),
persName(**mgr**(CODE), NAME)
WHERE CODE NOT IN
(SELECT CODE FROM D_4)



Def.: Semantics of mappings

We say that \mathcal{I} **satisfies** $\Phi(\vec{x}) \rightsquigarrow \Psi(f(\vec{x}, \vec{x}))$ wrt a database \mathcal{S} , if for every tuple of values \vec{v} such that \vec{v} is the answer of the SQL query $\Phi(\vec{x})$ over \mathcal{S} , and for each ground atom X in $\Psi[f(\vec{v}, \vec{v})]$, we have that:

- if X has the form $A(s)$, then $s^{\mathcal{I}} \in A^{\mathcal{I}}$;
- if X has the form $P(s_1, s_2)$, then $(s_1^{\mathcal{I}}, s_2^{\mathcal{I}}) \in P^{\mathcal{I}}$.

Def.: Semantics of ontologies with mappings

An interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ is a **model** of $\mathcal{O}_m = \langle \mathcal{T}, \mathcal{S}, \mathcal{M} \rangle$ if:

- \mathcal{I} is a model of \mathcal{T} ;
- \mathcal{I} satisfies \mathcal{M} w.r.t. \mathcal{S} , i.e., satisfies every assertion in \mathcal{M} w.r.t. \mathcal{S} .

An ontology with mappings is **satisfiable** if it admits at least one model.



Given a (U)CQ q and $\mathcal{O} = \langle \mathcal{T}, \mathcal{S}, \mathcal{M} \rangle$ (assumed to be satisfiable), we compute **certain answers** $cert(q, \mathcal{O})$ as follows:

- 1 Using \mathcal{T} , **reformulate** CQ q as a union $r_{q, \mathcal{T}}$ of CQs.
- 2 Using \mathcal{M} , **unfold** $r_{q, \mathcal{T}}$ to obtain a union $unfold(r_{q, \mathcal{T}})$ of CQs.
- 3 **Evaluate** $unfold(r_{q, \mathcal{T}})$ directly over \mathcal{S} using RDBMS technology.

Correctness of this algorithm shows FOL-reducibility of query answering.

↪ Query answering can again be done using **RDBMS technology**.

↪ Prototype system implemented: **Quonto+Integration Module**



\mathcal{M} can be encoded in the following portion of a logic program:

<code>TempEmp(pers(<i>s</i>))</code>	<code>← Aux₁₁(<i>s</i>)</code>
<code>worksFor(pers(<i>s</i>), proj(<i>p</i>))</code>	<code>← Aux₁₂(<i>s</i>, <i>p</i>)</code>
<code>projName(proj(<i>p</i>), <i>p</i>)</code>	<code>← Aux₁₃(<i>p</i>)</code>
<code>until(pers(<i>s</i>), <i>d</i>)</code>	<code>← Aux₁₄(<i>s</i>, <i>d</i>)</code>
<code>Employee(pers(<i>s</i>))</code>	<code>← Aux₂₁(<i>s</i>)</code>
<code>persName(pers(<i>s</i>), <i>n</i>)</code>	<code>← Aux₂₂(<i>s</i>, <i>n</i>)</code>
<code>Manager(pers(<i>s</i>))</code>	<code>← Aux₃₁(<i>s</i>)</code>
<code>persName(pers(<i>s</i>), <i>n</i>)</code>	<code>← Aux₃₂(<i>s</i>, <i>n</i>)</code>
<code>Manager(mgr(<i>c</i>))</code>	<code>← Aux₄₁(<i>c</i>)</code>
<code>persName(mgr(<i>c</i>), <i>n</i>)</code>	<code>← Aux₄₂(<i>c</i>, <i>n</i>)</code>

where Aux_{ij} is a predicate denoting the result of the evaluation over \mathcal{S} of the query $\Phi_{m_{ij}}$ in the left-hand side of the mapping $M_{m_{ij}}$.



Example (cont'd)

Consider $q(x) \leftarrow \text{worksFor}(x, y)$, whose reformulation $Q' = r_{q, \mathcal{T}}$ is:

$$Q'(x) \leftarrow \text{worksFor}(x, y)$$

$$Q'(x) \leftarrow \text{until}(x, y)$$

$$Q'(x) \leftarrow \text{TempEmp}(x)$$

$$Q'(x) \leftarrow \text{Employee}(x)$$

$$Q'(x) \leftarrow \text{Manager}(x)$$

To compute the unfolding of Q' , we unify each of its atoms with the left-hand side of the logic program rules corresponding to the mapping assertions in \mathcal{M} , and we obtain the following *partial evaluation* of Q' :

$$q(\mathbf{pers}(s)) \leftarrow \text{Aux}_{12}(s, p)$$

$$q(\mathbf{pers}(s)) \leftarrow \text{Aux}_{14}(s, d)$$

$$q(\mathbf{pers}(s)) \leftarrow \text{Aux}_{11}(s)$$

$$q(\mathbf{pers}(s)) \leftarrow \text{Aux}_{21}(s)$$

$$q(\mathbf{pers}(s)) \leftarrow \text{Aux}_{31}(s, n)$$

$$q(\mathbf{mgr}(c)) \leftarrow \text{Aux}_{41}(c, n)$$



Example (cont'd)

From the above formulation, it is now possible to derive the corresponding SQL query Q'' that can be directly issued over the database \mathcal{S} :

```
SELECT concat(concat('pers (' ,SSN),'))  
FROM  $D_1$   
UNION  
SELECT concat(concat('pers (' ,SSN),'))  
FROM  $D_2$   
UNION  
SELECT concat(concat('pers (' ,SSN),'))  
FROM  $D_3, D_4$   
WHERE  $D_3.CODE=D_4.CODE$   
UNION  
SELECT concat(concat('mgr (' ,CODE),'))  
FROM  $D_3$   
WHERE CODE NOT IN (SELECT CODE FROM  $D_4$ )
```



Theorem

Query answering in a $DL-Lite_{\mathcal{A}}$ ontology with mappings $\mathcal{O} = \langle \mathcal{T}, \mathcal{S}, \mathcal{M} \rangle$ is

- 1 **NP-complete** in the size of the query.
- 2 **PTime** in the size of the **TBox** \mathcal{T} and the **mappings** \mathcal{M} .
- 3 **LogSpace** in the size of the **database** \mathcal{S} , in fact FOL-rewritable.

*Can we move to LAV or GLAV mappings?
No, if we want to stay in LOGSPACE.*



People involved in this work:

- Domenico Lembo
- Maurizio Lenzerini
- Marco Ruzzi
- Antonella Poggi
- Mariano Rodriguez Muro
- Riccardo Rosati
- Several master and PhD students, including the ones currently active:
 - Claudio Corona
 - Emma Di Pasquale
 - Fabio Savo



- 1 Introduction
- 2 Querying data through ontologies
- 3 *DL-Lite_A*: an ontology language for accessing data
- 4 Ontology-based data integration
- 5 References



- [1] D. Berardi, D. Calvanese, and G. De Giacomo.
Reasoning on UML class diagrams.
Artificial Intelligence, 168(1–2):70–118, 2005.
- [2] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, A. Poggi, and R. Rosati.

Linking data to ontologies: The description logic $DL-Lite_A$.

In *Proc. of the 2nd Int. Workshop on OWL: Experiences and Directions (OWLED 2006)*, volume 216 of *CEUR Electronic Workshop Proceedings*, <http://ceur-ws.org/Vol-216/>, 2006.

- [3] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati.
Tailoring OWL for data intensive ontologies.
In *Proc. of the 1st Int. Workshop on OWL: Experiences and Directions (OWLED 2005)*, volume 188 of *CEUR Electronic Workshop Proceedings*, <http://ceur-ws.org/Vol-188/>, 2005.



- [4] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati.
DL-Lite: Tractable description logics for ontologies.
In Proc. of the 20th Nat. Conf. on Artificial Intelligence (AAAI 2005), pages 602–607, 2005.
- [5] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati.
Data complexity of query answering in description logics.
In Proc. of the 10th Int. Conf. on the Principles of Knowledge Representation and Reasoning (KR 2006), pages 260–270, 2006.
- [6] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati.
Tractable reasoning and efficient query answering in description logics: The *DL-Lite* family.
J. of Automated Reasoning, 39(3):385–429, 2007.
- [7] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati.
Path-based identification constraints in description logics.
In Proc. of the 11th Int. Conf. on the Principles of Knowledge Representation and Reasoning (KR 2008), pages 231–241, 2008.



- [8] D. Calvanese and M. Rodríguez.
An extension of DIG 2.0 for handling bulk data.
In *Proc. of the 3rd Int. Workshop on OWL: Experiences and Directions (OWLED 2007)*, volume 258 of *CEUR Electronic Workshop Proceedings*, <http://ceur-ws.org/Vol-258/>, 2007.
- [9] A. Poggi, D. Lembo, D. Calvanese, G. De Giacomo, M. Lenzerini, and R. Rosati.
Linking data to ontologies.
J. on Data Semantics, X:133–173, 2008.
- [10] A. Poggi, M. Rodriguez, and M. Ruzzi.
Ontology-based database access with DIG-Mastro and the OBDA Plugin for Protégé.
In K. Clark and P. F. Patel-Schneider, editors, *Proc. of the OWL: Experiences and Directions 2008 (OWLED 2008 DC) Workshop*, 2008.



[11] M. Rodriguez-Muro, L. Lubyte, and D. Calvanese.

Realizing ontology based data access: A plug-in for Protégé.

In *Proc. of the 24th Int. Conf. on Data Engineering Workshops (ICDE 2008)*, pages 286–289, 2008.

