# ONTOLOGY-MEDIATED QUERY ANSWERING

# Harnessing Knowledge to Get More From Data

Meghyn Bienvenu (CNRS, University of Montpellier, Inria)

# ONTOLOGY-MEDIATED QUERY ANSWERING (OMQA)



## ONTOLOGY-MEDIATED QUERY ANSWERING (OMQA)



# In computer science:

a formal specification of the knowledge of a particular domain, thereby making it amenable to machine processing

Such a specification consists of:

- $\cdot$  terminology (or vocabulary) of the domain
- · semantic relationships between terms
  - · relations of inclusion, equivalence, disjointness, ...

To standardize the terminology of an application domain

- · meaning of terms is constrained, so fewer misunderstandings
- · by adopting a common vocabulary, easy to share information

To present an intuitive and unified view of data sources

- ontology can be used to enrich the data vocabulary, making it easier for users to formulate their queries
- $\cdot$  especially useful when integrating multiple data sources

#### To support automated reasoning

- · uncover implicit connections between terms, errors in modelling
- exploit knowledge in the ontology during query answering, to get back a more complete set of answers to queries

General medical ontologies: SNOMED CT ( $\sim$  300,000 terms!), GALEN Specialized ontologies: FMA (anatomy), NCI (cancer), ...



Querying & exchanging medical records (find patients for medical trials)

· myocardial infarction vs. MI vs. heart attack vs. 410.0

Supports tools for annotating and visualizing patient data (scans, x-rays)

## **Hundreds of ontologies** at BioPortal (http://bioportal.bioontology.org/): Gene Ontology (GO), Cell Ontology, Pathway Ontology, Plant Anatomy, ...



Help scientists share, query, & visualize experimental data

#### APPLICATIONS OF OMQA: ENTREPRISE INFORMATION SYSTEMS

Companies and organizations have lots of data

need easy and flexible access to support decision-making



Example industrial projects:

- · Public debt data: Sapienza Univ. & Italian Department of Treasury
- · Energy sector: Optique EU project (several univ, StatOil, Siemens)

## Description logics (DLs):

- · popular means for specifying ontologies
- $\cdot$  basis of the web ontology language OWL (W3C standard)

Formally: decidable fragments of first-order logic

- $\cdot\,$  inherit well-defined semantics
- $\cdot$  succinct, variable-free syntax

## Description logics (DLs):

- · popular means for specifying ontologies
- $\cdot$  basis of the web ontology language OWL (W3C standard)

Formally: decidable fragments of first-order logic

- $\cdot\,$  inherit well-defined semantics
- $\cdot$  succinct, variable-free syntax

Computational properties well understood (decidability, complexity)

Many implemented reasoners and tools available for use

# Introduction to DLs & OMQA

# Query Rewriting: Limits and Possibilities

Inconsistency Handling in OMQA

# INTRODUCTION TO DLS & OMQA

## Building blocks:

- concept names (unary predicates, classes)
- · role names (binary predicates, properties)

Faculty □ ¬Prof

∃Teaches.GradCourse

Prof Fellow Course Teaches HeadOf

Teaches<sup>-</sup>

## Building blocks:

- **concept names** (unary predicates, classes)
- · role names (binary predicates, properties)



Constructors to build complex descriptions  $\Box, \Box, \neg, \forall, \exists, ...$ 

Faculty □ ¬Prof ∃Tead

∃Teaches.GradCourse

Teaches<sup>—</sup>

### Ontology = set of axioms

concept inclusions

 Prof ⊑ Faculty
 Prof ⊑ ¬Fellow
 ∃Teaches.GradCourse ⊑ Prof

 · role inclusions
 TaughtBy ⊑ Teaches<sup>−</sup>
 HeadOf ⊑ MemberOf

Note: allowed constructors and axioms depends on chosen DL

## Instance queries (IQs): find instances of a given concept or role



## Instance queries (IQs): find instances of a given concept or role

# Faculty(x) Teaches(x, y)

**Conjunctive queries (CQs)** ~ SPJ queries in SQL, BGPs in SPARQL conjunctions of atoms, some variables can be existentially quantified

 $\exists y. Faculty(x) \land Teaches(x, y)$ 

(find all faculty members that teach something)

## Instance queries (IQs): find instances of a given concept or role

# Faculty(x) Teaches(x, y)

**Conjunctive queries (CQs)** ~ SPJ queries in SQL, BGPs in SPARQL conjunctions of atoms, some variables can be existentially quantified

 $\exists y. Faculty(x) \land Teaches(x, y)$ 

(find all faculty members that teach something)

Unions of conjunctive queries (UCQs): disjunction of CQs

#### ONTOLOGY-MEDIATED QUERY ANSWERING

#### Answering CQs in database setting



#### **ONTOLOGY-MEDIATED QUERY ANSWERING**

#### Answering CQs in database setting



#### Answering CQs in the presence of an ontology



Ontology, expressed in DL-Lite: (lightweight DL designed for OMQA)  $$\sim$ OWL 2 QL$$ 

Prof ⊑ FacultyFellow ⊑ FacultyProf ⊑ ¬FellowProf ⊑ ∃Teaches∃Teaches ⊑ Faculty∃Teaches ⊑ Course

Dataset:

 $\mathcal{D}_1 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}) \}$ 

Query:  $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Ontology, expressed in DL-Lite: (lightweight DL designed for OMQA)  $$\sim$ OWL 2 QL$$ 

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

Dataset:

 $\mathcal{D}_1 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}) \}$ 

Query:  $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Get the following answers:

anna Prof(anna) + Prof ⊑ Faculty + Prof ⊑ ∃Teaches
 tom Fellow(tom) + Fellow ⊑ Faculty + Teaches(tom, cs101)

# QUERY REWRITING: LIMITS AND POSSIBILITIES

Idea: reduce OMQA to database query evaluation

· rewriting step: ontology  $\mathcal{O}$  + query  $q \rightsquigarrow$  first-order (FO) query q'

FO queries  $\sim$  SQL queries

· evaluation step: evaluate query q' over dataset

Advantage: harness efficiency of relational database (DB) systems

Idea: reduce OMQA to database query evaluation

· rewriting step: ontology  $\mathcal{O}$  + query  $q \rightsquigarrow$  first-order (FO) query q'

FO queries  $\sim$  SQL queries

· evaluation step: evaluate query q' over dataset

Advantage: harness efficiency of relational database (DB) systems

**FO-rewriting** of *q* w.r.t.  $\mathcal{O}$ : **FO-query** *q'* such that **for every dataset**  $\mathcal{D}$ :

evaluating q' over  $\mathcal{D}$  (viewed as DB) gives correct result

Same ontology, data, and query as earlier:

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

 $\mathcal{D}_1 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Same ontology, data, and query as earlier:

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

 $\mathcal{D}_1 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

FO-rewriting of  $q_1(x)$  w.r.t.  $\mathcal{O}$ : (disjuncts = different ways to satisfy  $q_1$ ) (Faculty(x)  $\land \exists y$ .Teaches(x, y))  $\lor$  (Fellow(x)  $\land \exists y$ .Teaches(x, y))  $\lor$  Prof(x) Same ontology, data, and guery as earlier:

Prof  $\Box$  Faculty Fellow  $\Box$  Faculty Prof  $\Box$   $\neg$  Fellow  $Prof \Box \exists Teaches \exists Teaches \Box Faculty \exists Teaches^{-} \Box Course$ 

 $\mathcal{D}_1 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

FO-rewriting of  $q_1(x)$  w.r.t.  $\mathcal{O}$ : (disjuncts = different ways to satisfy  $q_1$ )  $(Faculty(x) \land \exists y.Teaches(x, y)) \lor (Fellow(x) \land \exists y.Teaches(x, y)) \lor Prof(x)$ 

Evaluating the rewriting over  $\mathcal{D}_1$  yields:

anna (matches 3rd disjunct) tom (matches 2nd disjunct)

Lots of **implemented rewriting algorithms** for DL-Lite

• many produce rewritings in the form of UCQs (like in our example)

Lots of implemented rewriting algorithms for DL-Lite

• many produce rewritings in the form of UCQs (like in our example)

Experiments showed that such rewritings can be huge!

 $\cdot\,$  can be difficult / impossible to generate and evaluate

Lots of implemented rewriting algorithms for DL-Lite

• many produce rewritings in the form of UCQs (like in our example)

Experiments showed that such rewritings can be huge!

 $\cdot$  can be difficult / impossible to generate and evaluate

Easy to show smallest UCQ-rewriting may be exponentially large:

Lots of implemented rewriting algorithms for DL-Lite

• many produce rewritings in the form of UCQs (like in our example)

Experiments showed that such rewritings can be huge!

· can be difficult / impossible to generate and evaluate

Easy to show smallest UCQ-rewriting may be exponentially large: Query:  $A_1^0(x) \land \ldots \land A_n^0(x)$  Ontology:  $A_i^1 \sqsubseteq A_i^0$   $(i = 1, \ldots, n)$ Rewriting:  $\bigvee_{(i_1, \ldots, i_n) \in \{0, 1\}} A_1^{i_1}(x) \land A_1^{i_1}(x) \land \ldots \land A_1^{i_1}(x)$ 

Lots of implemented rewriting algorithms for DL-Lite

• many produce rewritings in the form of UCQs (like in our example)

Experiments showed that such rewritings can be huge!

· can be difficult / impossible to generate and evaluate

Easy to show smallest UCQ-rewriting may be exponentially large: Query:  $A_1^0(x) \land \ldots \land A_n^0(x)$  Ontology:  $A_i^1 \sqsubseteq A_i^0$   $(i = 1, \ldots, n)$ Rewriting:  $\bigvee_{(i_1, \ldots, i_n) \in \{0, 1\}} A_1^{i_1}(x) \land A_1^{i_1}(x) \land \ldots \land A_1^{i_1}(x)$ 

But equivalent to polysize query  $\bigwedge_{i=1}^{n} (A_{i}^{0}(x) \lor A_{i}^{1}(x))$  (DNF vs CNF)

Different shapes of rewritings:

- · UCQs
- · positive existential (PE) queries
- non-recursive datalog (NDL) queries
- · first-order (FO) queries

disjunction of conjunctions can mix  $\lor$  and  $\land$ allows structure sharing can also use  $\neg, \forall$  Different shapes of rewritings:

- · UCQs
- · positive existential (PE) queries
- non-recursive datalog (NDL) queries
- first-order (FO) queries

sjunction of conjunctions can mix ∨ and ∧ allows structure sharing can also use ¬, ∀

Question: Do we always have polysize (PE / NDL / FO) rewriting?

Different shapes of rewritings:

- · UCQs
- · positive existential (PE) queries
- non-recursive datalog (NDL) queries
- · first-order (FO) queries

disjunction of conjunctions can mix  $\lor$  and  $\land$ allows structure sharing can also use  $\neg$ ,  $\forall$ 

Question: Do we always have polysize (PE / NDL / FO) rewriting? no When are polysize (PE / NDL / FO) rewritings possible?

## Natural restrictions:

- · ontology: bound existential depth
- query: bound treewidth / number of leaves (tree-shaped queries)

#### SUCCINCTNESS LANDSCAPE



Key technique: link rewriting size to circuit complexity

For most other ontology languages: FO-rewritings may not exist!

Take for example the lightweight DL *EL* basis for OWL 2 EL

no FO-rewriting of A(x) w.r.t.  $\mathcal{O} = \{\exists R.A \sqsubseteq A\}$ 

For most other ontology languages: FO-rewritings may not exist!

Take for example the **lightweight DL**  $\mathcal{EL}$  basis for OWL 2 EL no FO-rewriting of A(x) w.r.t.  $\mathcal{O} = \{\exists R.A \sqsubseteq A\}$ 

Hope: FO-rewritings do exist for typical queries and ontologies

Question: how to identify these good cases?(extend applicability)• does query q have an FO-rewriting w.r.t. ontology O?

For most other ontology languages: FO-rewritings may not exist!

Take for example the **lightweight DL**  $\mathcal{EL}$  basis for OWL 2 EL no FO-rewriting of A(x) w.r.t.  $\mathcal{O} = \{\exists R.A \sqsubseteq A\}$ 

Hope: FO-rewritings do exist for typical queries and ontologies

Question: how to **identify these good cases**? (extend applicability) • does query *q* have an FO-rewriting w.r.t. ontology *O*?

Tackled in series of recent papers

(incl: IJCAI'16 paper)

- $\cdot\,$  decision procedures for variety of DLs, query languages
- · practical methods for building rewritings promising first results

## INCONSISTENCY HANDLING IN OMQA

In realistic settings, can expect some errors in the data

 $\cdot$  data likely to be **inconsistent** with the ontology

Standard semantics: everything is implied - not informative!

In realistic settings, can expect some errors in the data

 $\cdot$  data likely to be **inconsistent** with the ontology

Standard semantics: everything is implied - not informative!

Two approaches to inconsistency handling:

- resolve the inconsistencies
  - · preferable, but not always applicable!
- $\cdot\,$  live with the inconsistencies adopt alternative semantics
  - meaningful answers to queries despite inconsistencies

### **Repair**: $\subseteq$ -maximal subset of the data consistent with the ontology

 $\cdot$  ways to achieve consistency, keeping as much information as possible

#### **Repair**: ⊆-maximal subset of the data consistent with the ontology

 $\cdot$  ways to achieve consistency, keeping as much information as possible

Plausible answers: hold no matter which repair is chosen

**Repair**: ⊆-maximal subset of the data consistent with the ontology

 $\cdot$  ways to achieve consistency, keeping as much information as possible

Plausible answers: hold no matter which repair is chosen

AR semantics: query each repair separately, intersect results



Same ontology and query as before, add Prof(tom) to dataset:

Prof⊑ Faculty	Fellow ⊑ Faculty	Prof⊑ ¬Fellow
Prof ⊑ ∃Teaches	$\exists Teaches \sqsubseteq Faculty$	$\exists Teaches^- \sqsubseteq Course$

 $D_2 = \{ Prof(anna), Fellow(tom), Teaches(tom, cs101), Prof(tom) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Same ontology and query as before, add Prof(tom) to dataset:

Prof⊑ Faculty	Fellow ⊑ Faculty	Prof⊑ ¬Fellow
Prof⊑∃Teaches	$\exists Teaches \sqsubseteq Faculty$	∃Teaches <sup>–</sup> ⊑ Course

 $\mathcal{D}_2 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}), \mathsf{Prof}(\mathsf{tom}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Two repairs of  $\mathcal{D}_2$  w.r.t.  $\mathcal{O}$ :

 $\mathcal{R}_1 = \{ Prof(anna), Fellow(tom), Teaches(tom, cs101) \} \}$ 

drop Prof(tom)

 $\mathcal{R}_2 = \{ Prof(anna), Prof(tom), Teaches(tom, cs101) \}$ 

drop Fellow(tom)

Same ontology and query as before, **add** Prof(tom) to dataset:

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

 $\mathcal{D}_2 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}), \mathsf{Prof}(\mathsf{tom}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Two repairs of  $\mathcal{D}_2$  w.r.t.  $\mathcal{O}$ :

 $\mathcal{R}_1 = \{ Prof(anna), Fellow(tom), Teaches(tom, cs101) \} \}$ 

drop Prof(tom)

 $\mathcal{R}_2 = \{ Prof(anna), Prof(tom), Teaches(tom, cs101) \}$ 

drop Fellow(tom)

Under AR semantics:

- · anna and tom are both answers to  $q_1$
- · cs101 is not an answer

**Repair**: ⊆-maximal subset of the data consistent with the ontology

 $\cdot$  ways to achieve consistency, keeping as much information as possible

Plausible answers: hold no matter which repair is chosen

AR semantics: query each repair separately, intersect results



Bad news: query answering under AR semantics is intractable (coNP-hard in the size of the data)

Worse: intractable even in very restricted settings ( $\mathcal{O} = \{A \sqsubseteq \neg B\}$ )

Brave semantics

possible answers

· answer required to hold w.r.t. at least one repair

**Brave semantics** 

possible answers

· answer required to hold w.r.t. at least one repair

**IAR semantics** 

surest answers

• query the intersection of all repairs (surest facts)

#### **Brave semantics**

· answer required to hold w.r.t. at least one repair

#### **IAR semantics**

surest answers

possible answers

· query the intersection of all repairs (surest facts)

Relationship between the semantics:

IAR answers $\subseteq$ AR answers $\subseteq$ brave answers

#### **Brave semantics**

· answer required to hold w.r.t. at least one repair

#### **IAR semantics**

· query the intersection of all repairs (surest facts)

Relationship between the semantics:

IAR answers $\subseteq$ AR answers $\subseteq$ brave answers

Good news: these semantics are tractable for DL-Lite ontologies

possible answers

surest answers

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

 $\mathcal{D}_2 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}), \mathsf{Prof}(\mathsf{tom}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Repairs of  $\mathcal{D}_2$  w.r.t.  $\mathcal{O}$ :

 $\mathcal{R}_1 = \{ Prof(anna), Fellow(tom), Teaches(tom, cs101) \} \}$ 

 $\mathcal{R}_2 = \{ Prof(anna), Prof(tom), Teaches(tom, cs101) \}$ 

Intersection of repairs: {Prof(anna), Teaches(tom, cs101)}

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

 $\mathcal{D}_2 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}), \mathsf{Prof}(\mathsf{tom}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Repairs of  $\mathcal{D}_2$  w.r.t.  $\mathcal{O}$ :

 $\mathcal{R}_1 = \{ Prof(anna), Fellow(tom), Teaches(tom, cs101) \} \}$ 

 $\mathcal{R}_2 = \{ Prof(anna), Prof(tom), Teaches(tom, cs101) \}$ 

Intersection of repairs: {Prof(anna), Teaches(tom, cs101)}

Brave answers: anna tom

Prof  $\sqsubseteq$  FacultyFellow  $\sqsubseteq$  FacultyProf  $\sqsubseteq$  ¬FellowProf  $\sqsubseteq$  3Teaches  $\sqsubseteq$  Faculty3Teaches  $\neg$  Course

 $\mathcal{D}_2 = \{ \mathsf{Prof}(\mathsf{anna}), \mathsf{Fellow}(\mathsf{tom}), \mathsf{Teaches}(\mathsf{tom}, \mathsf{cs101}), \mathsf{Prof}(\mathsf{tom}) \}$ 

 $q_1(x) = \exists y. Faculty(x) \land Teaches(x, y)$ 

Repairs of  $\mathcal{D}_2$  w.r.t.  $\mathcal{O}$ :

 $\mathcal{R}_1 = \{ Prof(anna), Fellow(tom), Teaches(tom, cs101) \} \}$ 

 $\mathcal{R}_2 = \{ Prof(anna), Prof(tom), Teaches(tom, cs101) \}$ 

Intersection of repairs: {Prof(anna), Teaches(tom, cs101)}

Brave answers: anna tom IAR answers: anna

CQAPri first system for AR query answering in DL-Lite

Implements hybrid approach:

- compute IAR and brave answers
  - · gives upper and lower **bounds on AR answers**
- · use SAT solvers to identify remaining AR answers
- three categories of answers : possible, likely, (almost) sure

#### Interaction with user

- explaining query results
  - why a possible answer? why not a sure answer?
- · query-driven repairing
  - · exploit user feedback to improve data quality

(IJCAI'16)

polytime

# CONCLUSION & OUTLOOK

## Promising approach, important applications

Significant recent advances, but still lots left to do!

- efficiency: more expressive ontology and query languages
- · robustness: inconsistencies, uncertainty, vagueness, exceptions
- usability: ontology construction, explanation, query formulation, ...

Connections to other disciplines: databases, Semantic Web, theory

#### Collaboration with other AI areas?

- natural language processing
- $\cdot$  machine learning

# QUESTIONS ?

SEE PAPER FOR REFERENCES & POINTERS TO THE LITERATURE

BASED ON JOINT WORK WITH:

Camille Bourgaux, Balder ten Cate, François Goasdoué, Peter Hansen, Carsten Lutz, Stanislav Kikot, Roman Kontchakov, Magdalena Ortiz, Vladimir Podolskii, Riccardo Rosati, Mantas Šimkus, Frank Wolter, Guohui Xiao, Michael Zakharyschev