

Answering Queries in Description Logics: Theory and Applications to Data Management

Diego Calvanese¹, Michael Zakharyashev²

¹ KRDB Research Centre
Free University of Bozen-Bolzano

² Birbeck College, London

ESLLI 2010, August 14–20, 2010
Copenhagen, Denmark

Overview of the Course

- 1 Introduction and background
 - 1 Ontology-based data management
 - 2 Brief introduction to computational complexity
 - 3 Query answering in databases
 - 4 Querying databases and ontologies
- 2 Lightweight description logics
 - 5 Introduction to description logics
 - 6 DLs for conceptual data modeling: the *DL-Lite* family
 - 7 The \mathcal{EL} family of tractable description logics
- 3 Query answering in the *DL-Lite* family
 - 8 Query answering in description logics
 - 9 Lower bounds for more expressive description logics
 - 10 Query answering by rewriting
- 4 The combined approach to query answering
 - 11 Query answering in *DL-Lite*: data completion
 - 12 Query rewriting in \mathcal{EL}
- 5 Linking ontologies to relational data
 - 13 The impedance mismatch problem
 - 14 Query answering in Ontology-Based Data Access systems
- 6 Conclusions and references

Outline of Lecture 1

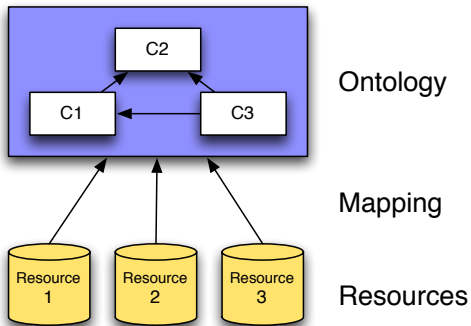
- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References



Outline of Lecture 1

- 1 **Ontology-based data management**
 - Challenges in data management
 - Description logics and ontologies for data management
 - Requirements on description logics for data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References

Ontologies at the core of information systems



The usage of all system resources (data and services) is done through the domain conceptualization.

Issue 3: Tools

- According to the principle that “there is no meaning without a language with a formal semantics”, the formal semantics becomes the solid basis for dealing with ontologies.
- Hence every kind of access to an ontology (to extract information, to modify it, etc.), requires to **fully** take into account its semantics.
- We need tools that perform reasoning over the ontology that is **sound and complete** wrt the semantics.
- The tools have to be “efficient”, especially wrt the size of the data.

In this course:

- We discuss the requirements, the principles, and the theoretical foundations for ontology-based data access tools.
- We briefly present a tool for querying data sources through ontologies that has been built according to those principles.

Outline of Lecture 1

- 1 **Ontology-based data management**
 - Challenges in data management
 - Description logics and ontologies for data management
 - **Requirements on description logics for data management**
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References



Which is the “right” expressive power?

What should an ontology language / description logic be able to express in order to be well suited for data management applications?

Let's start with an exercise!

Exercise

Requirements: We are interested in building a software application to manage filmed scenes for realizing a movie, by following the so-called “Hollywood Approach”.

Every **scene** is identified by a code (a string) and is described by a text in natural language.

Every scene is filmed from different positions (at least one), each of this is called a **setup**. Every setup is characterized by a code (a string) and a text in natural language where the photographic parameters are noted (e.g., aperture, exposure, focal length, filters, etc.). Note that a setup is related to a single scene.

For every setup, several **takes** may be filmed (at least one). Every take is characterized by a (positive) natural number, a real number representing the number of meters of film that have been used for shooting the take, and the code (a string) of the reel where the film is stored. Note that a take is associated to a single setup.

Scenes are divided into **internals** that are filmed in a theater, and **externals** that are filmed in a **location** and can either be “day scene” or “night scene”. Locations are characterized by a code (a string) and the address of the location, and a text describing them in natural language.


Write a precise specification of this domain using any formalism you like!

What do we need to express?

The formalism we use should allow us to express the following:

- domain partitioned into **classes** of objects (e.g., **Scene**, **Location**, ...)
- objects belonging to a class have specific (local) **properties** (e.g., **code** and **text** for a scene, ...)
- **relationships** between objects (e.g., scenes are **filmedFrom** setups, ...)
- **inclusions** and **hierarchies** between classes (e.g., **Internal** and **External Scenes**)
- **domain** and **range** of relations (e.g., the relation **filmedFrom** has **Scene** as domain and **Setup** as range)
- **mandatory participation** to relations (e.g., every Scene is filmedFrom some Setup)
- **functionality** of relations and attributes (e.g., every Setup is for at most one Scene), and more generally, **numeric constraints**

In addition, we may require:

- **inclusions** and **hierarchies** between relationships
- additional **properties of relationships**, such as transitivity, symmetry,  **unibz.it**

Solution 1: Use logic!!!

Alphabet: $Scene(x)$, $Setup(x)$, $Take(x)$, $Internal(x)$, $External(x)$, $Location(x)$,
 $filmedFrom(x, y)$, $tkOfStp(x, y)$, $located(x, y)$,

Axioms:

$\forall x, y. code_{Scene}(x, y) \rightarrow Scene(x) \wedge String(y)$
 $\forall x, y. description(x, y) \rightarrow Scene(x) \wedge Text(y)$
 $\forall x, y. code_{Setup}(x, y) \rightarrow Setup(x) \wedge String(y)$
 $\forall x, y. photographicPars(x, y) \rightarrow Setup(x) \wedge Text(y)$
 $\forall x, y. nbr(x, y) \rightarrow Take(x) \wedge Integer(y)$
 $\forall x, y. filmedMeters(x, y) \rightarrow Take(x) \wedge Real(y)$
 $\forall x, y. reel(x, y) \rightarrow Take(x) \wedge String(y)$
 $\forall x, y. theater(x, y) \rightarrow Internal(x) \wedge String(y)$
 $\forall x, y. nightScene(x, y) \rightarrow External(x) \wedge Boolean(y)$
 $\forall x, y. name(x, y) \rightarrow Location(x) \wedge String(y)$
 $\forall x, y. address(x, y) \rightarrow Location(x) \wedge String(y)$
 $\forall x, y. description(x, y) \rightarrow Location(x) \wedge Text(y)$
 $\forall x. Scene(x) \rightarrow (1 \leq \#\{y \mid code_{Scene}(x, y)\} \leq 1)$
 $\forall x. Internal(x) \rightarrow Scene(x)$
 $\forall x. External(x) \rightarrow Scene(x)$
 $\forall x. Internal(x) \rightarrow \neg External(x)$
 $\forall x. Scene(x) \rightarrow Internal(x) \vee External(x)$

$\forall x, y. filmedFrom(x, y) \rightarrow$
 $Scene(x) \wedge Setup(y)$
 $\forall x, y. tkOfStp(x, y) \rightarrow$
 $Take(x) \wedge Setup(y)$
 $\forall x, y. located(x, y) \rightarrow$
 $External(x) \wedge Location(y)$
 $\forall x. Scene(x) \rightarrow$
 $(1 \leq \#\{y \mid filmedFrom(x, y)\})$
 $\forall y. Setup(y) \rightarrow$
 $(1 \leq \#\{x \mid filmedFrom(x, y)\} \leq 1)$
 $\forall x. Take(x) \rightarrow$
 $(1 \leq \#\{y \mid tkOfStp(x, y)\} \leq 1)$
 $\forall x. Setup(y) \rightarrow$
 $(1 \leq \#\{x \mid tkOfStp(x, y)\})$
 $\forall x. External(x) \rightarrow$
 $(1 \leq \#\{y \mid located(x, y)\} \leq 1)$
 . . .

Solution 1: Use logic – Discussion

Good points:

- Precise semantics.
- Formal verification.
- Allows for query answering.
- Machine comprehensible.
- Virtually unlimited expressiveness (*).

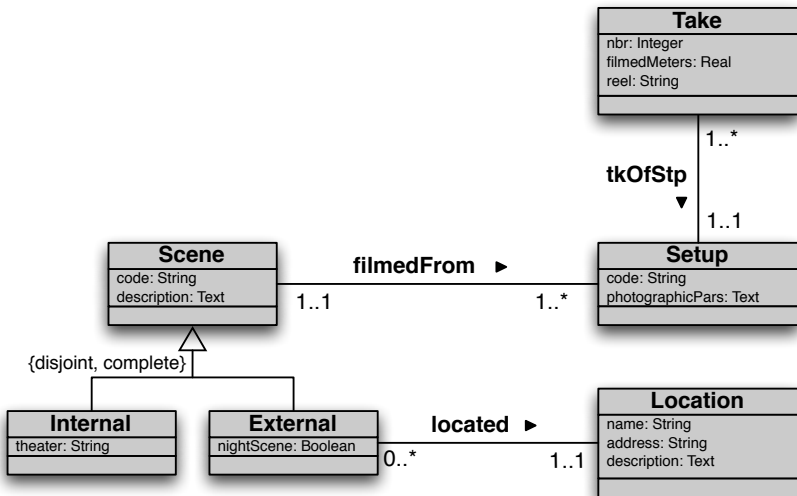
Bad points:

- Difficult to generate.
- Difficult to understand for humans.
- Too unstructured (making reasoning difficult), no well-established methodologies available.
- Automated reasoning may be impossible.

(*) *Not really a bad point, in fact.*



Solution 2: Use conceptual modeling diagrams (UML)



Solution 3: Use both!!! (cont'd)

Important:

The logical theories that are obtained from conceptual modeling diagrams are of a specific form.

- Their expressiveness is limited (or better, well-disciplined).
- One can exploit the particular form of the logical theory to simplify reasoning.
- The aim is getting:
 - decidability, and
 - reasoning procedures that match the intrinsic computational complexity of reasoning over the conceptual modeling diagrams.

Question

Which are the ontology formalisms / description logics that capture precisely such logical theories?

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity**
 - Basic definitions
 - Hardness and completeness
 - Most important complexity classes
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity**
 - **Basic definitions**
 - Hardness and completeness
 - Most important complexity classes
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References

Computational complexity (1/2)

[J.E. Hopcroft, 2007; Papadimitriou, 1994]

Computational complexity theory aims at understanding how difficult it is to solve specific problems.

- A **problem** is considered as an (in general infinite) set of instances of the problem, each encoded in some meaningful (i.e., compact) way.
- Standard complexity theory deals with **decision problems**: i.e., problems that admit a yes/no answer.
- **Algorithm** that solves a decision problem:
 - input: an instance of the problem
 - output: yes or no
- The difficulty (complexity) is measured in terms of the amount of **resources** (time, space) that the algorithm needs to solve the problem.
~> complexity of the algorithm, or **upper bound**
- To measure the complexity of the problem, we consider the best possible algorithm that solves it.
~> **lower bound**



Complexity classes

To achieve robustness wrt encoding issues, usually one does not consider specific complexity functions f , but rather families \mathcal{C} of complexity functions, giving rise to complexity classes.

Def.: A time/space **complexity class** \mathcal{C}

... is the set of all problems P such that an instance of P of size n can be solved in time/space at most $C(n)$.

Note: Consider a (decision) problem P , and an encoding of the instances of P into strings over some alphabet Σ .

Once we fix such an encoding, the problem actually corresponds to a language L_P , namely the set of strings encoding those instances of the problem for which the answer is yes.

Hence, in the technical sense, a complexity class is actually a set of languages.

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
 - Basic definitions
 - **Hardness and completeness**
 - Most important complexity classes
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References

Reductions

To establish lower bounds on the complexity of problems, we make use of the notion of reduction:

Def.: A **reduction** from a problem P_1 to a problem P_2

... is a function R (the reduction) from instance of P_1 to instances of P_2 such that:

- 1 R is efficiently computable (i.e., in logarithmic space), and
- 2 An instance I of P_1 has answer yes iff $R(I)$ has answer yes.

P_1 **reduces to** P_2 if there is a reduction R from P_1 to P_2 .

Intuition: If P_1 reduces to P_2 , then P_2 is at least as difficult as P_1 , since we can solve an instance I of P_1 by reducing it to the instance $R(I)$ of P_2 and then solve $R(I)$.

Hardness and completeness

Def.: A problem P is **hard** for a complexity class \mathcal{C}

... if every problem in \mathcal{C} can be reduced to P .

Def.: A problem P is **complete** for a complexity class \mathcal{C} if

- 1 it is hard for \mathcal{C} , and
- 2 it belongs to \mathcal{C}

Intuitively, a problem that is complete for \mathcal{C} is among the hardest problems in \mathcal{C} .

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
 - Basic definitions
 - Hardness and completeness
 - **Most important complexity classes**
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References

Tractability and intractability: P_{TIME} and NP

Def.: P_{TIME}

Set of problems solvable in **polynomial time by a deterministic TM.**

- These problems are considered **tractable**, i.e., solvable for large inputs.
- Is a robust class (P_{TIME} computations compose).

Def.: NP

Set of problems solvable in **polynomial time by a non-deterministic TM.**

- These problems are believed **intractable**, i.e., unsolvable for large inputs.
- The best known algorithms actually require exponential time.
- Corresponds to a large class of practical problems, for which the following type of algorithm can be used:
 - 1 Non-deterministically guess a possible solution of polynomial size.
 - 2 Check in polynomial time that the guessed solutions is good.

Complexity classes above NP

Def.: PSPACE

Set of problems solvable in **polynomial space** by a deterministic TM.

- Polynomial space is “not really good”, since these problems may require exponential time.
- These problems are considered to be more difficult than NP problems.
- Practical algorithms and heuristics work less well than for NP problems.

Def.: EXPTIME

Set of problems solvable in **exponential time by a deterministic TM**.

- This is the first provably intractable complexity class.
- These problems are considered to be very difficult.

Complexity classes below PTIME

Def.: LOGSPACE and NLOGSPACE

Set of problems solvable in logarithmic space by a (non-)deterministic TM.

- Note: when measuring the space complexity, the size of the input does not count, and only the working memory (TM tape) is considered.
- Note 2: logarithmic space computations compose (this is not trivial).
- Correspond to reachability in undirected and directed graphs, respectively.

Def.: AC⁰

Set of problems solvable in constant time using a polynomial number of processors.

- These problems are solvable efficiently even for very large inputs.
- Corresponds to the complexity of model checking a fixed FO formula when the input is the model only.

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases**
 - First-order logic queries
 - Conjunctive queries
 - Unions of conjunctive queries
- 4 Querying databases and ontologies
- 5 References



(Union of) Conjunctive queries – (U)CQs

(Unions of) **conjunctive queries** are an important class of queries:

- A (U)CQ is a FOL query using only conjunction, existential quantification (and disjunction).
- Hence, UCQs contain no negation, no universal quantification, and no function symbols besides constants.
- Correspond to SQL/relational algebra (**union**) **select-project-join (SPJ) queries** – the most frequently asked queries.
- (U)CQs exhibit nice computational and semantic properties, and have been studied extensively in database theory.
- They are important in practice, since relational database engines are specifically optimized for CQs.

Conjunctive queries and SQL – Example

Relational alphabet:

Person(name, age), Lives(person, city), Manages(boss, employee)

Query: return name and age of all persons that live in the same city as their boss.

Expressed in SQL:

```
SELECT P.name, P.age
FROM Person P, Manages M, Lives L1, Lives L2
WHERE P.name = L1.person AND P.name = M.employee AND
      M.boss = L2.person AND L1.city = L2.city
```

Expressed as a CQ: (the distinguished variables are the blue ones)

$$\exists b, e, p_1, c_1, p_2, c_2. \text{Person}(n, a) \wedge \text{Manages}(b, e) \wedge \text{Lives}(p_1, c_1) \wedge \text{Lives}(p_2, c_2) \wedge n = p_1 \wedge n = e \wedge b = p_2 \wedge c_1 = c_2$$

Or simpler: $\exists b, c. \text{Person}(n, a) \wedge \text{Manages}(b, n) \wedge \text{Lives}(n, c) \wedge \text{Lives}(b, c)$



Datalog notation for CQs

A CQ $q = \exists \vec{y}. conj(\vec{x}, \vec{y})$ can also be written using **datalog notation** as

$$q(\vec{x}_1) \leftarrow conj'(\vec{x}_1, \vec{y}_1)$$

where $conj'(\vec{x}_1, \vec{y}_1)$ is the list of atoms in $conj(\vec{x}, \vec{y})$ obtained by equating the variables \vec{x}, \vec{y} according to the equalities in $conj(\vec{x}, \vec{y})$.

As a result of such an equality elimination, we have that \vec{x}_1 and \vec{y}_1 can contain constants and multiple occurrences of the same variable.

Def.: In the above query q , we call:

- $q(\vec{x}_1)$ the **head**;
- $conj'(\vec{x}_1, \vec{y}_1)$ the **body**;
- the variables in \vec{x}_1 the **distinguished variables**;
- the variables in \vec{y}_1 the **non-distinguished variables**.

Conjunctive queries – Example

- Consider the alphabet $\Sigma = \{E/2\}$ and an **interpretation** $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$. Note that $E^{\mathcal{I}}$ is a binary relation, i.e., \mathcal{I} is a directed graph.
- The following **CQ** q returns all nodes that participate to a triangle in the graph:

$$\exists y, z. E(x, y) \wedge E(y, z) \wedge E(z, x)$$

- The query q in **datalog notation** becomes:

$$q(x) \leftarrow E(x, y), E(y, z), E(z, x)$$

- The query q in **SQL** is (we use `Edge(f, s)` for $E(x, y)$):

```
SELECT E1.f
FROM Edge E1, Edge E2, Edge E3
WHERE E1.s = E2.f AND E2.s = E3.f AND E3.s = E1.f
```


CQ evaluation – Combined, data, and query complexity

Theorem (**Combined complexity** of CQ evaluation)

$\{\langle \mathcal{I}, \alpha, q \rangle \mid \mathcal{I}, \alpha \models q\}$ is **NP-complete** — see below for hardness.

- time: exponential
- space: polynomial

Theorem (**Data complexity** of CQ evaluation)

$\{\langle \mathcal{I}, \alpha \rangle \mid \mathcal{I}, \alpha \models q\}$ is **in AC^0**

- time: polynomial
- space: logarithmic

Theorem (**Query complexity** of CQ evaluation)

$\{\langle \alpha, q \rangle \mid \mathcal{I}, \alpha \models q\}$ is **NP-complete** — see below for hardness.

- time: exponential
- space: polynomial

NP-hardness of CQ evaluation

The previous reduction immediately gives us the hardness for combined complexity.

Theorem

CQ evaluation is NP-hard in combined complexity.

Note: in the previous reduction, the interpretation does not depend on the actual graph. Hence, the reduction provides also the lower-bound for query complexity.

Theorem

CQ evaluation is NP-hard in query (and combined) complexity.



Datalog notation for UCQs

A union of conjunctive queries

$$q = \bigvee_{i=1, \dots, n} \exists \vec{y}_i. \text{conj}_i(\vec{x}, \vec{y}_i)$$

is written in **datalog notation** as

$$\left\{ \begin{array}{l} q(\vec{x}) \leftarrow \text{conj}'_1(\vec{x}, \vec{y}'_1) \\ \vdots \\ q(\vec{x}) \leftarrow \text{conj}'_n(\vec{x}, \vec{y}'_n) \end{array} \right\}$$

where each element of the set is the datalog expression corresponding to the conjunctive query $q_i = \exists \vec{y}_i. \text{conj}_i(\vec{x}, \vec{y}_i)$.

Note: in general, we omit the set brackets.

Evaluation of UCQs

From the definition of FOL query we have that:

$$\mathcal{I}, \alpha \models \bigvee_{i=1, \dots, n} \exists \vec{y}_i. conj_i(\vec{x}, \vec{y}_i)$$

if and only if

$$\mathcal{I}, \alpha \models \exists \vec{y}_i. conj_i(\vec{x}, \vec{y}_i) \quad \text{for some } i \in \{1, \dots, n\}.$$

Hence to evaluate a UCQ q , we simply evaluate a number (linear in the size of q) of conjunctive queries in isolation.

Hence, **evaluating UCQs has the same complexity as evaluating CQs.**

UCQ evaluation – Combined, data, and query complexity

Theorem (**Combined complexity** of UCQ evaluation)

$\{\langle \mathcal{I}, \alpha, q \rangle \mid \mathcal{I}, \alpha \models q\}$ is **NP-complete**.

- time: exponential
- space: polynomial

Theorem (**Data complexity** of UCQ evaluation)

$\{\langle \mathcal{I}, q \rangle \mid \mathcal{I}, \alpha \models q\}$ is **in AC^0** (query q fixed).

- time: polynomial
- space: logarithmic

Theorem (**Query complexity** of UCQ evaluation)

$\{\langle \alpha, q \rangle \mid \mathcal{I}, \alpha \models q\}$ is **NP-complete** (interpretation \mathcal{I} fixed).

- time: exponential
- space: polynomial



Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies**
 - Query answering in traditional databases
 - Query answering in ontologies
 - Query answering in ontology-based data access
- 5 References



Query answering

In ontology-based data access we are interested in a reasoning service that is not typical in ontologies (or in a FOL theory, or in UML class diagrams, or in a knowledge base) but it is very common in databases: **query answering**.

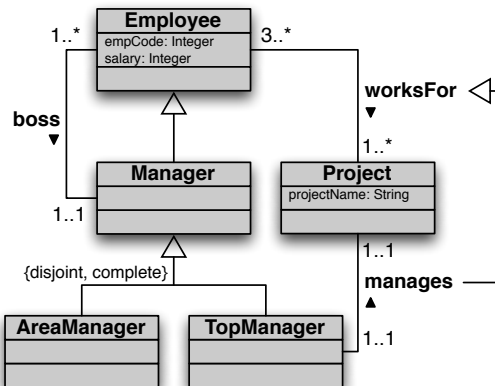
Def.: Query

Is an expression at the intensional level denoting a set of tuples of individuals satisfying a given condition.

Def.: Query Answering

Is the reasoning service that actually computes the answer to a query.

Example of query



$$q(\mathit{ce}, \mathit{cm}, \mathit{sa}) \leftarrow \exists e, p, m. \text{worksFor}(e, p) \wedge \text{manages}(m, p) \wedge \text{boss}(m, e) \wedge \text{empCode}(e, \mathit{ce}) \wedge \text{empCode}(m, \mathit{cm}) \wedge \text{salary}(e, \mathit{sa}) \wedge \text{salary}(m, \mathit{sa})$$

Query answering under different assumptions

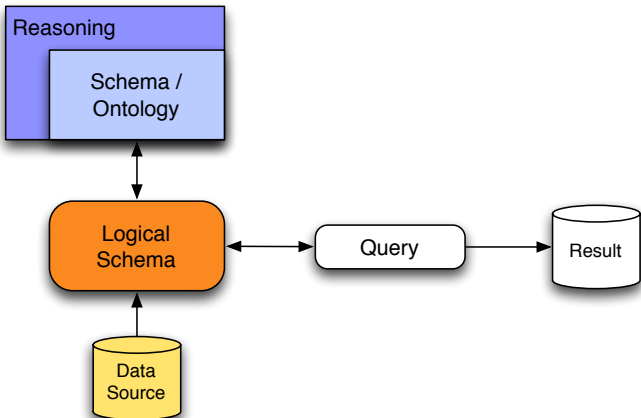
There are two fundamentally different assumptions when addressing query answering:

- **Complete information** on the data, as in traditional databases.
- **Incomplete information** on the data, as in ontologies, but also information integration in databases.

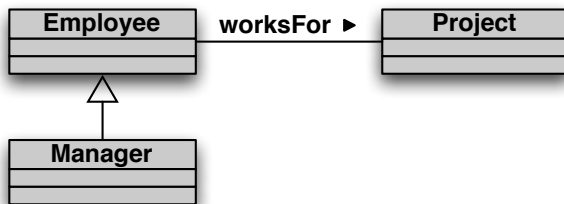
Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies**
 - Query answering in traditional databases
 - Query answering in ontologies
 - Query answering in ontology-based data access
- 5 References

Query answering in traditional databases (cont'd)



Query answering in traditional databases – Example



For each concept/relationship we have a (complete) table in the DB.

DB: Employee = { john, mary, nick }
Manager = { john, nick }
Project = { prA, prB }
worksFor = { (john,prA), (mary,prB) }

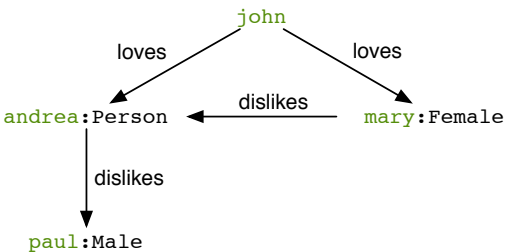
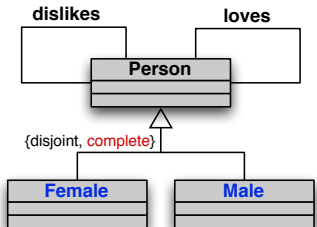
Query: $q(x) \leftarrow \exists p. \text{Manager}(x) \wedge \text{Project}(p) \wedge \text{worksFor}(x, p)$

Answer: { john }

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies**
 - Query answering in traditional databases
 - Query answering in ontologies**
 - Query answering in ontology-based data access
- 5 References

QA in ontologies – Andrea's Example (cont'd)



$$q(x) \leftarrow \exists y, z. \text{loves}(x, y) \wedge \text{Female}(y) \wedge \text{dislikes}(y, z) \wedge \text{Male}(z)$$

Answer: { **john** }

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

Questions that need to be addressed

In the context of ontology-based data access:

- 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** take into account these issues?

Which language to use for querying ontologies?

Two borderline 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 tradeoff is to use (unions of) **conjunctive queries**.

Outline of Lecture 1

- 1 Ontology-based data management
- 2 Brief introduction to computational complexity
- 3 Query answering in databases
- 4 Querying databases and ontologies
- 5 References

References II

[Vardi, 1982] Moshe Y. Vardi.

The complexity of relational query languages.

In *Proc. of the 14th ACM SIGACT Symp. on Theory of Computing (STOC'82)*, pages 137–146, 1982.