

An Introduction to Description Logics and Query Rewriting

Roman Kontchakov and Michael Zakharyashev

Department of Computer Science and Information Systems,
Birkbeck, University of London, U.K.

Abstract. This chapter gives an overview of the description logics underlying the OWL 2 Web Ontology Language and its three tractable profiles, OWL 2 RL, OWL 2 EL and OWL 2 QL. We consider the syntax and semantics of these description logics as well as main reasoning tasks and their computational complexity. We also discuss the semantical foundations for first-order and datalog rewritings of conjunctive queries over knowledge bases given in the OWL 2 profiles, and outline the architecture of the ontology-based data access system Ontop.

1 Introduction

The first aim of this chapter is to introduce and discuss logic-based formalisms that underpin the OWL 2 Web Ontology Language and its three tractable profiles: RL, EL and QL. OWL 2 is a rather involved language that was designed to represent knowledge about various domains of interest in a machine-accessible form. The diagram in Fig. 1, taken from the official W3C document,¹ shows the general structure of OWL 2. As follows from the diagram, there are (at least) five syntaxes for OWL 2 (various tools can use their own versions). In this chapter, we consider a sixth one, the language of Description Logic (DL). Although not covering all the bells and whistles of the full OWL 2, it will allow the reader to quickly grasp the meaning of the main modelling constructs that OWL 2 provides. The language of DL is elegant and concise because it stems from the formalisms that have been developed in mathematical logic since the middle of the 19th century. It is underpinned by the crystal-clear model-theoretic semantics developed by A. Tarski since the mid 1930s (as shown in the diagram in Fig. 1, OWL 2 has two semantics: RDF-based² and Direct Semantics³; the latter is based on the model-theoretic semantics of DLs).

We will introduce, in Section 2, most important modelling constructs of OWL 2 and their semantics in terms of the description logics *ALCHI* and *SROIQ* and the model-theoretic semantics. We then explain fundamental reasoning tasks such as checking consistency (satisfiability), concept and role subsumption, instance checking and conjunctive query answering, and discuss their computational complexity. Our focus in Section 3 is on the DLs underlying the three

¹ www.w3.org/TR/owl2-overview

² www.w3.org/TR/owl2-rdf-based-semantics

³ www.w3.org/TR/owl2-direct-semantics

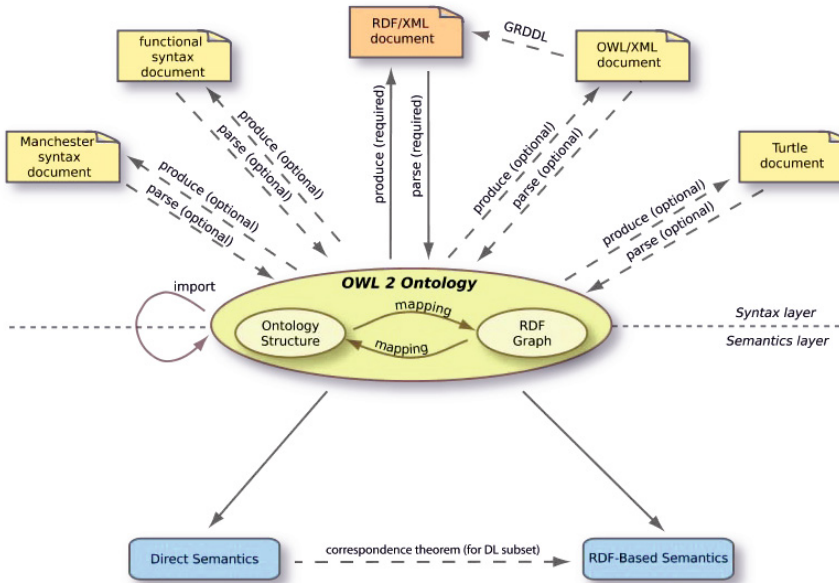


Fig. 1. General structure of OWL 2

profiles (or fragments) of OWL2 that were identified by the OWL2 working group to ensure tractability of reasoning at the expense of the expressive power. The Euler diagram in Fig. 2 gives a general overview of the DLs considered in this chapter (RDFS denotes the RDFS fragment of OWL2 DL under the Direct Semantics).

The second aim of the chapter is to explain, in Section 4, the semantical foundations for first-order and datalog rewritings of conjunctive queries over knowledge bases given in the OWL 2 profiles. Query rewriting is a fundamental technique underlying ontology-based data access: it reduces answering queries over knowledge bases to answering first-order or datalog queries over plain data, which can be done using conventional database management systems or, respectively, datalog engines. For more expressive languages, ontology-based data access will be discussed in Chapter 6.

Finally, in Section 5, we give an overview of the architecture of the ontology-based data access system *Ontop*⁴.

2 Description Logics

Description Logic is an area of knowledge representation and reasoning in Artificial Intelligence and the Semantic Web that studies logic-based formalisms whose languages operate with *concepts* to represent classes of individuals in an

⁴ ontop.inf.unibz.it

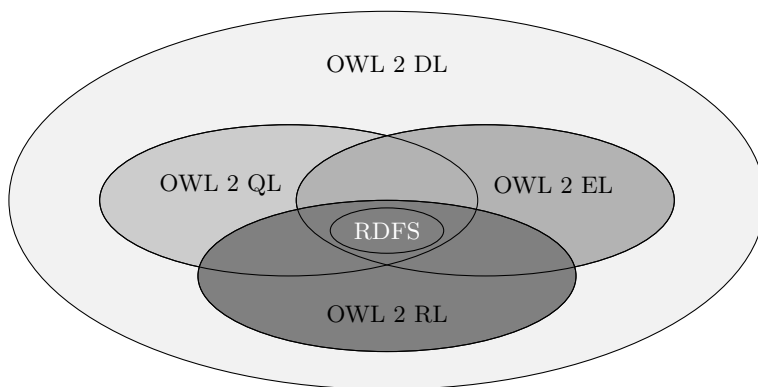


Fig. 2. Relationships between the DL fragments of OWL 2

application domain, and *roles* to represent binary relations between the individuals. Each concrete formalism, called a *description logic* (DL, for short), is characterised by its set of constructs that can be used to build complex concepts and roles from primitive ones.

The zoo of DLs is very big. We begin by defining one particular representative, which goes in the zoo under the moniker of *ALC^HI*. The example below shows a simple knowledge base (or an ontology), which is given in the OWL functional-style syntax (FSS) and the more terse syntax of *ALC^HI*. The reader is invited to decipher the meaning of the knowledge base before consulting the formal definitions.

Example 1. The following three statements are written in the FSS:

- SubClassOf(ObjectIntersectionOf(Person, (1)
ObjectSomeValuesFrom(takesCourse, Course)), Student),
- SubObjectPropertyOf(mastersDegreeFrom, degreeFrom), (2)
- SubClassOf(ObjectSomeValuesFrom(
ObjectInverseOf(takesCourse), owl:Thing), Course), (3)
- ClassAssertion(Student, john), (4)
- ObjectPropertyAssertion(takesCourse, john, sw). (5)

Using the syntax of *ALC^HI*, the same statements can be expressed as follows:

- Person $\sqcap \exists \text{takesCourse}.\text{Course} \sqsubseteq \text{Student}$, (1')
- mastersDegreeFrom $\sqsubseteq \text{degreeFrom}$, (2')
- takesCourse⁻. $\top \sqsubseteq \text{Course}$, (3')
- Student(john), (4')
- takesCourse(john, sw). (5')

2.1 Syntax

The alphabet of \mathcal{ALCHI} contains three pairwise disjoint and countably infinite sets: *concept names* A_1, A_2, \dots , *role names* P_1, P_2, \dots and *individual names* a_1, a_2, \dots . An \mathcal{ALCHI} role, R , is either a role name P_i or its *inverse* P_i^- . \mathcal{ALCHI} concepts, C , are constructed from two special primitive concepts, \top ('top') and \perp ('bottom'), concept names and roles using the following grammar:

$$C ::= A_i \mid \top \mid \perp \mid \neg C \mid C_1 \sqcap C_2 \mid C_1 \sqcup C_2 \mid \exists R.C \mid \forall R.C.$$

An \mathcal{ALCHI} *terminological box* (or TBox), \mathcal{T} , is a finite set of *concept* and *role inclusion axioms* of the form

$$C_1 \sqsubseteq C_2 \quad \text{and} \quad R_1 \sqsubseteq R_2,$$

where C_1, C_2 are concepts and R_1, R_2 roles. An \mathcal{ALCHI} *assertion box* (or ABox), \mathcal{A} , is a finite set of *concept* and *role assertions* of the form

$$C(a) \quad \text{and} \quad R(a, b),$$

where C is a concept, R a role and a, b are individual names. Given an ABox \mathcal{A} , we denote by $\text{ind}(\mathcal{A})$ the set of individual names that occur in \mathcal{A} . Taken together, \mathcal{T} and \mathcal{A} comprise an \mathcal{ALCHI} *knowledge base* (or KB) $\mathcal{K} = (\mathcal{T}, \mathcal{A})$.

Example 1 shows that, in many respects, the DL syntax is simply a less verbose form of the FSS. Thus, \sqsubseteq stands for `SubClassOf` and `SubObjectPropertyOf`, and P^- for `ObjectInverseOf(P)`; the correspondences for concept constructs are listed below:

DL	FSS
\perp	<code>owl:Nothing</code>
\top	<code>owl:Thing</code>
$\neg C$	<code>ObjectComplementOf(C)</code>
$C_1 \sqcap C_2$	<code>ObjectIntersectionOf(C₁, C₂)</code>
$C_1 \sqcup C_2$	<code>ObjectUnionOf(C₁, C₂)</code>
$\exists R.C$	<code>ObjectSomeValuesFrom(R, C)</code>
$\forall R.C$	<code>ObjectAllValuesFrom(R, C)</code>

In Section 2.4, we shall see how most of the OWL constructs can be mapped into the DL syntax. But before that, we need to define the meaning of DL constructs.

2.2 Semantics

As well as all DLs, \mathcal{ALCHI} is equipped with a Tarski-style semantics defined in terms of interpretations (which are a simplified form of interpretations in the Direct Semantics of OWL). An *interpretation* \mathcal{I} is a pair $(\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ that consists of a non-empty *domain of interpretation* $\Delta^{\mathcal{I}}$ and an *interpretation function* $\cdot^{\mathcal{I}}$. The latter assigns

- an element $a_i^{\mathcal{I}} \in \Delta^{\mathcal{I}}$ to each individual name a_i ;
- a subset $A_i^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$ to each concept name A_i ;
- a binary relation $P_i^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$ to each role name P_i .

(In \mathcal{ALCHL} , distinct individuals are usually assumed to be interpreted by distinct domain elements—this is called the *unique name assumption*, or UNA. In this chapter we follow the convention and assume the UNA for all of our DLs, which cannot distinguish between models with and without UNA. Note, however, that OWL 2 and its three profiles *do not* adopt the UNA and use constructs like `SameIndividual` and `DifferentIndividuals` instead.) We extend inductively the interpretation function $\cdot^{\mathcal{I}}$ to complex roles and concepts by taking

$$\begin{aligned}
 (P^-)^{\mathcal{I}} &= \{(v, u) \mid (u, v) \in P^{\mathcal{I}}\}, \\
 \top^{\mathcal{I}} &= \Delta^{\mathcal{I}}, \\
 \perp^{\mathcal{I}} &= \emptyset, \\
 (\neg C)^{\mathcal{I}} &= \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}}, \\
 (C_1 \sqcap C_2)^{\mathcal{I}} &= C_1^{\mathcal{I}} \cap C_2^{\mathcal{I}}, \\
 (C_1 \sqcup C_2)^{\mathcal{I}} &= C_1^{\mathcal{I}} \cup C_2^{\mathcal{I}}, \\
 (\exists R.C)^{\mathcal{I}} &= \{u \mid \text{there is } v \in C^{\mathcal{I}} \text{ such that } (u, v) \in R^{\mathcal{I}}\}, \\
 (\forall R.C)^{\mathcal{I}} &= \{u \mid v \in C^{\mathcal{I}}, \text{ for all } v \text{ with } (u, v) \in R^{\mathcal{I}}\}.
 \end{aligned}$$

Having fixed the interpretation of individual names, concepts and roles, we now define the *satisfaction relation* \models for inclusions and assertions:

$$\begin{aligned}
 \mathcal{I} \models C_1 \sqsubseteq C_2 &\quad \text{if and only if} \quad C_1^{\mathcal{I}} \subseteq C_2^{\mathcal{I}}, \\
 \mathcal{I} \models R_1 \sqsubseteq R_2 &\quad \text{if and only if} \quad R_1^{\mathcal{I}} \subseteq R_2^{\mathcal{I}}, \\
 \mathcal{I} \models C(a) &\quad \text{if and only if} \quad a^{\mathcal{I}} \in C^{\mathcal{I}}, \\
 \mathcal{I} \models R(a, b) &\quad \text{if and only if} \quad (a^{\mathcal{I}}, b^{\mathcal{I}}) \in R^{\mathcal{I}}.
 \end{aligned}$$

We say that an interpretation \mathcal{I} is a *model* of a knowledge base $\mathcal{K} = (\mathcal{T}, \mathcal{A})$ if it satisfies all concept and roles inclusions of \mathcal{T} and all concept and role assertions of \mathcal{A} . In this case we write $\mathcal{I} \models \mathcal{K}$ (and also $\mathcal{I} \models \mathcal{T}$ and $\mathcal{I} \models \mathcal{A}$).

It is to be remembered that (unlike databases) the choice of domains and interpretation functions is not fixed in the DL semantics, so that every knowledge base can have many models. This reflects the *open world assumption*, or OWA, adopted in DL (and OWL), according to which no single agent can possess complete knowledge. Thus, we have to consider all possible assignments of truth-values to assertions—as long as they do not contradict the given knowledge base. (Databases adopt the *closed world assumption*, or CWA, that defines everything unknown as false.)

Example 2. Consider the following knowledge base $\mathcal{K} = (\mathcal{T}, \mathcal{A})$:

$$\begin{aligned}\mathcal{T} &= \{ \text{GraduateStudent} \sqsubseteq \text{Student}, \\ &\quad \text{GraduateStudent} \sqsubseteq \exists \text{takesCourse}.\text{GraduateCourse} \}, \\ \mathcal{A} &= \{ \text{GraduateStudent}(\text{john}) \}.\end{aligned}$$

Denote by \mathcal{I}_1 an interpretation with domain $\Delta^{\mathcal{I}_1} = \{\text{john}, \text{sw}\}$ such that

$$\begin{aligned}\text{john}^{\mathcal{I}_1} &= \text{john}, \\ \text{GraduateStudent}^{\mathcal{I}_1} &= \{\text{john}\}, & \text{Student}^{\mathcal{I}_1} &= \{\text{john}\}, \\ \text{GraduateCourse}^{\mathcal{I}_1} &= \{\text{sw}\}, & \text{takesCourse}^{\mathcal{I}_1} &= \{(\text{john}, \text{sw})\}.\end{aligned}$$

The reader can readily check that \mathcal{I}_1 is a model of \mathcal{K} ; that is, the ‘world’ described by \mathcal{I}_1 satisfies the knowledge and data given in \mathcal{K} . Now, take another interpretation \mathcal{I}_2 with domain $\Delta^{\mathcal{I}_2} = \{a\}$ in which

$$\begin{aligned}\text{john}^{\mathcal{I}_2} &= a, \\ \text{GraduateStudent}^{\mathcal{I}_2} &= \{a\}, & \text{Student}^{\mathcal{I}_2} &= \{a\}, \\ \text{GraduateCourse}^{\mathcal{I}_2} &= \{a\}, & \text{takesCourse}^{\mathcal{I}_2} &= \{(a, a)\}.\end{aligned}$$

This interpretation does not make much sense from the modelling point of view (because a takes course a in the world described by \mathcal{I}_2). Nevertheless, \mathcal{I}_2 satisfies all of the inclusions in \mathcal{T} and assertions in \mathcal{A} , and so is a model of \mathcal{K} . Yet another interpretation, \mathcal{I}_3 , with domain $\Delta^{\mathcal{I}_3} = \{\text{john}\}$ and

$$\begin{aligned}\text{john}^{\mathcal{I}_3} &= \text{john}, \\ \text{GraduateStudent}^{\mathcal{I}_3} &= \{\text{john}\}, & \text{Student}^{\mathcal{I}_3} &= \{\text{john}\}, \\ \text{GraduateCourse}^{\mathcal{I}_3} &= \emptyset, & \text{takesCourse}^{\mathcal{I}_3} &= \emptyset\end{aligned}$$

satisfies the assertions in \mathcal{A} and the first concept inclusion in \mathcal{T} but fails to satisfy the second concept inclusion, and so is not a model of \mathcal{K} .

Intuitively, everything that takes place in *each and every* model of a knowledge base is a logical consequence of the KB, not necessarily explicitly presented in it. Finding logical consequences is usually referred to as *reasoning*.

2.3 Reasoning Problems

A very basic reasoning problem is to decide whether a given knowledge base is *consistent* in the sense that it does not imply mutually contradicting statements. Formally, we can call a knowledge base \mathcal{K} *satisfiable* (or *consistent*) if there exists at least one model of \mathcal{K} (which is obviously enough to guarantee that \mathcal{K} contains no contradictions).

Example 3. Let \mathcal{T} be a TBox containing the following concept inclusions:

$$\begin{aligned} \text{UndergraduateStudent} &\sqsubseteq \forall \text{takesCourse}.\text{UndergraduateCourse}, \\ \text{UndergraduateCourse} \sqcap \text{GraduateCourse} &\sqsubseteq \perp, \end{aligned}$$

and \mathcal{A} an ABox with the following assertions:

$$\begin{aligned} &\text{UndergraduateStudent}(\text{john}), \\ &\text{takesCourse}(\text{john}, \text{sw}), \\ &\text{GraduateCourse}(\text{sw}). \end{aligned}$$

If we assume that $(\mathcal{T}, \mathcal{A})$ has a model, then the undergraduate student John in it will only be able to take undergraduate courses that are not graduate courses such as SW. However, according to the ABox, John takes SW, which is a contradiction. Thus, $(\mathcal{T}, \mathcal{A})$ is inconsistent.

Another important reasoning problem is entailment. We say that a concept inclusion $C_1 \sqsubseteq C_2$ is *entailed by* a knowledge base \mathcal{K} and write $\mathcal{K} \models C_1 \sqsubseteq C_2$ if $\mathcal{I} \models C_1 \sqsubseteq C_2$ for all models \mathcal{I} of \mathcal{K} (entailment for role inclusions and concept and role assertions is defined similarly).

Example 4. Consider a TBox \mathcal{T} with the following two concept inclusions:

$$\begin{aligned} \forall \text{takesCourse}.\text{UndergraduateCourse} &\sqsubseteq \text{UndergraduateStudent}, \\ \text{FirstYearStudent} &\sqsubseteq \exists \text{takesCourse}.\text{UndergraduateCourse}. \end{aligned}$$

In the model \mathcal{I}_1 of \mathcal{T} given below, $\text{FirstYearStudent} \sqsubseteq \text{UndergraduateStudent}$ holds true. However, \mathcal{T} *does not entail* this concept inclusion because there is another model, \mathcal{I}_2 , where it is not satisfied:

$$\begin{array}{ll} \Delta^{\mathcal{I}_1} = \{j, s\}, & \Delta^{\mathcal{I}_2} = \{j, s, \ell\}, \\ \text{takesCourse}^{\mathcal{I}_1} = \{(j, s)\}, & \text{takesCourse}^{\mathcal{I}_2} = \{(j, s), (j, \ell)\}, \\ \text{FirstYearStudent}^{\mathcal{I}_1} = \{j\}, & \text{FirstYearStudent}^{\mathcal{I}_2} = \{j\}, \\ \text{UndergraduateStudent}^{\mathcal{I}_1} = \{j\}, & \text{UndergraduateStudent}^{\mathcal{I}_2} = \emptyset, \\ \text{UndergraduateCourse}^{\mathcal{I}_1} = \{s\}, & \text{UndergraduateCourse}^{\mathcal{I}_2} = \{s\}. \end{array}$$

Intuitively, in \mathcal{I}_1 , the individual j is a first-year student who takes only one undergraduate course, and so must be an undergraduate student (to satisfy the first concept inclusion). In contrast, in \mathcal{I}_2 , the individual j also takes another course, ℓ , which is not an undergraduate course, and therefore, j does not have to be an undergraduate student (still satisfying the first concept inclusion); see Fig. 3.

If $C_1 \sqsubseteq C_2$ is entailed by \mathcal{K} , then we also say that C_2 *subsumes* C_1 with respect to \mathcal{K} (or that C_1 is *subsumed* by C_2 with respect to \mathcal{K}). We note in passing that if \mathcal{K} is inconsistent then any concept is subsumed by any other

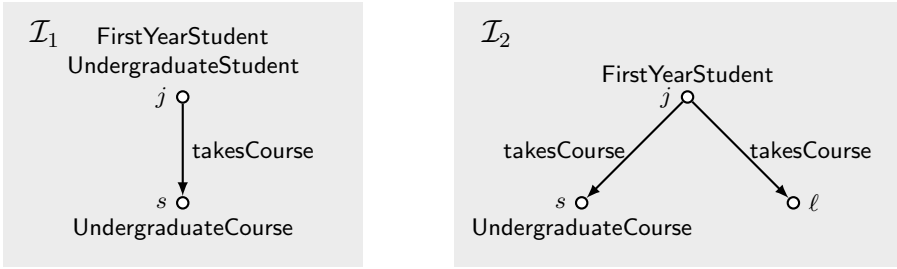


Fig. 3. Interpretations for \mathcal{T} in Example 4

concept: indeed, an arbitrary inclusion $C_1 \sqsubseteq C_2$ is trivially true in every model of \mathcal{K} simply because \mathcal{K} has no models. The following proposition, the proof of which is left to the reader as a simple exercise, shows that concept subsumption is in fact reducible to consistency.

Proposition 5. $(\mathcal{T}, \mathcal{A}) \models C_1 \sqsubseteq C_2$ if and only if $(\mathcal{T}, \mathcal{A} \cup \{C_1(a), \neg C_2(a)\})$ is not satisfiable, for a fresh individual name a (not occurring in \mathcal{A}).

(Note as a warning that not every DL allows negation and complex concepts in the ABoxes, in which case the reduction may be not so trivial.)

If $C(a)$ is entailed by \mathcal{K} (that is, $C(a)$ holds in every model of \mathcal{K}), then we also say that a is an *instance* of C with respect to \mathcal{K} . The problem of checking whether a is an instance of a given C with respect to \mathcal{K} is called *instance checking*. This problem is also reducible to knowledge base consistency (provided that complex concepts are allowed in the ABoxes):

Proposition 6. $(\mathcal{T}, \mathcal{A}) \models C(a)$ if and only if $(\mathcal{T}, \mathcal{A} \cup \{\neg C(a)\})$ is not satisfiable.

A more general reasoning task is answering conjunctive queries over knowledge bases. A *conjunctive query* (CQ for short) $\mathbf{q}(\mathbf{x})$ is an expression of the form $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$, where $\varphi(\mathbf{x}, \mathbf{y})$ is a conjunction of atoms such as $A(z)$ and $P(z_1, z_2)$, for a concept name A , role name P and terms z, z_1 and z_2 , which are individual names or variables from \mathbf{x} and \mathbf{y} . The variables x_i in $\mathbf{x} = (x_1, \dots, x_n)$ are called *answer variables* and the variables in \mathbf{y} *existentially quantified variables*. Given a tuple $\mathbf{a} = (a_1, \dots, a_n)$ of individual names from \mathcal{A} , we denote by $\mathbf{q}(\mathbf{a})$ the result of replacing each answer variable x_i in $\exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$ with the respective a_i from \mathbf{a} . A tuple \mathbf{a} of individual names from \mathcal{A} is a *certain answer* to $\mathbf{q}(\mathbf{x})$ over $(\mathcal{T}, \mathcal{A})$ if, for any model \mathcal{I} of $(\mathcal{T}, \mathcal{A})$, the sentence $\mathbf{q}(\mathbf{a})$ is true in \mathcal{I} ($\mathcal{I} \models \mathbf{q}(\mathbf{a})$, in symbols). We write $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ to indicate that \mathbf{a} is a certain answer to $\mathbf{q}(\mathbf{x})$ over $(\mathcal{T}, \mathcal{A})$. A CQ \mathbf{q} without answer variables is called *Boolean*, in which case a certain answer to \mathbf{q} over $(\mathcal{T}, \mathcal{A})$ is ‘yes’ if $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}$ and ‘no’ otherwise. The problem of answering Boolean CQs is known as *CQ entailment*.

Example 7. (Andrea’s example (Schaerf, 1993)) Suppose a TBox \mathcal{T} contains the inclusions

$$\top \sqsubseteq \text{Male} \sqcup \text{Female}, \quad \text{Male} \sqcap \text{Female} \sqsubseteq \perp$$

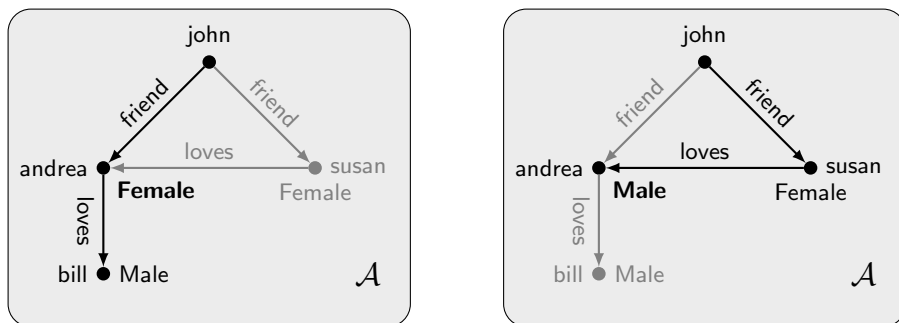


Fig. 4. Two representative models of $(\mathcal{T}, \mathcal{A})$ in Example 7

and an ABox \mathcal{A} contains the assertions

- | | |
|-----------------------|-----------------------|
| friend(john, susan), | friend(john, andrea), |
| loves(susan, andrea), | loves(androia, bill), |
| Female(susan), | Male(bill). |

Consider the CQ

$$q(x) = \exists y, z (\text{friend}(x, y) \wedge \text{Female}(y) \wedge \text{loves}(y, z) \wedge \text{Male}(z)),$$

which asks to find every individual (in the ABox) with a female friend who is in love with a male. Note that the same CQ can be expressed in the query language SPARQL⁵ as follows:

```

SELECT ?x
WHERE {
    ?x :friend ?y.
    ?y a :Female.
    ?y :loves ?z.
    ?z a :Male.
}
    
```

(Here **a** is an abbreviation of `rdf:type` and can be read as ‘is a.’) We invite the reader to check that the only certain answer to $q(x)$ over $(\mathcal{T}, \mathcal{A})$ is $x \mapsto \text{john}$. (Hint: in every model \mathcal{I} of $(\mathcal{T}, \mathcal{A})$, either $\text{androia}^{\mathcal{I}} \in \text{Female}^{\mathcal{I}}$ or $\text{androia}^{\mathcal{I}} \in \text{Male}^{\mathcal{I}}$; see Fig. 4.)

Observe that this particular CQ $q(x)$ can also be represented as an *instance query* $C(x)$, where C is the complex concept

$$\exists \text{friend} . (\text{Female} \sqcap \exists \text{loves} . \text{Male}).$$

⁵ www.w3.org/TR/sparql11-query

It is readily seen that, for any individual name a , we have $(\mathcal{T}, \mathcal{A}) \models C(a)$ if and only if $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(a)$. However, in general, CQs are not necessarily tree-shaped, can contain more than one answer variable, and so are more expressive than instance queries.

2.4 From OWL to DL

The language of OWL 2 contains more constructs than any of the DLs, with many of these constructs being just shortcuts for certain DL expressions. For example,

ObjectPropertyDomain(takesCourse, Student)

can be represented as the concept inclusion

$$\exists \text{takesCourse}.\top \sqsubseteq \text{Student}.$$

Similarly,

ObjectPropertyRange(takesCourse, Course)

can be represented as

$$\exists \text{takesCourse}^{\neg}.\top \sqsubseteq \text{Course}.$$

Note that the latter concept inclusion can only be written in the DLs with *inverse roles* (such as the RL or QL profiles of OWL 2 to be discussed below). If, however, inverse roles are not available in a DL (for instance, \mathcal{ALC}), then one can use a universal restriction:

$$\top \sqsubseteq \forall \text{takesCourse}.\text{Course}.$$

One can verify that the two concept inclusions are equivalent in the sense that they are satisfied by precisely the same interpretations. We have collected standard equivalencies of this sort in the following proposition:

Proposition 8. *The following pairs of (sets of) concept inclusions have the same models:*

$$C_1 \sqsubseteq \forall P.C_2 \quad \text{and} \quad \exists P^{\neg}.C_1 \sqsubseteq C_2, \quad (6)$$

$$C_1 \sqcup C_2 \sqsubseteq C \quad \text{and} \quad \{C_i \sqsubseteq C \mid i = 1, 2\}, \quad (7)$$

$$C \sqsubseteq C_1 \sqcap C_2 \quad \text{and} \quad \{C \sqsubseteq C_i \mid i = 1, 2\}, \quad (8)$$

$$C_1 \sqsubseteq \neg C_2 \quad \text{and} \quad C_1 \sqcap C_2 \sqsubseteq \perp. \quad (9)$$

OWL 2 also has a shortcut for equivalent classes and properties: **EquivalentClasses**, **EquivalentObjectProperties** and **EquivalentDataProperties**. These can be represented using \sqsubseteq : for example,

EquivalentClasses(C_1, C_2, \dots, C_n)

can be written in the DL parlance as n concept inclusions

$$C_1 \sqsubseteq C_2, \quad C_2 \sqsubseteq C_3, \quad \dots, \quad C_{n-1} \sqsubseteq C_n, \quad C_n \sqsubseteq C_1.$$

Another way is to use a common DL abbreviation \equiv , which is defined by taking $A \equiv B$ if and only if $A \sqsubseteq B$ and $B \sqsubseteq A$.

Equivalence (9) provides two alternative ways of expressing *disjointness* of concepts C_1 and C_2 in DL. OWL2 offers a shortcut for pairwise disjointness of n classes:

$$\text{DisjointClasses}(C_1, C_2, \dots, C_n),$$

which can be represented in DLs as $n(n-1)/2$ concept inclusions

$$C_i \sqcap C_j \sqsubseteq \perp, \quad \text{for all } i, j \text{ with } 1 \leq i < j \leq n.$$

Yet another shortcut in OWL2 allows one to say that class C is the disjoint union of C_1, \dots, C_n :

$$\text{DisjointUnion}(C, C_1, C_2, \dots, C_n),$$

which combines pairwise disjointness of C_1, \dots, C_n (see above) with

$$C \equiv C_1 \sqcup \dots \sqcup C_n.$$

Object properties can be declared symmetric in OWL2 by using axioms of the form $\text{SymmetricObjectProperty}(P)$. The same effect can be achieved in DLs with the help of the role inclusion

$$P^- \sqsubseteq P.$$

More expressive description logics contain additional constructs such as number restrictions, $\exists R.\text{Self}$, transitive roles, role chains, etc. In particular, the DL subset of OWL2, known as OWL2DL, is based on the description logic called *SROIQ* (Horrocks et al., 2006).

2.5 Complexity of Reasoning

Having formulated the reasoning problems, we are facing the following fundamental questions:

- Are these problems *decidable* in the sense that there exist algorithms which always halt and return correct answers?
- How *complex* are such algorithms in terms of the time or memory space they require?
- Are these algorithms (reasonably) *efficient* in real-world applications?

The next theorem, which is a compilation of various results from (Tobies, 2001; Eiter et al., 2009; Horrocks et al., 2006),⁶ gives answers to the first two questions:

⁶ See also the DL Complexity Navigator at www.cs.man.ac.uk/~ezolin/dl

Theorem 9. (i) *The satisfiability problem is EXPTIME-complete for \mathcal{ALCHI} KBs and N2EXPTIME-complete for \mathcal{SROIQ} KBs.*

(ii) *Concept and role subsumption and instance checking are EXPTIME- and CON2EXPTIME-complete for, respectively, \mathcal{ALCHI} and \mathcal{SROIQ} KBs.*

(iii) *CQ entailment over \mathcal{ALCHI} KBs is 2EXPTIME-complete.*

Note that full OWL 2 under the RDF-based semantics is undecidable (Motik, 2007), while OWL 2 DL under the direct (model-theoretic) semantics is decidable. There is, however, a price to pay for the additional expressive power of \mathcal{SROIQ} underlying OWL 2 DL: satisfiability is harder than in \mathcal{ALCHI} and CQ entailment is not even known to be decidable.

EXPTIME-completeness of satisfiability in \mathcal{ALCHI} means, in particular, two things: first, there exists a satisfiability-checking algorithm that runs in at most exponential time in the size of the input knowledge base and, second, *no* algorithm can check satisfiability of *any* given knowledge base in polynomial time. This complexity-theoretic result does not suggest that reasoning algorithms can be efficient in practice. Fortunately, practical reasoners for OWL 2 DL have been implemented: FaCT++ (Tsarkov and Horrocks, 2006), Hermit (Horrocks et al., 2012), Pellet (Sirin et al., 2007), Konclude (Steigmiller et al., 2014). Their efficiency in typical real-world applications (rather than worst-case scenarios) relies upon sophisticated optimisation techniques and the empirical fact that ontologies designed by humans for such applications are often simple enough for basic reasoning techniques. On the other hand, answering instance queries and, more generally, conjunctive queries over knowledge bases in expressive languages has not become practical so far, especially for large data sets.

3 Description Logics for the OWL Profiles

The current W3C recommendation OWL 2 of the Web Ontology Language identifies three *profiles* (fragments or sub-languages) specifically designed to ensure efficiency (tractability) of reasoning at the expense of expressive power. In this section, we introduce these profiles in the form of DLs (sacrificing some features for the clarity of presentation).

We begin with an observation that in order to ensure tractability of reasoning (that is, the existence of a *deterministic polynomial-time* algorithm for checking consistency) one has to avoid using \sqcup on the right-hand side of concept inclusions. Intuitively, this construct requires (non-deterministic) reasoning about possible cases such as in Example 7 (case 1: Andrea is a female; case 2: Andrea is a male), which can be NP-hard; for an introduction on the computational complexity, consult classical (Garey and Johnson, 1979) or more recent (Kozen, 2006; Arora and Barak, 2009).

We illustrate non-deterministic reasoning by encoding the NP-complete *graph 3-colourability* problem: given an (undirected) graph $G = (V, E)$, decide whether each of its vertices can be painted in one of the given three colours in such a way that no pair of adjacent vertices has the same colour. We represent the input graph G by means of an ABox \mathcal{A}_G comprising assertions

$$\text{edge}(v_1, v_2), \quad \text{for each } \{v_1, v_2\} \in E.$$

Consider a TBox \mathcal{T} containing the following concept inclusions:

$$\begin{aligned} \top &\sqsubseteq C_1 \sqcup C_2 \sqcup C_3, \\ C_i \sqcap C_j &\sqsubseteq \perp, & 1 \leq i < j \leq 3, \\ C_i \sqcap \exists \text{edge}.C_i &\sqsubseteq \perp, & 1 \leq i \leq 3, \end{aligned}$$

where C_1 , C_2 and C_3 are concept names representing the given three colours. It is not hard to see that each model of $(\mathcal{T}, \mathcal{A}_G)$ gives rise to a 3-colouring of G and, conversely, each 3-colouring of G can be encoded in a model of $(\mathcal{T}, \mathcal{A}_G)$. In other words, $(\mathcal{T}, \mathcal{A}_G)$ is satisfiable if and only if G is 3-colourable. This means that the satisfiability problem for knowledge bases in any DL able to express this TBox is NP-hard (that is, not tractable).

3.1 OWL 2 RL

The OWL 2 RL profile⁷ is aimed at applications requiring scalable reasoning that can be done by rule-based implementations (such as datalog engines). Its design was inspired by the so-called Description Logic Programs (Grosz et al., 2003) and pD* (ter Horst, 2005). OWL 2 RL is supported by Oracle Database 11g, OWLIM, BaseVISor, ELLY, Jena and RDFox (for details and references, see www.w3.org/2001/sw/wiki/OWL/Implementations).

A key feature of OWL 2 RL is that it does not allow existential quantifiers on the right-hand side of concept inclusions. Therefore, when reasoning with an OWL 2 RL knowledge base, we do not have to deal with individuals that are not explicitly present in the knowledge base ABox.

In this section, we consider a somewhat simplified version of OWL 2 RL, which will be called *RL*. Concept and role inclusions in *RL* take the form

$$B \sqsubseteq A, \quad R_1 \sqsubseteq R_2 \quad \text{and} \quad B \sqsubseteq \perp,$$

where R_1 and R_2 are roles (role names or their inverses), A is a concept name and B a concept defined by the following grammar:

$$B ::= A \mid \exists R.\top \mid \exists R.B \mid B_1 \sqcap B_2.$$

Observe that there is no \sqcup in the syntax, and the existential quantifiers occur only on the left-hand side of concept inclusions. On the other hand, by Proposition 8, universal restrictions, intersection and complement on the right-hand side of concept inclusions and union on the left-hand side of concept inclusions are simply syntactic sugar.

An *RL knowledge base* $(\mathcal{T}, \mathcal{A})$ comprises a finite set \mathcal{T} of inclusions introduced above and a *simple* ABox \mathcal{A} , which contains only assertions of the form $A(a)$ and $P(a, b)$, for a concept name A and a role name P .

To understand reasoning in *RL* (and the other two profiles of OWL 2), it is useful to represent its concept and role inclusions as first-order sentences.

⁷ www.w3.org/TR/owl2-profiles/#OWL_2_RL

Example 10. The *RL* concept inclusions

$$\begin{aligned} \exists \text{takesCourse. UndergraduateCourse} &\sqsubseteq \text{UndergraduateStudent}, \\ \text{UndergraduateStudent} &\sqsubseteq \text{Student} \end{aligned}$$

are equivalent (have the same models) as the first-order sentences

$$\begin{aligned} \forall x \forall y (\text{takesCourse}(x, y) \wedge \text{UndergraduateCourse}(y) &\rightarrow \text{UndergraduateStudent}(x)), \\ \forall x (\text{UndergraduateStudent}(x) &\rightarrow \text{Student}(x)). \end{aligned}$$

More formally, we define a *standard translation* ST of concepts and roles by induction on their structure. First, for any concept name A and any role name P , we set

$$\begin{aligned} ST_x(A) &= A(x), \\ ST_{x,y}(P) &= P(x, y), \end{aligned}$$

where the subscript specifies the variables used as arguments of the predicates. After that, we extend the translation ST to complex roles and concepts by taking

$$\begin{aligned} ST_{x,y}(P^-) &= P(y, x), \\ ST_x(\exists R. \top) &= \exists y ST_{x,y}(R), \\ ST_x(\exists R. B) &= \exists y (ST_{x,y}(R) \wedge ST_y(B)), \\ ST_x(B_1 \sqcap B_2) &= ST_x(B_1) \wedge ST_x(B_2). \end{aligned}$$

Finally, we translate concept and role inclusions into universally quantified implications:

$$\begin{aligned} (B \sqsubseteq A)^* &= \forall x (ST_x(B) \rightarrow ST_x(A)), \\ (R_1 \sqsubseteq R_2)^* &= \forall x \forall y (ST_{x,y}(R_1) \rightarrow ST_{x,y}(R_2)), \\ (B \sqsubseteq \perp)^* &= \forall x (ST_x(B) \rightarrow \perp). \end{aligned}$$

As $\exists R. \top$ and $\exists R. B$ can occur only on the left-hand side of concept inclusions, the translation \cdot^* of any *RL* TBox contains only sentences of the form

$$\begin{aligned} \forall \mathbf{y} (\gamma_1(\mathbf{y}) \wedge \cdots \wedge \gamma_k(\mathbf{y}) &\rightarrow \gamma_0(\mathbf{y})), \\ \forall \mathbf{y} (\gamma_1(\mathbf{y}) \wedge \cdots \wedge \gamma_k(\mathbf{y}) &\rightarrow \perp), \end{aligned}$$

where $\gamma_0(\mathbf{y}), \dots, \gamma_k(\mathbf{y})$ are unary or binary predicates with variables in \mathbf{y} . Such sentences are examples of *Horn clauses* (sets of sentences of the first kind are also called *datalog programs*; see e.g., Ceri et al. (1989)). A very important property of Horn clauses known from logic and databases is that everything we may want to know about a knowledge base, whose TBox axioms are Horn clauses, can be found in the canonical model (or chase), which is constructed by ‘applying’ the clauses to the ABox. We first illustrate the construction by a simple example.

Example 11. Consider the TBox \mathcal{T} from Example 10 and the following ABox:

$$\mathcal{A} = \{ \text{takesCourse}(\text{john}, \text{sp1}), \text{UndergraduateCourse}(\text{sp1}) \}.$$

We begin the construction of the canonical model by representing the ABox as an interpretation \mathcal{I}_0 with the domain $\Delta^{\mathcal{I}_0} = \{\text{john}, \text{sp1}\}$ and the interpretation function given by

$$\begin{aligned} \text{takesCourse}^{\mathcal{I}_0} &= \{(\text{john}, \text{sp1})\}, \\ \text{UndergraduateCourse}^{\mathcal{I}_0} &= \{\text{sp1}\}, \\ \text{UndergraduateStudent}^{\mathcal{I}_0} &= \emptyset, \\ \text{Student}^{\mathcal{I}_0} &= \emptyset. \end{aligned}$$

The interpretation \mathcal{I}_0 does not satisfy the first concept inclusion in \mathcal{T} because john is related by takesCourse to sp1 , which is an $\text{UndergraduateCourse}$, but john is not an instance of $\text{UndergraduateStudent}$. To repair this ‘defect’, we apply (as a rule) the first concept inclusion to \mathcal{I}_0 and obtain an interpretation \mathcal{I}_1 with the same domain, $\Delta^{\mathcal{I}_1} = \Delta^{\mathcal{I}_0}$, and the interpretation function $\cdot^{\mathcal{I}_1}$ that expands $\cdot^{\mathcal{I}_0}$ with

$$\text{UndergraduateStudent}^{\mathcal{I}_1} = \{\text{john}\}$$

(all other symbols have the same interpretation as in \mathcal{I}_0). Now, \mathcal{I}_1 satisfies the first concept inclusion of \mathcal{T} but fails to satisfy the second one. We repair this defect by ‘applying’ the second concept inclusion to \mathcal{I}_1 and obtaining an interpretation \mathcal{I}_2 with the same domain and the interpretation function expanding $\cdot^{\mathcal{I}_1}$ with

$$\text{Student}^{\mathcal{I}_2} = \{\text{john}\}$$

(all other symbols keep their interpretation). It is readily seen that now \mathcal{I}_2 is a model of $(\mathcal{T}, \mathcal{A})$. Note that we constructed \mathcal{I}_2 by adding to the given ABox \mathcal{A} only those assertions— $\text{UndergraduateStudent}(\text{john})$ and $\text{Student}(\text{john})$ —that were required by the TBox \mathcal{T} . That is why \mathcal{I}_2 is referred to as the *minimal model* or *canonical model* of $(\mathcal{T}, \mathcal{A})$.

The procedure used in the example above is known as *forward chaining*. Formally, it can be described as follows. First, a simple ABox can be regarded as an interpretation:

Definition 12. The *standard model* $\mathcal{I}_{\mathcal{A}}$ of a simple ABox \mathcal{A} is defined by taking

$$\begin{aligned} \Delta^{\mathcal{I}_{\mathcal{A}}} &= \text{ind}(\mathcal{A}), \\ a^{\mathcal{I}_{\mathcal{A}}} &= a, && \text{for } a \in \text{ind}(\mathcal{A}), \\ A^{\mathcal{I}_{\mathcal{A}}} &= \{a \mid A(a) \in \mathcal{A}\}, && \text{for concept name } A, \\ P^{\mathcal{I}_{\mathcal{A}}} &= \{(a, b) \mid P(a, b) \in \mathcal{A}\}, && \text{for role name } P. \end{aligned}$$

Thus, the domain of $\mathcal{I}_{\mathcal{A}}$ is the set of individual names in \mathcal{A} , and each individual name is interpreted in $\mathcal{I}_{\mathcal{A}}$ by itself, which is often referred to as the *standard name assumption*.

Next, given an *RL* knowledge base $(\mathcal{T}, \mathcal{A})$, we construct a sequence of interpretations $\mathcal{I}_0, \mathcal{I}_1, \dots, \mathcal{I}_n$ by setting $\mathcal{I}_0 = \mathcal{I}_{\mathcal{A}}$ and then applying the following rules to each \mathcal{I}_k to obtain \mathcal{I}_{k+1} :

- (c) if $a \in B^{\mathcal{I}_k}$, $B \sqsubseteq A \in \mathcal{T}$ but $a \notin A^{\mathcal{I}_k}$, then we add a to $A^{\mathcal{I}_{k+1}}$;
- (r) if $(a, b) \in R_1^{\mathcal{I}_k}$, $R_1 \sqsubseteq R_2 \in \mathcal{T}$ but $(a, b) \notin R_2^{\mathcal{I}_k}$, then we add (a, b) to $R_2^{\mathcal{I}_{k+1}}$;
- (b) if $a \in B^{\mathcal{I}_k}$ and $B \sqsubseteq \perp \in \mathcal{T}$, then the process terminates.

Since the domains of the \mathcal{I}_k are finite and all coincide with the set of individuals in \mathcal{A} , the process terminates after a finite number of steps either because neither (c) nor (r) is applicable or because (b) applies. In the former case the resulting interpretation satisfies all the inclusions in \mathcal{T} and all the assertions in \mathcal{A} ; it is called the *canonical model of $(\mathcal{T}, \mathcal{A})$* and denoted by $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. In the latter case, $(\mathcal{T}, \mathcal{A})$ is inconsistent.

Let \mathcal{T}_+ be the *positive* part of an *RL* TBox \mathcal{T} that consists of all role inclusions in \mathcal{T} and all concept inclusions of the form $B \sqsubseteq A$ in \mathcal{T} . By definition, $(\mathcal{T}_+, \mathcal{A})$ is consistent for any \mathcal{A} , and so its canonical model $\mathcal{C}_{\mathcal{T}_+, \mathcal{A}}$ is defined and can be used to check consistency of the full knowledge base (see, e.g., Cali et al. (2012)):

Theorem 13. *An *RL* knowledge base $(\mathcal{T}, \mathcal{A})$ is consistent if and only if we have $\mathcal{C}_{\mathcal{T}_+, \mathcal{A}} \models B \sqsubseteq \perp$, for all $B \sqsubseteq \perp$ in \mathcal{T} .*

Similar results hold for the other two profiles of OWL 2 (or indeed for any Horn DL), and therefore in the following we do not consider negative concept inclusions (although they are part of both OWL 2 EL and OWL 2 QL).

The canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ is *universal* in the sense that it can be homomorphically mapped into any other model of $(\mathcal{T}, \mathcal{A})$ (this notion formalises the minimality mentioned above).

Definition 14. A *homomorphism* h from an interpretation \mathcal{I}_1 to an interpretation \mathcal{I}_2 is a map from $\Delta^{\mathcal{I}_1}$ to $\Delta^{\mathcal{I}_2}$ such that

- $h(a^{\mathcal{I}_1}) = a^{\mathcal{I}_2}$, for each individual name a ,
- $h(u) \in A^{\mathcal{I}_2}$, for any $u \in A^{\mathcal{I}_1}$ and any concept name A ,
- $(h(u), h(v)) \in P^{\mathcal{I}_2}$, for any $(u, v) \in P^{\mathcal{I}_1}$ and any role name P .

The universality of the canonical models means, in particular, that checking subsumption and answering CQs over *RL* knowledge bases can be done by analysing their canonical models:

Theorem 15. *Let $(\mathcal{T}, \mathcal{A})$ be a consistent *RL* knowledge base. Then $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models (\mathcal{T}, \mathcal{A})$. In addition, we have the following:*

- (i) $(\mathcal{T}, \mathcal{A}) \models B \sqsubseteq A$ if and only if $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models B \sqsubseteq A$;
- (ii) $(\mathcal{T}, \mathcal{A}) \models R_1 \sqsubseteq R_2$ if and only if $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models R_1 \sqsubseteq R_2$;
- (iii) $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ if and only if $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models \mathbf{q}(\mathbf{a})$.

Observe that the forward chaining procedure only requires a polynomial number of steps to construct the canonical model (more precisely, the number of steps is bounded by $O(m^2 \times s)$, where m is the number of individual names in the ABox and s is the number of concept and role names). Therefore, knowledge base consistency, subsumption and instance checking are tractable in RL (the matching lower bound will be discussed in Section 4):

Theorem 16. *The problems of knowledge base consistency, concept and role subsumption and instance checking are P-complete in RL . The problem of CQ entailment in RL is NP-complete.*

Forward chaining can be regarded as a bottom-up saturation procedure. In datalog, however, more common is the top-down procedure, where a proof is constructed for a given statement; see e.g., (Ceri et al., 1989). We shall return to this topic in Section 4.3.

It is important to note that, since the clauses in the standard translation of RL inclusions contain only universally quantified variables, the canonical model of any RL knowledge base is *finite* (forward chaining does not add any new elements to the individuals in the ABox). We now turn to the profiles where this is not necessarily the case.

3.2 OWL 2 EL

The design of the OWL 2 EL profile⁸ of OWL 2 (Baader et al., 2008, 2005) was based on the observation that biomedical ontologies such as SNOMED CT,⁹ NCI¹⁰ and GO¹¹ essentially use only conjunctions and existential quantifiers.

In this section, we consider a somewhat simplified version of OWL 2 EL, which will be called EL . Concept and role inclusions in EL look as follows:

$$C_1 \sqsubseteq C_2 \quad \text{and} \quad P_1 \sqsubseteq P_2,$$

where P_1 and P_2 are role names (role inverses are not allowed in EL) and C_1 , C_2 are concepts defined by the following grammar:

$$C ::= A \mid \exists P.T \mid \exists P.C \mid C_1 \sqcap C_2$$

(again, existential restrictions contain only role names). The following are typical concept inclusions from the SNOMED CT ontology:

Pericardium \sqsubseteq Tissue \sqcap containedIn.Heart,
 Pericarditis \sqsubseteq Inflammation \sqcap hasLocation.Pericardium,
 Inflammation \sqsubseteq Disease \sqcap actsOn.Tissue,
 Disease \sqcap hasLocation.containedIn.Heart \sqsubseteq HeartDisease \sqcap NeedsTreatment.

⁸ www.w3.org/TR/owl2-profiles/#OWL_2_EL

⁹ www.ihtsdo.org/snomed-ct

¹⁰ www.obofoundry.org/cgi-bin/detail.cgi?id=ncithesaurus

¹¹ www.geneontology.org

ABoxes in EL are simple; see Section 3.1. The language we defined above is almost identical to \mathcal{ELH} : the only difference is that we do not allow \top to occur outside the scope of existential restrictions. Thus, the EL concepts here are exactly those RL concepts from Section 3.1 that do not contain inverse roles. Note that EL extended with inverse roles is not tractable; moreover reasoning becomes as complex as in \mathcal{ALCHL} , that is, EXPTIME-complete (Baader et al., 2008).

Since the existential restrictions can also occur on the right-hand side of concept inclusions, the result of translating a given EL TBox into first-order logic (see Section 3.1) is no longer a datalog program. Instead, it belongs to an extension of datalog called datalog[±] (Calì et al., 2012) or existential rules (Baget et al., 2011)¹²; in database theory, this language has been known since the early 1980s under the name of tuple-generating dependencies (Abiteboul et al., 1995). More precisely, EL TBoxes can be translated into sets of sentences of the form

$$\forall \mathbf{y} (\gamma_1(\mathbf{y}) \wedge \cdots \wedge \gamma_k(\mathbf{y}) \rightarrow \exists \mathbf{x} \gamma_0(\mathbf{x}, \mathbf{y})),$$

where $\gamma_1(\mathbf{y}), \dots, \gamma_k(\mathbf{y})$ contain only universally quantified variables \mathbf{y} whereas $\gamma_0(\mathbf{x}, \mathbf{y})$ contains both universally quantified variables \mathbf{y} and existentially quantified variables \mathbf{x} (note that the standard translations of EL TBoxes are in fact more restricted than this general form, but it suffices for our explanations). Since these sentences are Horn clauses, we can again apply the forward chaining procedure (chase). However, in the case of EL TBoxes the chase has to ‘invent’ new domain elements for the existential quantifiers. Following the terminology of database theory, these fresh (previously not existing) domain elements will be called *labelled nulls*. In general, the forward chaining procedure may require infinitely many fresh labelled nulls resulting in a possibly infinite canonical model.

The forward chaining construction of the canonical model for an EL knowledge base $(\mathcal{T}, \mathcal{A})$ can be defined by taking the standard model $\mathcal{I}_{\mathcal{A}}$ of the ABox as \mathcal{I}_0 (see Definition 12) and applying inductively the following rules to obtain \mathcal{I}_{k+1} from \mathcal{I}_k :

- (c′) if $d \in C^{\mathcal{I}_k}$ and $C \sqsubseteq A \in \mathcal{T}$, then we add d to $A^{\mathcal{I}_{k+1}}$;
- (r′) if $(d, d') \in P_1^{\mathcal{I}_k}$ and $P_1 \sqsubseteq P_2 \in \mathcal{T}$, then we add (d, d') to $P_2^{\mathcal{I}_{k+1}}$;
- (e) if $d \in C^{\mathcal{I}_k}$ and $C \sqsubseteq \exists P.D \in \mathcal{T}$, where D is a concept name or \top , then we take a *fresh* labelled null, d' , and add d' to $D^{\mathcal{I}_{k+1}}$ and (d, d') to $P^{\mathcal{I}_{k+1}}$

(here, we assume that only $A, \exists P.A$ and $\exists P.\top$ can occur on the right-hand side of concept inclusions—this restriction is inessential, as we shall see in Theorem 20). Note that rules (c′) and (r′) are similar to the rules from Section 3.1 except that they are applied without checking whether d and (d, d') are already present in the interpretation of the concept and role, respectively. This modification of the chase procedure is usually called the *oblivious chase* (Johnson and Klug, 1984); see also (Calì et al., 2013).

¹² See also Chapter 6 in this volume.

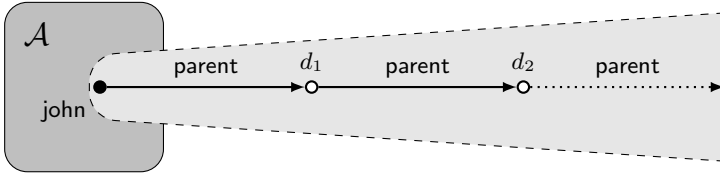


Fig. 5. The infinite canonical model for the KB from Example 17

Example 17. Consider an *EL* TBox \mathcal{T} with the following concept inclusion:

$$\text{Person} \sqsubseteq \exists \text{parent}.\text{Person}.$$

Let us apply the forward chaining procedure to the ABox $\mathcal{A} = \{\text{Person}(\text{john})\}$. We begin by setting $\mathcal{I}_0 = \mathcal{I}_{\mathcal{A}}$:

$$\Delta^{\mathcal{I}_0} = \{\text{john}\}, \quad \text{Person}^{\mathcal{I}_0} = \{\text{john}\} \quad \text{and} \quad \text{parent}^{\mathcal{I}_0} = \emptyset.$$

By the concept inclusion in the TBox, there must be some d_1 that is *parent*-related to *john*. So, we expand \mathcal{I}_0 to \mathcal{I}_1 by taking

$$\Delta^{\mathcal{I}_1} = \{\text{john}, d_1\}, \quad \text{Person}^{\mathcal{I}_1} = \{\text{john}, d_1\} \quad \text{and} \quad \text{parent}^{\mathcal{I}_1} = \{(\text{john}, d_1)\}.$$

Now, the concept inclusion is satisfied for *john* but fails for d_1 , since there must be some d_2 that is *parent*-related to d_1 . So, we expand \mathcal{I}_1 to \mathcal{I}_2 by taking

$$\Delta^{\mathcal{I}_2} = \{\text{john}, d_1, d_2\}, \quad \text{Person}^{\mathcal{I}_2} = \{\text{john}, d_1, d_2\}, \quad \text{parent}^{\mathcal{I}_2} = \{(\text{john}, d_1), (d_1, d_2)\}.$$

If we take all the newly introduced labelled nulls d_i to be distinct then, clearly, this process will continue *ad infinitum*; see Fig. 5. (A possibility of making some of the d_i identical will be discussed later on in this section.)

Since our simplified definition of *EL* does not involve \perp and \neg , every *EL* knowledge base $(\mathcal{T}, \mathcal{A})$ is consistent and the (possibly infinite) forward chaining procedure constructs the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ of $(\mathcal{T}, \mathcal{A})$. As in the case of *RL*, the resulting canonical model is universal: it can be homomorphically mapped into any other model of the knowledge base. Note, however, that universal models are not uniquely defined (for instance, one can take two identical fresh labelled nulls instead of one) and, in contrast to *RL*, some knowledge bases may only have *infinite* universal models. For instance, the knowledge base $(\mathcal{T}, \mathcal{A})$ from Example 17 has no finite universal model. Indeed, suppose, for the sake of contradiction, that there is a finite universal model \mathcal{U} of $(\mathcal{T}, \mathcal{A})$. Then \mathcal{U} must contain a sequence of domain elements connected by *parent* into a cycle. However, the (canonical) model constructed in Example 17 does not contain a homomorphic image of such a cycle, contrary to our assumption.

On the other hand, if we only want to check concept and role subsumption or find answers to instance queries, then we do not have to consider infinite models. The ‘folding’ construction we are going describe below re-uses the labelled nulls and reminds of the filtration technique known from modal logic; see, e.g., (Chagrov and Zakharyashev, 1997).

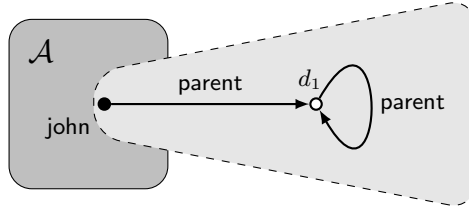


Fig. 6. The small canonical model for the KB from Example 17

Example 18. Consider $(\mathcal{T}, \mathcal{A})$ from Example 17 and suppose that, on step 2 of forward chaining, instead of introducing a fresh labelled null, d_2 , we take d_1 instead. This will result in the following interpretation \mathcal{I}_* (see Fig. 6):

$$\begin{aligned} \Delta^{\mathcal{I}_*} &= \{\text{john}, d_1\}, \\ \text{Person}^{\mathcal{I}_*} &= \{\text{john}, d_1\}, \\ \text{parent}^{\mathcal{I}_*} &= \{(\text{john}, d_1), (d_1, d_1)\}. \end{aligned}$$

Although the interpretation \mathcal{I}_* makes little sense from the modelling point of view (it states that d_1 is its own parent), it is clearly a model of $(\mathcal{T}, \mathcal{A})$. Moreover, as we shall see below, this small model is good enough for checking subsumption and answering instance queries.

To define such small canonical models (also known as generating models) formally, we first convert *EL* TBoxes to a normal form.

Definition 19. An *EL* TBox is said to be in *normal form* if any of its concept inclusions looks as follows:

$$A_1 \sqcap A_2 \sqsubseteq A, \quad \exists P.D \sqsubseteq A \quad \text{or} \quad A \sqsubseteq \exists P.D,$$

where P is a role name, A , A_1 and A_2 are concept names and D is either a concept name or \top .

By introducing abbreviations for complex concepts, one can transform any *EL* TBox to an equivalent one in normal form.

Theorem 20 (Baader et al. (2005)). *Every EL TBox \mathcal{T} can be transformed into a TBox \mathcal{T}' in normal form such that the size of \mathcal{T}' is linear in the size of \mathcal{T} , and \mathcal{T} and \mathcal{T}' are equivalent in the following sense:*

- every model of \mathcal{T} can be extended to a model of \mathcal{T}' by defining interpretations of the fresh concept names,
- every model of \mathcal{T}' is a model of \mathcal{T} .

Example 21. Given a concept inclusion $\exists P.A \sqcap B \sqsubseteq \exists R.\exists P.A$, we first take a fresh concept name C that will stand for $\exists P.A$ and obtain

$$C \sqcap B \sqsubseteq \exists R.C, \quad \underbrace{\exists P.A \sqsubseteq C, \quad C \sqsubseteq \exists P.A.}_{\text{'C is equivalent to } \exists P.A \text{'}}$$

Next, we take another fresh concept name D to replace $\exists R.C$, which results in the following TBox in normal form:

$$C \sqcap B \sqsubseteq D, \quad \exists P.A \sqsubseteq C, \quad C \sqsubseteq \exists P.A, \quad D \sqsubseteq \exists R.C, \quad \exists R.C \sqsubseteq D.$$

We are now in a position to define generating models formally. Let $(\mathcal{T}, \mathcal{A})$ be an *EL* knowledge base in normal form. We say that a concept *occurs positively in* \mathcal{T} if it occurs on the right-hand side of a concept inclusion in \mathcal{T} . For each concept $\exists S.D$ occurring positively in \mathcal{T} (where D is either a concept name or \top), we introduce a *witness* $w_{\exists S.D}$ and define a *generating relation* $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ on the set of these witnesses together with $\text{ind}(\mathcal{A})$ by taking:

$$\begin{aligned} a \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists S.D} & \quad \text{if} \quad a \in \text{ind}(\mathcal{A}) \text{ and } (\mathcal{T}, \mathcal{A}) \models (\exists S.D)(a), \\ w_{\exists P.A} \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists S.D} & \quad \text{if} \quad \mathcal{T} \models A \sqsubseteq \exists S.D. \end{aligned}$$

The *generating model* $\mathcal{G}_{\mathcal{T}, \mathcal{A}}$ for $(\mathcal{T}, \mathcal{A})$ is defined as follows:

$$\begin{aligned} \Delta^{\mathcal{G}_{\mathcal{T}, \mathcal{A}}} &= \text{ind}(\mathcal{A}) \cup \{w_{\exists S.D} \mid \exists S.D \text{ occurs positively in } \mathcal{T}\}, \\ a^{\mathcal{G}_{\mathcal{T}, \mathcal{A}}} &= a, \quad \text{for } a \in \text{ind}(\mathcal{A}), \\ A^{\mathcal{G}_{\mathcal{T}, \mathcal{A}}} &= \{a \in \text{ind}(\mathcal{A}) \mid (\mathcal{T}, \mathcal{A}) \models A(a)\} \cup \\ & \quad \{w_{\exists S.D} \mid \mathcal{T} \models D \sqsubseteq A\}, \quad \text{for a concept name } A, \\ P^{\mathcal{G}_{\mathcal{T}, \mathcal{A}}} &= \{(a, b) \mid S(a, b) \in \mathcal{A}, \mathcal{T} \models S \sqsubseteq P\} \cup \\ & \quad \{(w, w_{\exists S.D}) \mid w \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists S.D}, \mathcal{T} \models S \sqsubseteq P\}, \quad \text{for a role name } P. \end{aligned}$$

It should be clear that $\mathcal{G}_{\mathcal{T}, \mathcal{A}}$ can be constructed in polynomial number of steps by the modified forward chaining procedure that does not invent fresh labelled nulls for $\exists S.D$ but instead re-uses the existing element $w_{\exists S.D}$ in the domain. The following theorem shows that $\mathcal{G}_{\mathcal{T}, \mathcal{A}}$ is indeed a model of $(\mathcal{T}, \mathcal{A})$ and it provides enough information about all concept and role subsumptions and about instance queries.

Theorem 22 (Baader et al. (2005)). *Let $(\mathcal{T}, \mathcal{A})$ be an *EL* knowledge base. Then $\mathcal{G}_{\mathcal{T}, \mathcal{A}} \models (\mathcal{T}, \mathcal{A})$. In addition, we have the following:*

- (i) $(\mathcal{T}, \mathcal{A}) \models C_1 \sqsubseteq C_2$ if and only if $\mathcal{G}_{\mathcal{T}, \mathcal{A}} \models C_1 \sqsubseteq C_2$;
- (ii) $(\mathcal{T}, \mathcal{A}) \models P_1 \sqsubseteq P_2$ if and only if $\mathcal{G}_{\mathcal{T}, \mathcal{A}} \models P_1 \sqsubseteq P_2$;
- (iii) $(\mathcal{T}, \mathcal{A}) \models C(a)$ if and only if $\mathcal{G}_{\mathcal{T}, \mathcal{A}} \models C(a)$;
- (iv) $(\mathcal{T}, \mathcal{A}) \models P(a, b)$ if and only if $\mathcal{G}_{\mathcal{T}, \mathcal{A}} \models P(a, b)$.

Since the generating model $\mathcal{G}_{\mathcal{T}, \mathcal{A}}$ can be constructed in polynomial time in the size of $(\mathcal{T}, \mathcal{A})$, concept and role subsumption are tractable; the same concerns instance checking (matching lower bounds will be discussed in Section 4).

Theorem 23. *The problems of concept and role subsumption and instance checking are P-complete in EL.*

However, the following example shows that $\mathcal{G}_{\mathcal{T},\mathcal{A}}$ cannot be directly used to compute answers to conjunctive queries.¹³

Example 24. Consider $(\mathcal{T}, \mathcal{A})$ from Example 17 and the CQ

$$\mathbf{q} = \exists x \text{parent}(x, x).$$

It should be clear that the answer to \mathbf{q} over $(\mathcal{T}, \mathcal{A})$ is ‘no’ because $\mathcal{C}_{\mathcal{T},\mathcal{A}} \not\models \mathbf{q}$; see Fig. 5. On the other hand, since the generating model $\mathcal{G}_{\mathcal{T},\mathcal{A}}$ contains a loop (see Fig. 6), we have $\mathcal{G}_{\mathcal{T},\mathcal{A}} \models \mathbf{q}$.

We call $\mathcal{G}_{\mathcal{T},\mathcal{A}}$ generating because the standard canonical model $\mathcal{C}_{\mathcal{T},\mathcal{A}}$ of $(\mathcal{T}, \mathcal{A})$ can be generated by unravelling cycles of the generating relation $\rightsquigarrow_{\mathcal{T},\mathcal{A}}$ into infinite trees: for example, $\mathcal{G}_{\mathcal{T},\mathcal{A}}$ in Fig. 6 can be unravelled into $\mathcal{C}_{\mathcal{T},\mathcal{A}}$ in Fig. 5. We remark that the generating model defined here was initially represented as a pair functions by Brandt (2004) and later called the canonical model; see e.g., (Lutz and Wolter, 2007). We prefer the term generating model to avoid confusion with the (possibly infinite) canonical model (the chase).

In Section 4, we shall return to the problem of answering CQs over *EL* knowledge bases and revisit the canonical model construction. In the meantime, we consider the third profile of OWL 2.

3.3 OWL 2 QL

The OWL 2 QL profile¹⁴ of OWL 2 was designed for *ontology-based data access* via query rewriting, where answering CQs over a knowledge base is reduced to answering first-order queries over a database storing the ABox of the KB (this will be discussed in Section 4). OWL 2 QL is based on the logics of the DL-Lite family (Calvanese et al., 2007; Artale et al., 2009).

In this section, we consider a somewhat simplified version of OWL 2 QL, which will be called *QL*. (It is almost identical to what is known as *DL-Lite_R* (Calvanese et al., 2007) or *DL-Lite_{core}^H* (Artale et al., 2009).) Concept and role inclusions in *QL* are of the form

$$B \sqsubseteq C \quad \text{and} \quad R_1 \sqsubseteq R_2,$$

where R_1 and R_2 are roles (role names or their inverses) and B and C are concepts defined by the following grammar:

$$\begin{aligned} B & ::= A \quad | \quad \exists R.T, \\ C & ::= A \quad | \quad \exists R.T \quad | \quad \exists R.C. \end{aligned}$$

¹³ The generating model can be used for answering CQs but the given query has to be modified to take account of identifications of the labelled nulls. This is known as the *combined approach* to query answering, see (Lutz et al., 2009) for the case of *EL* and (Kontchakov et al., 2011) for the case of *QL*. Alternatively, a special procedure has to filter out spurious answers resulting from identification (Lutz et al., 2013).

¹⁴ www.w3.org/TR/owl2-profiles/#OWL_2_QL

ABoxes in QL are simple; see Section 3.1. Note that the universal restrictions are not allowed at all and the existential restrictions on the left-hand side of concept inclusions must have \top as their filler (such existential restrictions are called *unqualified* and the \top symbol is often omitted, making $\exists R$ out of $\exists R.\top$); however, the existential restrictions on the right-hand side of concept inclusions can be qualified. Similarly to the case of EL , the standard first-order translations of QL TBoxes require existential quantification, which causes the canonical models to be infinite.

Analogously to EL TBoxes, one can transform any QL TBox to an equivalent one in normal form (see Theorem 20).

Definition 25. A QL TBox is said to be in *normal form* if any concept inclusion in it looks as follows:

$$A' \sqsubseteq A, \quad \exists R \sqsubseteq A \quad \text{or} \quad A \sqsubseteq \exists R.D,$$

where R is a role, A and A' are concept names and D is either a concept name or \top .

Given a QL knowledge base $(\mathcal{T}, \mathcal{A})$ with \mathcal{T} in normal form, we can find all answers to a CQ q over this KB by evaluating q over the (possibly infinite) canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$, which can be constructed using forward chaining (cf. Section 3.2). We begin by taking the standard model $\mathcal{I}_{\mathcal{A}}$ of the ABox as \mathcal{I}_0 (see Definition 12) and apply inductively the following rules to obtain \mathcal{I}_{k+1} from \mathcal{I}_k :

- (c') if $d \in B^{\mathcal{I}_k}$ and $B \sqsubseteq A \in \mathcal{T}$, then we add d to $A^{\mathcal{I}_{k+1}}$;
- (r') if $(d, d') \in R_1^{\mathcal{I}_k}$ and $R_1 \sqsubseteq R_2 \in \mathcal{T}$, then we add (d, d') to $R_2^{\mathcal{I}_{k+1}}$;
- (e) if $d \in B^{\mathcal{I}_k}$ and $B \sqsubseteq \exists R.D \in \mathcal{T}$, where D is a concept name or \top , then we take a *fresh* labelled null, d' , and add d' to $D^{\mathcal{I}_{k+1}}$ and (d, d') to $R^{\mathcal{I}_{k+1}}$.

(These rules are similar to the rules from Section 3.2 except that they refer to roles with inverses and QL concepts).

The canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ constructed using rules (c'), (r') and (e) in a bottom-up fashion can alternatively be defined by unravelling the generating structure, which is closer to the top-down approach and will be required for query rewriting in Section 4. There are two key observations that lead us to the alternative definition: first, fresh labelled nulls can only be added by applying (e), and, second, if two labelled nulls, d_1 and d_2 , are introduced by applying (e) with the same concept inclusion $B \sqsubseteq \exists R.D$ then the same rules will be applicable to d_1 and d_2 in the continuation of the forward chaining procedure. So, each labelled null d' resulting from applying (e) to some $B \sqsubseteq \exists R.D$ on a domain element d can be identified with a pair of the form $(d, \exists R.D)$. More formally, for each concept $\exists R.D$ that occurs positively in \mathcal{T} , we introduce a *witness* $w_{\exists R.D}$ and define a *generating relation* $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ on the set of these witnesses together with $\text{ind}(\mathcal{A})$ by taking:

$$\begin{aligned} a \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R.D} & \quad \text{if} \quad a \in \text{ind}(\mathcal{A}), \mathcal{I}_{\mathcal{A}} \models B(a) \text{ and } \mathcal{T} \models B \sqsubseteq \exists R.D, \\ w_{\exists S.B} \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R.D} & \quad \text{if} \quad \mathcal{T} \models \exists S^- \sqsubseteq \exists R.D \quad \text{or} \quad \mathcal{T} \models B \sqsubseteq \exists R.D. \end{aligned}$$

A $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ -path σ is a finite sequence $aw_{\exists R_1.D_1} \cdots w_{\exists R_n.D_n}$, $n \geq 0$, such that $a \in \text{ind}(\mathcal{A})$ and, if $n > 0$, then

$$a \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R_1.D_1} \quad \text{and} \quad w_{\exists R_i.D_i} \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R_{i+1}.D_{i+1}}, \quad \text{for } i < n.$$

Thus, a path of the form $\sigma w_{\exists R.D}$ represents the fresh labelled null introduced by applying (e) to some $B \sqsubseteq \exists R.D$ on the domain element σ (and which corresponds to the pair $(\sigma, \exists R.D)$ mentioned above). Denote by $\text{tail}(\sigma)$ the last element in σ ; as we noted above, the last element in σ uniquely determines all the subsequent rule applications. The *canonical model* $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ is defined by taking $\Delta^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ to be the set of all $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ -paths and setting

$$\begin{aligned} a^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}} &= a, & \text{for } a \in \text{ind}(\mathcal{A}), \\ A^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}} &= \{a \in \text{ind}(\mathcal{A}) \mid \mathcal{I}_{\mathcal{A}} \models B(a) \text{ and } \mathcal{T} \models B \sqsubseteq A\} \cup \\ &\quad \{\sigma w_{\exists R.D} \mid \mathcal{T} \models \exists R^- \sqsubseteq A \text{ or } \mathcal{T} \models D \sqsubseteq A\}, \quad \text{for a concept name } A, \\ P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}} &= \{(a, b) \mid \mathcal{I}_{\mathcal{A}} \models R(a, b) \text{ and } \mathcal{T} \models R \sqsubseteq P\} \cup \\ &\quad \{(\sigma w_{\exists R.D}, \sigma) \mid \text{tail}(\sigma) \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R.D}, \mathcal{T} \models R \sqsubseteq P^-\} \cup \\ &\quad \{(\sigma, \sigma w_{\exists R.D}) \mid \text{tail}(\sigma) \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R.D}, \mathcal{T} \models R \sqsubseteq P\}, \quad \text{for a role name } P. \end{aligned}$$

Intuitively, by the definition of rule (c'), an ABox individual a belongs to $A^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ just in case there is a sequence of concepts B_0, B_1, \dots, B_n such that $\mathcal{I}_{\mathcal{A}} \models B_0(a)$, the $B_{i-1} \sqsubseteq B_i$ are in \mathcal{T} , for $1 \leq i \leq n$, and $B_n = A$; in other words, if $\mathcal{I}_{\mathcal{A}} \models B_0(a)$ and $\mathcal{T} \models B_0 \sqsubseteq A$, for some concept B_0 . Similarly, by the definition of rules (c') and (e), a labelled null of the form $\sigma w_{\exists R.D}$ belongs to $A^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ just in case $\mathcal{T} \models \exists R^- \sqsubseteq A$ or $\mathcal{T} \models D \sqsubseteq A$. For a role name P , rules (r') and (e) provide an analogous argument. More precisely, by the definition of rule (r'), a pair (d, d') of domain elements belongs to $P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ just in case there is a sequence of roles R_0, \dots, R_n such that $(d, d') \in R_0^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$, the $R_{i-1} \sqsubseteq R_i$ are in \mathcal{T} , for $1 \leq i \leq n$, and $R_n = P$; in other words, $(d, d') \in R_0^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ and $\mathcal{T} \models R_0 \sqsubseteq P$, for some role R_0 . It then follows from the definition of rule (e) that a pair (d, d') belongs to $P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ just in three cases: (i) both elements of the pair are ABox individuals with $\mathcal{I}_{\mathcal{A}} \models R(d, d')$ and $\mathcal{T} \models R \sqsubseteq P$ (ii) the first component of the pair is created by an application of (e) to the second component of the pair: $d = \sigma w_{\exists R.D}$, $d' = \sigma$ and $\mathcal{T} \models R \sqsubseteq P^-$ or (iii) the second component of the pair is created by an application of (e) to the first component: $d = \sigma$, $d' = \sigma w_{\exists R.D}$ and $\mathcal{T} \models R \sqsubseteq P$. These three cases are reflected in the three sets in the union in the definition of $P^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$.

Example 26. Consider a QL TBox \mathcal{T} with the following concept and role inclusions:

$$\begin{aligned} \text{RA} &\sqsubseteq \exists \text{worksOn.Project}, \\ \text{Project} &\sqsubseteq \exists \text{isManagedBy.Prof}, \\ \text{worksOn}^- &\sqsubseteq \text{involves}, \\ \text{isManagedBy} &\sqsubseteq \text{involves} \end{aligned}$$

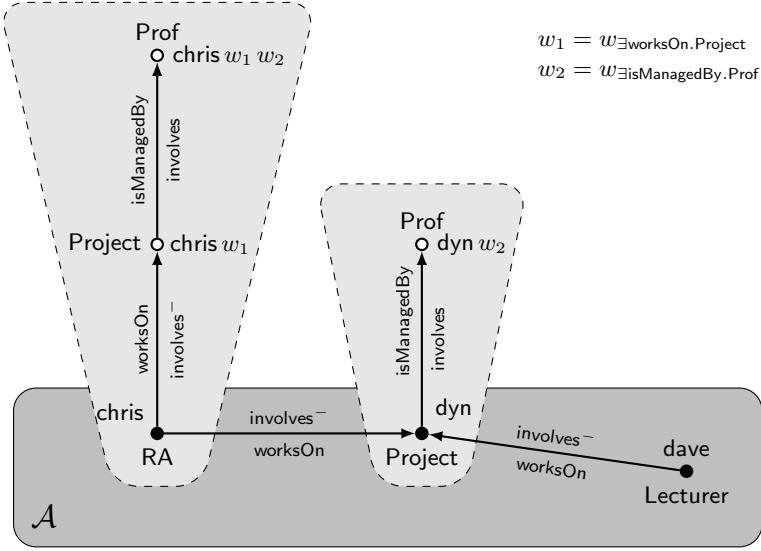


Fig. 7. The canonical model for the knowledge base in Example 26

and an ABox \mathcal{A} comprising the following assertions:

$$\begin{aligned} & \text{RA}(\text{chris}), \quad \text{worksOn}(\text{chris}, \text{dyn}), \quad \text{Project}(\text{dyn}), \\ & \text{Lecturer}(\text{dave}), \quad \text{worksOn}(\text{dave}, \text{dyn}). \end{aligned}$$

Two concepts occur positively in \mathcal{T} : $\exists \text{worksOn. Project}$ and $\exists \text{isManagedBy. Prof}$. For brevity, the witnesses for them will be denoted by w_1 and w_2 , respectively. The generating relation $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ is then defined by taking

$$\text{chris} \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_1, \quad \text{dyn} \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_2, \quad w_1 \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_2.$$

This gives the following $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ -paths:

$$\text{chris}, \quad \text{chris } w_1, \quad \text{chris } w_1 w_2, \quad \text{dyn}, \quad \text{dyn } w_2 \quad \text{and} \quad \text{dave}$$

(note that if the graph of $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ contains a cycle then the set of $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ -paths is infinite; cf. Example 17). The resulting canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ is depicted in Fig. 7.

For any QL knowledge base, the defined canonical model is universal, and thus contains all the necessary information for checking concept and role subsumptions, answering instance queries and, more generally, for answering CQs:

Theorem 27. *Let $(\mathcal{T}, \mathcal{A})$ be a QL knowledge base. Then $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models (\mathcal{T}, \mathcal{A})$ and, for any CQ $q(\mathbf{x})$ and any tuple \mathbf{a} in $\text{ind}(\mathcal{A})$, we have*

$$(\mathcal{T}, \mathcal{A}) \models q(\mathbf{a}) \quad \text{if and only if} \quad \mathcal{C}_{\mathcal{T}, \mathcal{A}} \models q(\mathbf{a}).$$

Although the canonical models in *EL* and *QL* contain all the information required to compute answers to conjunctive queries, these canonical models are not necessarily finite (in contrast to *RL*), and therefore cannot be simply materialised, say, by a datalog engine, a triple store or a database management system. In the next section, we analyse the problem of answering CQs and develop practical techniques for dealing with infinite canonical models.

4 Conjunctive Query Answering via Query Rewriting

Database management systems implement sophisticated algorithms for evaluating SQL queries over relational data instances. The (theoretically and empirically supported) fact that databases have been very efficient in practice suggests the following approach to answering queries over knowledge bases. We store the assertions of a given ABox \mathcal{A} in a relational database. Given a CQ $q(\mathbf{x})$, we use the inclusions of the TBox \mathcal{T} to ‘rewrite’ $q(\mathbf{x})$ into another query $q'(\mathbf{x})$ that would return, when evaluated over the data instance \mathcal{A} , all the certain answers to $q(\mathbf{x})$ over $(\mathcal{T}, \mathcal{A})$. It is important to emphasise here that the rewriting $q'(\mathbf{x})$ must only depend on the TBox \mathcal{T} and the given query $q(\mathbf{x})$, and so should work for *all possible* ABoxes \mathcal{A} . Thus, we arrive to the following definition.

We call a CQ $q(\mathbf{x})$ and a TBox \mathcal{T} *first-order rewritable* (FO-rewritable, for short) if there exists a first-order formula $q'(\mathbf{x})$ such that, for any ABox \mathcal{A} and any tuple \mathbf{a} of individuals in \mathcal{A} , we have

$$(\mathcal{T}, \mathcal{A}) \models q(\mathbf{a}) \quad \text{if and only if} \quad \mathcal{I}_{\mathcal{A}} \models q'(\mathbf{a}), \quad (10)$$

where $\mathcal{I}_{\mathcal{A}}$ is the standard model of \mathcal{A} (see Definition 12). The formula $q'(\mathbf{x})$ is called an *FO-rewriting of q and \mathcal{T}* .

This idea of reducing CQ answering over knowledge bases to database query answering, first formulated by Calvanese et al. (2005), may sound too good to be applicable in all cases. In fact, there are a few issues in realising this idea.

A minor one is that, from a practical point of view, if an FO-rewriting $q'(\mathbf{x})$ is to be executed by a relational database engine then $q'(\mathbf{x})$ must be a domain-independent query (Abiteboul et al., 1995). This is the case, in particular, for FO-rewritings that contain only conjunction, disjunction and existential quantifiers. Such rewritings (and formulas) are called *positive existential rewritings* (PE-rewritings, for short).

A more serious issue is that not all DL constructs can guarantee FO-rewritability of all CQs. To understand why, let us recall (Libkin, 2004) that the problem of evaluating a first-order formula in a given interpretation belongs to the class AC^0 for *data complexity*. Data complexity is a complexity measure that only takes account of the size of the data (the interpretation in this case) and regards the query to be fixed. This measure was suggested by Vardi (1982) who tried to find a theoretical explanation of the practical efficiency of database management systems. It is also known from the complexity theory that AC^0 is a *proper* subclass of LOGSPACE and that $LOGSPACE \subseteq NLOGSPACE \subseteq P$ (Papadimitriou, 1994; Arora and Barak, 2009). It follows that if the problem ‘ $(\mathcal{T}, \mathcal{A}) \models q?$ ’ is

NLOGSPACE- or P-hard, for some fixed \mathbf{q} and \mathcal{T} , then these \mathbf{q} and \mathcal{T} cannot be FO-rewritable. This observation allows us to identify the DL constructs that can ruin FO-rewritability. Here we give two simple examples illustrating this technique; for more consult, e.g., (Calvanese et al., 2006; Artale et al., 2009).

Example 28. A typical example of an NLOGSPACE-complete problem is the *reachability problem* for directed graphs: given a directed graph $G = (V, E)$ with vertices V and arcs E and two distinguished vertices $s, t \in V$, decide whether there is a directed path from s to t in G . We represent the graph and the target vertex by means of an ABox $\mathcal{A}_{G,t}$ comprising

$$\begin{aligned} &\text{edge}(v_1, v_2), \quad \text{for } (v_1, v_2) \in E, \\ &\text{ReachableFromTarget}(t). \end{aligned}$$

Consider now a TBox containing concept inclusion

$$\exists \text{edge}.\text{ReachableFromTarget} \sqsubseteq \text{ReachableFromTarget}$$

and the Boolean CQ

$$\mathbf{q} = \text{ReachableFromTarget}(s).$$

It is readily seen that $(\mathcal{T}, \mathcal{A}_{G,t}) \models \mathbf{q}$ if and only if there is a directed path from s to t in G . Therefore, the problem ‘ $(\mathcal{T}, \mathcal{A}_{G,t}) \models \mathbf{q}$?’ is NLOGSPACE-hard. Since \mathbf{q} and \mathcal{T} are fixed (and do not depend on G or t), \mathbf{q} and \mathcal{T} cannot be FO-rewritable. In other words, TBoxes capable of computing the transitive closure of some relations in ABoxes do not allow FO-rewritability.

Example 29. Next, we consider an example of a P-complete problem. Let (V, E) be a pair that consists of a finite set of vertices V and a relation $E \subseteq V \times V \times V$. A vertex $v \in V$ is said to be *accessible* from a set $S \subseteq V$ of source vertices in (V, E) if either $v \in S$ or $(v_1, v_2, v) \in E$, for some v_1 and v_2 that are accessible from S in (V, E) (v_1 and v_2 are called inputs and v the output; such a triple can also be thought of as an implication of the form $v_1 \wedge v_2 \rightarrow v$). The *path system accessibility* problem is defined as follows: given (V, E) as above, source vertices $S \subseteq V$ and a terminal vertex $t \in V$, decide whether t is accessible from S in (V, E) . This problem is known to be P-complete (Garey and Johnson, 1979). The path system (V, E) and source vertices S can be encoded by an ABox $\mathcal{A}_{V,E,S}$ in the following way:

$$\begin{aligned} &\text{Accessible}(v), && \text{for } v \in S, \\ &\text{input}_1(e, v_1), \text{ input}_2(e, v_2) \text{ and } \text{output}(v, e), && \text{for } e = (v_1, v_2, v) \in E. \end{aligned}$$

Consider now a TBox \mathcal{T} containing

$$\begin{aligned} &\exists \text{input}_1.\text{Accessible} \sqcap \exists \text{input}_2.\text{Accessible} \sqsubseteq \text{BothInputsAccessible}, \\ &\exists \text{output}.\text{BothInputsAccessible} \sqsubseteq \text{Accessible} \end{aligned}$$

and the Boolean CQ

$$\mathbf{q} = \text{Accessible}(t).$$

It should be clear that $(\mathcal{T}, \mathcal{A}_{V,E,S}) \models \text{Accessible}(v)$ if and only if v is accessible from S in (V, E) ; and $(\mathcal{T}, \mathcal{A}) \models \text{BothInputsAccessible}(e)$ if and only if *both* inputs of e are accessible, that is, both are instances of `Accessible`. Therefore, the answer to \mathbf{q} is ‘yes’ if and only if t is accessible from S in (V, E) . Thus, the problem ‘ $(\mathcal{T}, \mathcal{A}_{V,E,S}) \models \mathbf{q}$ ’ is P-hard, and so these fixed \mathbf{q} and \mathcal{T} cannot be FO-rewritable.

The OWL2QL profile of OWL2 was designed so that problems such as graph reachability and path system accessibility above could not be expressed in it (on the other hand, observe that both of the TBoxes above belong to *RL* and *EL*, which proves the lower complexity bounds in Theorems 16 and 23). A number of various rewriting techniques have been proposed and implemented for OWL2QL: PerfectRef (Poggi et al., 2008), Presto / Prexto (Rosati and Almatelli, 2010; Rosati, 2012), as well as for extensions of OWL2QL to datalog[±] and existential rules: Nyaya (Gottlob et al., 2011) and PURE (König et al., 2012) and more expressive DLs: Requiem / Blackout (Pérez-Urbina et al., 2009, 2012), Rapid (Chortaras et al., 2011) and Clipper (Eiter et al., 2012), which go beyond FO-rewritability.

In this section, we discuss the tree-witness rewriting (Kikot et al., 2012b). We require the following definitions in the sequel. Whenever convenient, we write $S(\mathbf{z})$ for either a unary atom $A(z_1)$ or a binary atom $P(z_1, z_2)$; we also identify $P^-(z_2, z_1)$ with $P(z_1, z_2)$. Any CQ $\mathbf{q}(\mathbf{x}) = \exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$ is regarded as the set of atoms in φ , so we can write $S(\mathbf{z}) \in \mathbf{q}$ (when referring to the query as a set of atoms, we often omit the answer variables). Any set \mathbf{q} of atoms can also be viewed as an interpretation over the domain of its terms (variables and individual names) such that an atom $S(\mathbf{z})$ is true in this interpretation just in case $S(\mathbf{z}) \in \mathbf{q}$. We slightly abuse notation and denote by \mathbf{q} both the set of atoms and the corresponding interpretation. The reason behind these definitions and notations is as follows: it is not hard to see that $\mathcal{I} \models \mathbf{q}(\mathbf{a})$ if and only if there is a homomorphism from $\mathbf{q}(\mathbf{a})$ to \mathcal{I} (see Definition 14).

4.1 PE-Rewriting for Flat QL (and RDFS)

We first consider an important special case of *flat QL* TBoxes that do not contain existential quantifiers on the right-hand side of concept inclusions. In other words, flat *QL* TBoxes in normal form can only contain concept and role inclusions of the form

$$A \sqsubseteq A', \quad \exists R \sqsubseteq A \quad \text{and} \quad R \sqsubseteq R',$$

for concept names A and A' and roles R and R' . Note that the language of flat *QL* TBoxes differs from the language of RDFS¹⁵ only in that *QL* allows inverse roles in role inclusions whereas RDFS restricts role inclusions to just role names.

¹⁵ www.w3.org/TR/rdf-schema

Let \mathcal{T} be a flat *QL* TBox and $\mathbf{q}(\mathbf{x})$ a conjunctive query. By Theorem 27, $(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a})$ if and only if $\mathbf{q}(\mathbf{a})$ is true in the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. Since the TBox is flat, the generating relation $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ is empty, the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ contains no labelled nulls, and so, by the definition of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$, we have

$$\begin{aligned} \mathcal{C}_{\mathcal{T}, \mathcal{A}} \models A(a) & \quad \text{if and only if} \quad \mathcal{I}_{\mathcal{A}} \models B(a) \text{ and } \mathcal{T} \models B \sqsubseteq A, \text{ for some } B, \\ \mathcal{C}_{\mathcal{T}, \mathcal{A}} \models P(a, b) & \quad \text{if and only if} \quad \mathcal{I}_{\mathcal{A}} \models R(a, b) \text{ and } \mathcal{T} \models R \sqsubseteq P, \text{ for some } R. \end{aligned}$$

Define a PE-formula $\mathbf{q}_{\text{ext}}(\mathbf{x})$ to be the result of replacing every atom $A(z)$ in $\mathbf{q}(\mathbf{x})$ with $A_{\text{ext}}(z)$ and every atom $P(z_1, z_2)$ in $\mathbf{q}(\mathbf{x})$ with $P_{\text{ext}}(z_1, z_2)$, where

$$A_{\text{ext}}(u) = \bigvee_{\mathcal{T} \models B \sqsubseteq A} ST_u(B), \quad P_{\text{ext}}(u, v) = \bigvee_{\mathcal{T} \models R \sqsubseteq P} ST_{u,v}(R),$$

and ST is the standard translation defined in Section 3.1. It is not hard to see that, for any ABox \mathcal{A} and any tuple \mathbf{a} from $\text{ind}(\mathcal{A})$, we have $\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models \mathbf{q}(\mathbf{a})$ if and only if $\mathcal{I}_{\mathcal{A}} \models \mathbf{q}_{\text{ext}}(\mathbf{a})$.

Proposition 30. *For any CQ $\mathbf{q}(\mathbf{x})$ and any flat *QL* TBox \mathcal{T} , $\mathbf{q}_{\text{ext}}(\mathbf{x})$ is a PE-rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} .*

Thus, in the case of flat *QL* TBoxes, it is really easy to construct PE-rewritings.

Example 31. Consider the CQ

$$\mathbf{q}(x, y) = \text{Student}(x) \wedge \text{takesCourse}(x, y) \wedge \text{teacherOf}(p0, y)$$

and a flat *QL* TBox \mathcal{T} with the following concept and role inclusions:

$$\begin{aligned} \text{UndergraduateStudent} & \sqsubseteq \text{Student}, \\ \text{enrolledAt} & \sqsubseteq \text{Student}, \\ \text{teaches}^- & \sqsubseteq \text{teacherOf}. \end{aligned}$$

We then define the following formulas for concept and role names from the query:

$$\begin{aligned} \text{Student}_{\text{ext}}(u) & = \text{Student}(u) \vee \text{UndergraduateStudent}(u) \vee \\ & \quad \exists v \text{enrolledAt}(u, v), \\ \text{takesCourse}_{\text{ext}}(u, v) & = \text{takesCourse}(u, v), \\ \text{teacherOf}_{\text{ext}}(u, v) & = \text{teacherOf}(u, v) \vee \text{teaches}(v, u), \end{aligned}$$

and we obtain the following PE-rewriting of $\mathbf{q}(x, y)$ and \mathcal{T} :

$$\begin{aligned} \mathbf{q}_{\text{ext}}(x, y) & = (\text{Student}(x) \vee \text{UndergraduateStudent}(x) \vee \exists v \text{enrolledAt}(x, v)) \wedge \\ & \quad \text{takesCourse}(x, y) \wedge (\text{teacherOf}(p0, y) \vee \text{teaches}(y, p0)). \end{aligned}$$

Next, we turn to the case of general *QL* TBoxes.

4.2 Tree-Witness PE-Rewriting for Full QL

Suppose \mathcal{T} is a QL TBox in normal form. By Theorem 27, to compute certain answers to $\mathbf{q}(\mathbf{x})$ over $(\mathcal{T}, \mathcal{A})$, for some \mathcal{A} , it is enough to find answers to $\mathbf{q}(\mathbf{x})$ in the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. To do this, we have to check, for every tuple \mathbf{a} of elements in $\text{ind}(\mathcal{A})$, whether there exists a homomorphism from $\mathbf{q}(\mathbf{a})$ to $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. Thus, as in the case of flat TBoxes, the answer variables take values from $\text{ind}(\mathcal{A})$. However, the existentially quantified variables in $\mathbf{q}(\mathbf{x})$ can be mapped both to $\text{ind}(\mathcal{A})$ and to the labelled nulls in $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. In order to define the rewriting that does not depend on a particular ABox, we need to have a closer look at the structure of the canonical models $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ with fixed \mathcal{T} and varying \mathcal{A} .

Let $\exists R.D$ be a concept occurring positively in \mathcal{T} (recall that D is either a concept name or \top). For an individual name a , we define the $\exists R.D$ -generated \mathcal{T} -tree on a as the restriction of the canonical model of the KB $(\mathcal{T}, \{(\exists R.D)(a)\})$ to the domain that consists of a and the labelled nulls with the prefix $aw_{\exists R.D}$. We denote this tree by $\mathcal{C}_{\mathcal{T}}^{\exists R.D}(a)$. (Note that $\mathcal{C}_{\mathcal{T}}^{\exists R.D}(a)$ is not necessarily a model of \mathcal{T} .) Take now any ABox \mathcal{A} and any $a \in \text{ind}(\mathcal{A})$. By the definition of the canonical model, if $a \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists R.D}$ then $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ contains a sub-tree that is isomorphic to the $\exists R.D$ -generated \mathcal{T} -tree on a , excluding possibly the root. Moreover, such sub-trees may intersect only on their common root a . For instance, the canonical model in Fig. 7 contains the $\exists \text{worksOn.Project}$ -generated \mathcal{T} -tree on `chris` and the $\exists \text{isManagedBy.Prof}$ -generated \mathcal{T} -tree on `dyn`. The following example illustrates a more complex configuration.

Example 32. Consider a TBox \mathcal{T} with the following concept inclusions:

$$\begin{aligned} A &\sqsubseteq \exists R.D, & D &\sqsubseteq \exists P_1, & D &\sqsubseteq \exists P_2, \\ B &\sqsubseteq \exists S.C \end{aligned}$$

and suppose that an ABox \mathcal{A} contains $A(a), P_1(a, b), A(b), B(b), P_2(b, c)$. The canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ is depicted in Fig. 8. The individual a in this canonical model has a single \mathcal{T} -tree generated by $\exists R.D$. The individual b has two \mathcal{T} -trees, one generated by $\exists R.D$ and another by $\exists S.C$. These two \mathcal{T} -trees intersect only on their common root b .

For a more formal treatment, we require the following operation. Given interpretations \mathcal{I}_1 and \mathcal{I}_2 (under the standard name assumption), we define their *join*, $\mathcal{I}_1 \oplus \mathcal{I}_2$, by taking

$$\begin{aligned} \Delta^{\mathcal{I}_1 \oplus \mathcal{I}_2} &= \Delta^{\mathcal{I}_1} \cup \Delta^{\mathcal{I}_2}, \\ a^{\mathcal{I}_1 \oplus \mathcal{I}_2} &= a, && \text{for an individual } a \text{ in } \mathcal{I}_1 \text{ or } \mathcal{I}_2, \\ A^{\mathcal{I}_1 \oplus \mathcal{I}_2} &= A^{\mathcal{I}_1} \cup A^{\mathcal{I}_2}, && \text{for a concept name } A, \\ P^{\mathcal{I}_1 \oplus \mathcal{I}_2} &= P^{\mathcal{I}_1} \cup P^{\mathcal{I}_2}, && \text{for a role name } P. \end{aligned}$$

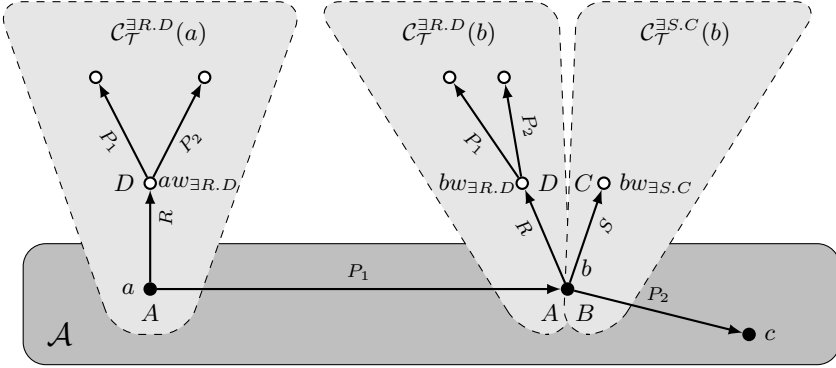


Fig. 8. The canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ from Example 32

Then the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ of any QL knowledge base $(\mathcal{T}, \mathcal{A})$, with \mathcal{T} in normal form, can be represented as the following join:

$$\mathcal{C}_{\mathcal{T}, \mathcal{A}} = \mathcal{I}_{\mathcal{A}}^* \oplus \bigoplus_{\substack{a \in \text{ind}(\mathcal{A}) \\ B \sqsubseteq \exists R.D \in \mathcal{T} \text{ with } \mathcal{I}_{\mathcal{A}} \models B(a)}} \mathcal{C}_{\mathcal{T}}^{\exists R.D}(a), \quad (11)$$

where $\mathcal{I}_{\mathcal{A}}^*$ is a model of \mathcal{A} with domain $\text{ind}(\mathcal{A})$, which will be called the *ABox part* of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$; the join of the $\mathcal{C}_{\mathcal{T}}^{\exists R.D}(a)$ will be called the *anonymous part* of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. We are now fully equipped to present the tree-witness rewriting of a CQ $q(\mathbf{x})$ and a QL TBox \mathcal{T} .

Following the divide and conquer strategy, we show how the process of constructing FO-rewritings can be split into two steps: the first step considers only the flat part of the TBox and uses the formulas $A_{\text{ext}}(u)$ and $P_{\text{ext}}(u, v)$ defined in Section 4.1; the second step (to be described below) takes account of the remaining part of the TBox, that is, inclusions of the form $B \sqsubseteq \exists R.D$.

H-completeness. Let \mathcal{T} be a (not necessarily flat) QL TBox. A simple ABox \mathcal{A} is said to be *H-complete with respect to \mathcal{T}* if, for all concept names A and role names P , we have

$$\begin{aligned} A(a) \in \mathcal{A} & \quad \text{if} \quad \mathcal{I}_{\mathcal{A}} \models B(a) \text{ and } \mathcal{T} \models B \sqsubseteq A, \text{ for some } B = A' \text{ or } B = \exists R, \\ P(a, b) \in \mathcal{A} & \quad \text{if} \quad \mathcal{I}_{\mathcal{A}} \models R(a, b) \text{ and } \mathcal{T} \models R \sqsubseteq P, \text{ for some } R. \end{aligned}$$

We say that a first-order formula $q'(\mathbf{x})$ is an *FO-rewriting of $q(\mathbf{x})$ and \mathcal{T} over H-complete ABoxes* if (10) holds for any ABox \mathcal{A} that is H-complete with respect to \mathcal{T} and any tuple \mathbf{a} from $\text{ind}(\mathcal{A})$; as before, a PE-rewriting is an FO-rewriting that contains only conjunction, disjunction and existential quantification.

Observe that if an ABox \mathcal{A} is H-complete with respect to \mathcal{T} then the ABox part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$, that is, $\mathcal{I}_{\mathcal{A}}^*$ in (11), coincides with $\mathcal{I}_{\mathcal{A}}$. Thus, if \mathcal{T} is flat then $q(\mathbf{x})$ itself is clearly a PE- and FO-rewriting of $q(\mathbf{x})$ and \mathcal{T} over H-complete ABoxes.

More generally, we can easily obtain rewritings (over arbitrary ABoxes) from rewritings over H-complete ABoxes:

Proposition 33. *If $q'(x)$ is an FO-rewriting (PE-rewriting) of $q(x)$ and \mathcal{T} over H-complete ABoxes, then $q'_{\text{ext}}(x)$ is an FO-rewriting (respectively, PE-rewriting) of $q(x)$ and \mathcal{T} .*

Thus, we can only focus on constructing PE-rewritings over H-complete ABoxes.

Tree Witnesses. Consider a CQ $q(x)$ and a knowledge base $(\mathcal{T}, \mathcal{A})$. Suppose that, for some tuple a in $\text{ind}(\mathcal{A})$, there is a homomorphism h from $q(a)$ to $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. Then h partitions $q(a)$ into the atoms mapped by h to the ABox part and atoms mapped by h to the $\exists R.D$ -generated \mathcal{T} -trees of the anonymous part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. The tree-witness rewriting of $q(x)$ and \mathcal{T} we are going to present now lists all possible partitions of the atoms of $q(x)$ into such subsets. We begin with an example illustrating this idea.

Example 34. Consider the QL TBox \mathcal{T} from Example 26 with the concept and role inclusions

$$\text{RA} \sqsubseteq \exists \text{worksOn.Project}, \quad (12)$$

$$\text{Project} \sqsubseteq \exists \text{isManagedBy.Prof}, \quad (13)$$

$$\text{worksOn}^- \sqsubseteq \text{involves}, \quad (14)$$

$$\text{isManagedBy} \sqsubseteq \text{involves} \quad (15)$$

and the CQ asking to find those who work with professors:

$$q(x) = \exists y, z (\text{worksOn}(x, y) \wedge \text{involves}(y, z) \wedge \text{Prof}(z)).$$

Recall that if the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$, for some ABox \mathcal{A} , contains some individuals $a \in \text{RA}^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ and $b \in \text{Project}^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$, then $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ must also contain the $\exists \text{worksOn.Project}$ -generated \mathcal{T} -tree on a and the $\exists \text{isManagedBy.Prof}$ -generated \mathcal{T} -tree on b ; see Fig. 7, where such an a is *chris* and such a b is *dyn*. Let us consider all possible ways of obtaining certain answers to the query, that is, all possible homomorphisms from atoms of $q(x)$ to $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ such that the answer variable x is mapped to $\text{ind}(\mathcal{A})$. First, x , y and z can be mapped to ABox individuals. Alternatively, x and y can be mapped to ABox individuals, a and b , and if b belongs to $\text{Project}^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$, then there is a homomorphism h_1 from the last two atoms of $q(a)$ to the anonymous part; see Fig. 9. Another option is to map only x to an ABox individual, a , and if a belongs to $\text{RA}^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ then the whole $q(a)$ can be homomorphically mapped to the anonymous part; see h_2 in Fig. 9. Finally, another homomorphism, h_3 in Fig. 9, maps both x and z to a provided that a is in $\text{RA}^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ and $\text{Prof}^{\mathcal{C}_{\mathcal{T}, \mathcal{A}}}$ at the same time. The possible ways of mapping subsets of a query to the anonymous part of the canonical model are called *tree*

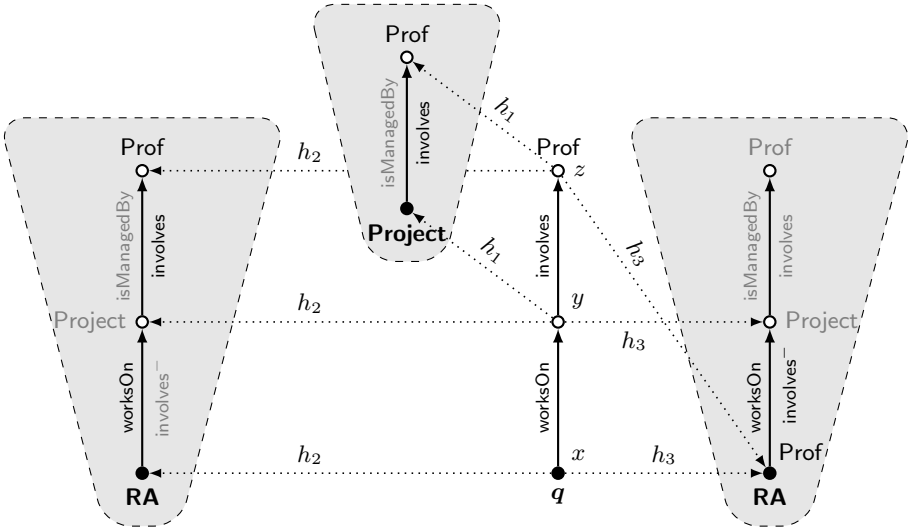


Fig. 9. Three homomorphisms from subsets of q to \mathcal{T} -trees

witnesses. The three tree witnesses for $q(x)$ and \mathcal{T} found above give rise to the following PE-rewriting $q_{tw}(x)$ of $q(x)$ and \mathcal{T} over H-complete ABoxes:

$$q_{tw}(x) = \exists y, z \left[(\text{worksOn}(x, y) \wedge \text{involves}(y, z) \wedge \text{Prof}(z)) \vee (\text{worksOn}(x, y) \wedge \text{Project}(y)) \vee \text{RA}(x) \vee (\text{RA}(x) \wedge \text{Prof}(z) \wedge (x = z)) \right].$$

We now give a general definition of the tree-witness rewriting over H-complete ABoxes. Let \mathcal{T} be a QL TBox in normal form and $q(x)$ a CQ. Consider a pair $t = (t_r, t_i)$ of disjoint sets of terms in $q(x)$, where t_i is non-empty and contains only existentially quantified variables (t_r , on the other hand, can be empty and can contain answer variables and individual names). Set

$$q_t = \{ S(z) \in q \mid z \subseteq t_r \cup t_i \text{ and } z \not\subseteq t_r \}.$$

We say that t is a tree witness for $q(x)$ and \mathcal{T} generated by $\exists R.D$ if the following two conditions are satisfied:

- (tw1) there exists a homomorphism h from q_t to $\mathcal{C}_{\mathcal{T}}^{\exists R.D}(a)$, for some a , such that $t_r = \{z \mid h(z) = a\}$ and t_i contains the remaining variables in q_t ,
- (tw2) q_t is a minimal subset of q such that, for any $y \in t_i$, every atom in q containing y belongs to q_t .

Note that unary atoms with arguments in t_r or binary atoms with both arguments in t_r do not belong to q_t and, therefore, condition (tw1) does not require them

to be homomorphically mapped into $\mathcal{C}_{\mathcal{T}}^{\exists R.D}(a)$. The terms in \mathfrak{t}_r (if any) are called the *roots* of \mathfrak{t} and the (existentially quantified) variables in \mathfrak{t}_i the *interior* of \mathfrak{t} . The homomorphism h in condition **(tw1)** is not necessarily unique; however, it is important that all roots are mapped to a and all variables of the interior are *not* mapped to a . Thus, $\mathfrak{q}_{\mathfrak{t}}$ can contain at most one individual name, a ; if $\mathfrak{q}_{\mathfrak{t}}$ does not contain an individual name then the choice of a is irrelevant. Condition **(tw2)** reflects the fact that if a homomorphism from \mathfrak{q} to the canonical model of $(\mathcal{T}, \mathcal{A})$, for some \mathcal{A} , maps a variable y of an atom $R(y, z)$ to a non-root of a tree $\mathcal{C}_{\mathcal{T}}^{\exists R.D}(a)$ then the other variable of the atom must be mapped to the same $\exists R.D$ -generated \mathcal{T} -tree on a .)

Let $\mathfrak{t} = (\mathfrak{t}_r, \mathfrak{t}_i)$ be a tree witness for $\mathfrak{q}(\mathbf{x})$ and \mathcal{T} . Consider the following formula

$$\text{tw}_{\mathfrak{t}} = \exists u \left[\bigwedge_{x \in \mathfrak{t}_r} (x = u) \quad \wedge \quad \bigvee_{\substack{B \sqsubseteq \exists R.D \in \mathcal{T} \\ \mathfrak{t} \text{ generated by } \exists R.D}} ST_u(B) \right], \quad (16)$$

whose free variables are the roots, \mathfrak{t}_r , of \mathfrak{t} . The formula $\text{tw}_{\mathfrak{t}}$ describes the ABox individuals that root the trees in the anonymous part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ into which the atoms $\mathfrak{q}_{\mathfrak{t}}$ of the tree witness \mathfrak{t} can be homomorphically mapped. More formally, if $\mathcal{I}_{\mathcal{A}} \models \text{tw}_{\mathfrak{t}}(a, \dots, a)$, for some $a \in \text{ind}(\mathcal{A})$, then $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ contains the $\exists R.D$ -generated \mathcal{T} -tree on a , and so there is a homomorphism from $\mathfrak{q}_{\mathfrak{t}}$ to $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ that maps all the roots of \mathfrak{t} to a . Conversely, if there is a homomorphism from $\mathfrak{q}_{\mathfrak{t}}$ to $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ such that all the roots of \mathfrak{t} are mapped to a (but all the variables from the interior, \mathfrak{t}_i , of \mathfrak{t} are mapped to labelled nulls) then $\mathcal{I}_{\mathcal{A}} \models \text{tw}_{\mathfrak{t}}(a, \dots, a)$.

Let $\Theta_{\mathcal{T}}^{\mathfrak{q}}$ be the set of tree witnesses for $\mathfrak{q}(\mathbf{x})$ and \mathcal{T} . Tree witnesses \mathfrak{t} and \mathfrak{t}' are said to be *conflicting* if $\mathfrak{q}_{\mathfrak{t}} \cap \mathfrak{q}_{\mathfrak{t}'} \neq \emptyset$ (in other words, the interior of one tree witness, say, \mathfrak{t} , contains a root or an interior variable of the other, \mathfrak{t}' , or the other way round, which makes it impossible to have both tree witnesses mapped into the anonymous part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ at the *same time*). A set $\Theta \subseteq \Theta_{\mathcal{T}}^{\mathfrak{q}}$ of tree witnesses is said to be *independent* if any two distinct tree witnesses in Θ are non-conflicting. If Θ is independent then we can ‘cut’ the query $\mathfrak{q}(\mathbf{x})$ into independent subqueries in the following way. Consider a homomorphism that, for each $\mathfrak{t} \in \Theta$, maps the subset $\mathfrak{q}_{\mathfrak{t}}$ of \mathfrak{q} to the $\exists R.D$ -generated \mathcal{T} -tree on some a (provided that \mathfrak{t} is generated by $\exists R.D$) and maps the remaining atoms in \mathfrak{q} to the ABox part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. By (11), such a homomorphism is possible if there is a tuple \mathbf{a} in $\text{ind}(\mathcal{A})$ such that the formula

$$\mathfrak{q}_{\text{cut}}^{\Theta}(\mathbf{x}) = \exists \mathbf{y} \left((\mathfrak{q} \setminus \mathfrak{q}_{\Theta}) \wedge \bigwedge_{\mathfrak{t} \in \Theta} \text{tw}_{\mathfrak{t}} \right) \quad (17)$$

holds in $\mathcal{I}_{\mathcal{A}}$ on \mathbf{a} , where $\mathfrak{q} \setminus \mathfrak{q}_{\Theta}$ is the conjunction of all the atoms in \mathfrak{q} that do not belong to $\mathfrak{q}_{\mathfrak{t}}$, for any $\mathfrak{t} \in \Theta$. (Recall that due to H-completeness of the ABox, $\mathcal{I}_{\mathcal{A}}^* = \mathcal{I}_{\mathcal{A}}$.) Conversely, if there is a homomorphism from $\mathfrak{q}(\mathbf{a})$ to $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ then there exists an independent set Θ of tree witnesses such that $\mathcal{I}_{\mathcal{A}} \models \mathfrak{q}_{\text{cut}}^{\Theta}(\mathbf{a})$.

The following PE-formula $\mathfrak{q}_{\text{tw}}(\mathbf{x})$ is called the *tree-witness rewriting of $\mathfrak{q}(\mathbf{x})$ and \mathcal{T} over H-complete ABoxes*:

$$\mathfrak{q}_{\text{tw}}(\mathbf{x}) = \bigvee_{\Theta \subseteq \Theta_{\mathcal{T}}^{\mathfrak{q}} \text{ independent}} \mathfrak{q}_{\text{cut}}^{\Theta}(\mathbf{x}). \quad (18)$$

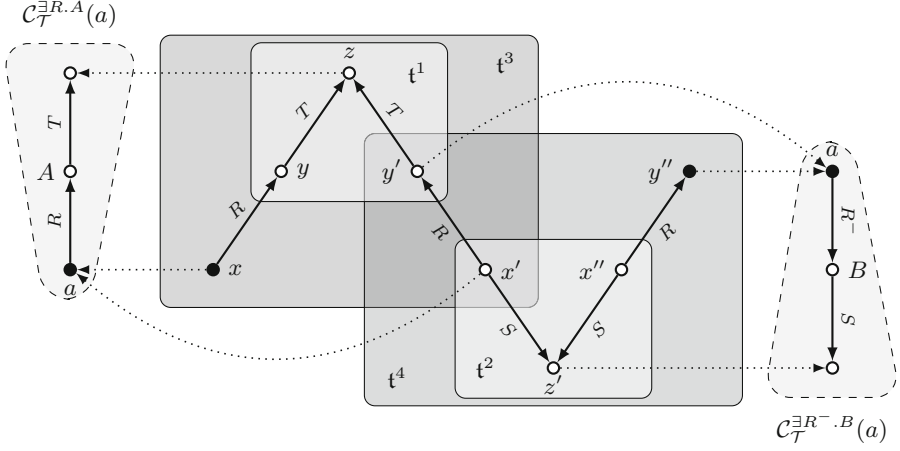


Fig. 10. Tree witnesses in Example 35

Example 35. Consider a TBox with the following concept inclusions:

$$\begin{array}{ll} A_0 \sqsubseteq \exists R.A, & A \sqsubseteq \exists T, \\ B_0 \sqsubseteq \exists R^-.B, & B \sqsubseteq \exists S \end{array}$$

and the following CQ

$$\mathbf{q}(x, y'') = \exists y, z, y', x', z', x'' (R(x, y) \wedge T(y, z) \wedge T(y', z) \wedge R(x', y') \wedge S(x', z') \wedge S(x'', z') \wedge R(x'', y''))$$

shown in Fig. 10 alongside the $\exists R.A$ -generated \mathcal{T} -tree and the $\exists R^-.B$ -generated \mathcal{T} -tree. There are four tree witnesses for $\mathbf{q}(x, y'')$ and \mathcal{T} :

- $\mathfrak{t}^1 = (\mathfrak{t}_r^1, \mathfrak{t}_i^1)$ generated by $\exists T$ with $\mathfrak{t}_r^1 = \{y, y'\}$ and $\mathfrak{t}_i^1 = \{z\}$ and

$$\mathbf{q}_{\mathfrak{t}^1} = \{T(y, z), T(y', z)\};$$

- $\mathfrak{t}^2 = (\mathfrak{t}_r^2, \mathfrak{t}_i^2)$ generated by $\exists S$ with $\mathfrak{t}_r^2 = \{x', x''\}$ and $\mathfrak{t}_i^2 = \{z'\}$ and

$$\mathbf{q}_{\mathfrak{t}^2} = \{S(x', z'), S(x'', z')\};$$

- $\mathfrak{t}^3 = (\mathfrak{t}_r^3, \mathfrak{t}_i^3)$ generated by $\exists R.A$ with $\mathfrak{t}_r^3 = \{x, x'\}$ and $\mathfrak{t}_i^3 = \{y, y', z\}$ and

$$\mathbf{q}_{\mathfrak{t}^3} = \{R(x, y), T(y, z), T(y', z), R(x', y')\};$$

- $\mathfrak{t}^4 = (\mathfrak{t}_r^4, \mathfrak{t}_i^4)$ generated by $\exists R^-.B$ with $\mathfrak{t}_r^4 = \{y', y''\}$ and $\mathfrak{t}_i^4 = \{x', x'', z'\}$ and

$$\mathbf{q}_{\mathfrak{t}^4} = \{R(x', y'), S(x', z'), S(x'', z'), R(x'', y'')\}.$$

Clearly, \mathfrak{t}^3 and \mathfrak{t}^4 are conflicting; \mathfrak{t}^3 is also conflicting with \mathfrak{t}^1 but not with \mathfrak{t}^2 despite sharing a common root; symmetrically, \mathfrak{t}^4 is conflicting with \mathfrak{t}^2 . Thus, we have the following 8 independent sets of tree witnesses:

$$\emptyset, \quad \{\mathfrak{t}^1\}, \quad \{\mathfrak{t}^2\}, \quad \{\mathfrak{t}^3\}, \quad \{\mathfrak{t}^4\}, \quad \{\mathfrak{t}^1, \mathfrak{t}^2\}, \quad \{\mathfrak{t}^1, \mathfrak{t}^4\}, \quad \{\mathfrak{t}^2, \mathfrak{t}^3\},$$

which result in a tree-witness rewriting of 8 subqueries with the following tree-witness formulas:

$$\begin{aligned} \text{tw}_{\mathfrak{t}^1}(y, y') &= \exists u ((u = y) \wedge (u = y') \wedge A(u)), \\ \text{tw}_{\mathfrak{t}^2}(x', x'') &= \exists u ((u = x') \wedge (u = x'') \wedge B(u)), \\ \text{tw}_{\mathfrak{t}^3}(x, x') &= \exists u ((u = x) \wedge (u = x') \wedge A_0(u)), \\ \text{tw}_{\mathfrak{t}^4}(y', y'') &= \exists u ((u = y') \wedge (u = y'') \wedge B_0(u)). \end{aligned}$$

Theorem 36 (Kikot et al. (2012b)). *Let \mathcal{T} be a QL TBox and $\mathbf{q}(\mathbf{x})$ a CQ. For any ABox \mathcal{A} that is H-complete with respect to \mathcal{T} and any tuple \mathbf{a} in $\text{ind}(\mathcal{A})$, we have*

$$\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models \mathbf{q}(\mathbf{a}) \quad \text{if and only if} \quad \mathcal{I}_{\mathcal{A}} \models \mathbf{q}_{\text{tw}}(\mathbf{a}).$$

By Proposition 33, to obtain a PE-rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} over *arbitrary* ABoxes, it is enough to take the tree-witness rewriting $\mathbf{q}_{\text{tw}}(\mathbf{x})$ over H-complete ABoxes and replace every atom $S(\mathbf{z})$ in it with $S_{\text{ext}}(\mathbf{z})$.

The number of tree witnesses, $|\Theta_{\mathcal{T}}^{\mathbf{q}}|$, is bounded by $3^{|\mathbf{q}|}$. On the other hand, there is a sequence of queries \mathbf{q}_n and ontologies \mathcal{T}_n with exponentially many (in $|\mathbf{q}_n|$) tree witnesses (Kikot et al., 2012b). The length of \mathbf{q}_{tw} is $O(2^{|\Theta_{\mathcal{T}}^{\mathbf{q}}|} \cdot |\mathbf{q}| \cdot |\mathcal{T}|)$. It is to be noted that there exist CQs \mathbf{q} and QL TBoxes \mathcal{T} any PE-rewritings of which are of exponential size in $|\mathbf{q}|$ provided that the rewritings use the same symbols as in \mathbf{q} and \mathcal{T} (Kikot et al., 2012a). One can always reduce the size of PE-rewritings to polynomial by employing two new constants that do not occur in \mathbf{q} ; for details and further references, consult (Gottlob et al., 2014).

If any two tree-witnesses for $\mathbf{q}(\mathbf{x})$ and \mathcal{T} are *compatible*—that is, they are either non-conflicting or one is included in the other—then $\mathbf{q}_{\text{tw}}(\mathbf{x})$ can be equivalently transformed into the PE-rewriting

$$\mathbf{q}'_{\text{tw}}(\mathbf{x}) = \exists \mathbf{y} \bigwedge_{S(\mathbf{z}) \in \mathbf{q}} (S(\mathbf{z}) \vee \bigvee_{\mathfrak{t} \in \Theta_{\mathcal{T}}^{\mathbf{q}} \text{ with } S(\mathbf{z}) \in \mathfrak{t}} \text{tw}_{\mathfrak{t}})$$

of size $O(|\Theta_{\mathcal{T}}^{\mathbf{q}}| \cdot |\mathbf{q}|^2 \cdot |\mathcal{T}|)$; for details we refer the reader to (Kikot et al., 2012b).

As we saw in Example 29, CQ entailment over *RL* and *EL* knowledge bases is P-hard for data complexity, and so some CQs over such knowledge bases do not have first-order rewritings. In the next two sections, we show that one can always rewrite CQs over *RL* or *EL* TBoxes into datalog queries of polynomial size.

4.3 Datalog Rewriting for *RL*

Recall from Section 3.1 that a *datalog program* is a set of Horn clauses with one positive literal, that is, universally quantified sentences of the form

$$\gamma_0(\mathbf{x}) \leftarrow \gamma_1(\mathbf{x}) \wedge \cdots \wedge \gamma_n(\mathbf{x}),$$

where each variable of the head, $\gamma_0(\mathbf{x})$, occurs in at least one of the atoms in the body, $\gamma_1(\mathbf{x}), \dots, \gamma_n(\mathbf{x})$. (Following the datalog tradition, we omit the universal quantifiers and write the implication from right to left.) Given a datalog program Π , a set of ground atoms D and a ground atom $Q(\mathbf{a})$, we write $(\Pi, D) \models Q(\mathbf{a})$ if $Q(\mathbf{a})$ is true in every interpretation satisfying D and all clauses of Π , or, equivalently, if $Q(\mathbf{a})$ is true in the minimal (or canonical) model of (Π, D) , which is constructed in the same way as the canonical models of *RL* knowledge bases (in general, however, the arity of predicates is not bounded by 2).

We say that a CQ $\mathbf{q}(\mathbf{x})$ and a TBox \mathcal{T} are *datalog-rewritable* if there exist a datalog program Π and a predicate $Q(\mathbf{x})$ such that, for any ABox \mathcal{A} and any tuple \mathbf{a} of individuals in \mathcal{A} , we have

$$(\mathcal{T}, \mathcal{A}) \models \mathbf{q}(\mathbf{a}) \quad \text{if and only if} \quad (\Pi, \mathcal{A}) \models Q(\mathbf{a}). \quad (19)$$

(In the following two sections we view any ABox as a set of ground atoms: each concept assertion is a unary ground atom and each role assertion is a binary ground atom.) In this case, the pair $(\Pi, Q(\mathbf{x}))$ is called a *datalog rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T}* .

Let \mathcal{T} be a *positive RL* TBox, that is, an *RL* TBox without negative concept inclusions of the form $B \sqsubseteq \perp$ (see Section 3.1). Denote by $\Pi_{\mathcal{T}}$ the datalog program that contains the standard translations of all concept and role inclusions in \mathcal{T} :

$$\begin{aligned} ST_u(A) &\leftarrow ST_u(B), & \text{for each } B \sqsubseteq A \text{ in } \mathcal{T}, \\ ST_{u,v}(R_2) &\leftarrow ST_{u,v}(R_1), & \text{for each } R_1 \sqsubseteq R_2 \text{ in } \mathcal{T}. \end{aligned}$$

(Note that B can be a complex *RL* concept, constructed using \sqcap and $\exists R.C$ with possibly inverse roles.) By Theorem 15, we then obtain the following result:

Corollary 37. *For any positive RL TBox \mathcal{T} and any CQ $\mathbf{q}(\mathbf{x}) = \exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$,*

$$(\Pi_{\mathcal{T}} \cup \{Q(\mathbf{x}) \leftarrow \varphi(\mathbf{x}, \mathbf{y})\}, Q(\mathbf{x}))$$

is a datalog rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} , where Q is a fresh predicate symbol.

Note that the size of this datalog rewriting is the sum of the sizes of the TBox and the query (thus, it is linear in both).

4.4 Tree-Witness Datalog Rewriting for *EL*

Let \mathcal{T} be an *EL* TBox in normal form (Definition 19). As we saw in Section 3.2, the canonical model of $(\mathcal{T}, \mathcal{A})$ can be defined by unravelling the cycles in the

generating model $\mathcal{G}_{\mathcal{T}, \mathcal{A}}$. So, the domain of the canonical model consists of $\rightsquigarrow_{\mathcal{T}, \mathcal{A}}$ -paths of the form $aw_{\exists P_1.D_1} \cdots w_{\exists P_n.D_n}$, $n \geq 0$, such that $a \in \text{ind}(\mathcal{A})$ and, if $n > 0$, then

$$a \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists P_1.D_1} \quad \text{and} \quad w_{\exists P_i.D_i} \rightsquigarrow_{\mathcal{T}, \mathcal{A}} w_{\exists P_{i+1}.D_{i+1}}, \quad \text{for } i < n.$$

Similarly to the case of QL , we can represent the canonical model as the join of its ABox and anonymous parts. The latter consists of the trees $\mathcal{C}_{\mathcal{T}}^{\exists P.D}(a)$ defined in precisely the same way as in Section 4.2: $\mathcal{C}_{\mathcal{T}}^{\exists P.D}(a)$ is the restriction of the canonical model of $(\mathcal{T}, \{\exists P.D\}(a))$ to a and the labelled nulls with the prefix $aw_{\exists P.D}$. More formally, $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ is represented as

$$\mathcal{C}_{\mathcal{T}, \mathcal{A}} = \mathcal{I}_{\mathcal{A}}^* \oplus \bigoplus_{\substack{a \in \text{ind}(\mathcal{A}) \\ A \sqsubseteq \exists P.D \in \mathcal{T} \text{ with } \mathcal{I}_{\mathcal{A}}^* \models A(a)}} \mathcal{C}_{\mathcal{T}}^{\exists P.D}(a), \quad (20)$$

where $\mathcal{I}_{\mathcal{A}}^*$ is a model of \mathcal{A} with domain $\text{ind}(\mathcal{A})$. We again follow the divide and conquer strategy and split the process of query rewiring in two steps: rewriting over H-complete ABoxes and tree witnesses.

H-completeness. Let \mathcal{T} be an EL TBox. A simple ABox \mathcal{A} is said to be *H-complete with respect to \mathcal{T}* if, for all concept names A and role names P , we have

$$\begin{aligned} A(a) \in \mathcal{A} & \quad \text{if } \mathcal{I}_{\mathcal{A}} \models C(a) \text{ and } \mathcal{T} \models C \sqsubseteq A, \text{ for some } EL \text{ concept } C, \\ P(a, b) \in \mathcal{A} & \quad \text{if } \mathcal{I}_{\mathcal{A}} \models R(a, b) \text{ and } \mathcal{T} \models R \sqsubseteq P, \text{ for some role name } R. \end{aligned}$$

(Note that the definition is the same as for QL except that now C can be an arbitrary EL concept, not only A' or $\exists R$.) We say that $(\Pi, Q(\mathbf{x}))$, for a datalog program Π and an atom $Q(\mathbf{x})$, is a *datalog rewriting of $q(\mathbf{x})$ and \mathcal{T} over H-complete ABoxes* if (19) holds for any ABox \mathcal{A} that is H-complete with respect to \mathcal{T} and any tuple \mathbf{a} from $\text{ind}(\mathcal{A})$.

Observe that if an ABox \mathcal{A} is H-complete with respect to \mathcal{T} then the ABox part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$, that is, $\mathcal{I}_{\mathcal{A}}^*$ in (20), coincides with $\mathcal{I}_{\mathcal{A}}$. Thus, if \mathcal{T} is a flat EL TBox in normal form (which does not contain concept inclusions of the form $A \sqsubseteq \exists P.D$) then

$$(\{Q(\mathbf{x}) \leftarrow \varphi(\mathbf{x}, \mathbf{y})\}, Q(\mathbf{x}))$$

is clearly a datalog rewriting of $q(\mathbf{x}) = \exists \mathbf{y} \varphi(\mathbf{x}, \mathbf{y})$ and \mathcal{T} over H-complete ABoxes. Now, any such TBox \mathcal{T} can also be seen as a positive RL TBox, and so, by Corollary 37, the datalog program $\Pi_{\mathcal{T}}$ defined in Section 4.3 describes precisely the conditions of H-completeness of \mathcal{A} :

$$\begin{aligned} A(a) \in \mathcal{A} & \quad \text{if and only if } (\Pi_{\mathcal{T}}, \mathcal{A}) \models A(a), \\ P(a, b) \in \mathcal{A} & \quad \text{if and only if } (\Pi_{\mathcal{T}}, \mathcal{A}) \models P(a, b). \end{aligned}$$

It follows then that we can easily obtain datalog rewritings (over arbitrary ABoxes) from datalog rewritings over H-complete ABoxes:

Proposition 38. *If $(\Pi, Q(\mathbf{x}))$ is a datalog rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} over H -complete ABoxes, then $(\Pi \cup \Pi_{\mathcal{T}}, Q(\mathbf{x}))$ is a datalog rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} over arbitrary ABoxes.*

Thus, we can only concentrate on constructing datalog rewritings over H -complete ABoxes.

Tree Witnesses. Let \mathcal{T} be an EL TBox in normal form and let $\mathbf{q}(\mathbf{x})$ be a CQ. Similarly to the case of QL , to construct a rewriting of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} , we need to consider all possible subsets of \mathbf{q} , the atoms of the query, that can be homomorphically mapped to the ABox part of the canonical model $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$ and subsets of \mathbf{q} that can be homomorphically mapped to the anonymous part of $\mathcal{C}_{\mathcal{T}, \mathcal{A}}$. Tree witnesses for $\mathbf{q}(\mathbf{x})$ and \mathcal{T} are defined in absolutely the same way as in Section 4.2. A major difference, however, is that EL does not contain inverse roles, and so each \mathcal{T} -tree in the anonymous part is generated by a concept of the form $\exists P.D$, where P is a role name. It follows that if a variable, say y , of any CQ belongs to the interior of some tree witness and the CQ contains an atom of the form $P(y, y')$, then y' must also be in the interior of the same tree witness. In particular, we have the following:

Proposition 39. *Let \mathcal{T} be an EL TBox and $\mathbf{q}(\mathbf{x})$ a CQ. Then, for any binary atom $P(z_1, z_2)$ in \mathbf{q} , where P is a role name,*

- there is no tree witness for $\mathbf{q}(\mathbf{x})$ and \mathcal{T} with $P(z_1, z_2) \in \mathbf{q}_{\mathfrak{t}}$, $z_2 \in \mathfrak{t}_{\mathfrak{r}}$, $z_1 \in \mathfrak{t}_{\mathfrak{i}}$;
- there is at most one tree witness \mathfrak{t} for $\mathbf{q}(\mathbf{x})$ and \mathcal{T} such that $P(z_1, z_2) \in \mathbf{q}_{\mathfrak{t}}$, $z_1 \in \mathfrak{t}_{\mathfrak{r}}$ and $z_2 \in \mathfrak{t}_{\mathfrak{i}}$.

These observations suggest a simple algorithm for constructing all tree witnesses for any given CQ and EL TBox, which we first illustrate by a concrete example.

Example 40. Let \mathcal{T} be an EL TBox containing the following concept and role inclusions:

$$A \sqsubseteq \exists S.B, \quad B \sqsubseteq \exists R.C, \quad R \sqsubseteq T.$$

Consider the CQ

$$\mathbf{q}(x) = \exists y, z, u, v (S(x, y) \wedge T(y, z) \wedge R(u, z) \wedge T(u, v) \wedge C(v)).$$

The algorithms begins with the smallest tree witnesses. By **(tw2)**, each tree witness \mathfrak{t} is uniquely defined by its interior (a non-empty set of existentially quantified variables): the roots are all the terms that do not belong to $\mathfrak{t}_{\mathfrak{i}}$ but occur in atoms with an argument in $\mathfrak{t}_{\mathfrak{i}}$. The smallest tree witnesses are thus induced by singleton sets $\mathfrak{t}_{\mathfrak{i}}$. For $\mathfrak{t}_{\mathfrak{i}} = \{v\}$, we have $\mathfrak{t}_{\mathfrak{r}} = \{u\}$, which gives the tree witness \mathfrak{t}^1 with

$$\mathbf{q}_{\mathfrak{t}^1} = \{T(u, v), C(v)\},$$

which is generated by $\exists R.C$; see h_1 in Fig. 11. For $\mathfrak{t}_{\mathfrak{i}} = \{u\}$, we have $\mathfrak{t}_{\mathfrak{r}} = \{v, z\}$ but, since EL contains no inverse roles, there can be no homomorphism h that

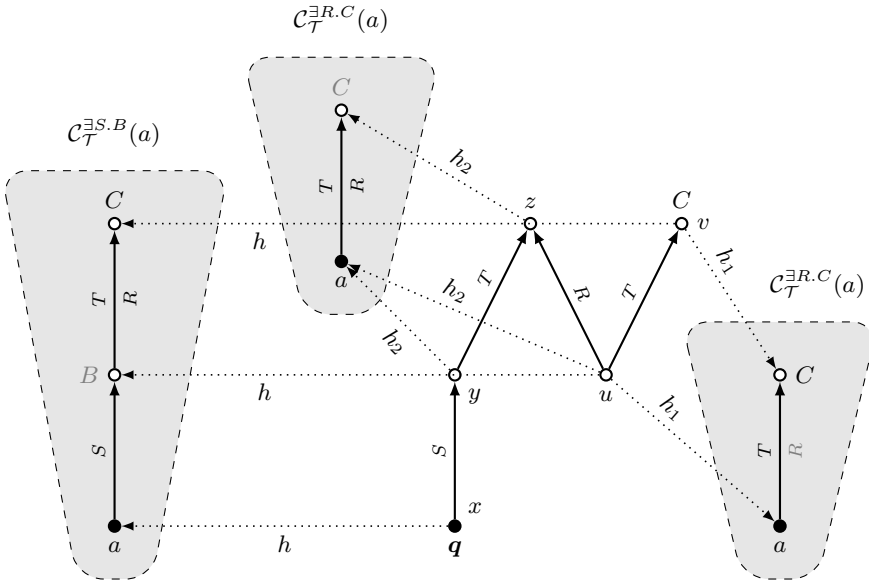


Fig. 11. Tree witnesses in Example 40

maps the set of atoms $\{T(y, z), R(u, z)\}$ into a \mathcal{T} -tree in such a way that h takes t_r to an ABox individual and t_i to a labelled null. The same argument applies to $t_i = \{y\}$. However, $t_i^2 = \{z\}$ gives rise to another tree witness, t^2 , with

$$q_{t^2} = \{T(y, z), R(u, z)\},$$

which is again generated by $\exists R.C$; see h_2 in Fig. 11 (note that both variables in $t_r^2 = \{y, u\}$ are mapped to a , the root of the $\exists R.C$ -generated \mathcal{T} -tree $C_T^{\exists R.C}(a)$). Thus, we have considered all singleton subsets of the existentially quantified variables and constructed all possible tree witnesses of *depth 1*. Next, we observe that larger tree witnesses must contain smaller tree witnesses in their interior. So, suppose that t_i contains both the roots and the interior of tree witness t^2 , that is, $\{y, z, u\} \subseteq t_i$. Then, since *EL* has no inverse roles, t_i must also contain v (and thus, the whole of tree witness t^1). In this way, we obtain a tree witness t of *depth 2* with $t_i = \{y, u, z, v\}$, $t_r = \{x\}$ and

$$q_t = \{S(x, y), T(y, z), R(u, z), T(u, v), C(v)\},$$

which is generated by $\exists S.B$; see h in Fig. 11. We have covered all possible subsets of the existentially quantified variables and in this way obtained all tree witnesses for $q(x)$ and \mathcal{T} .

A general algorithm constructing all tree witnesses for any given CQ $q(x)$ and *EL* TBox \mathcal{T} in normal form works as follows. It begins by identifying tree witnesses that have a single interior variable (these are tree witnesses of depth 1).

The algorithm then considers each set of tree witnesses $\mathfrak{t}^1, \dots, \mathfrak{t}^k$ sharing common roots (provided that all their roots are existentially quantified variables)—the roots and interiors will become the interior of a potential new, larger, tree witness. In order to satisfy **(tw2)**, the algorithm extends the set of atoms in $\mathbf{q}_0 = \mathbf{q}_{\mathfrak{t}^1} \cup \dots \cup \mathbf{q}_{\mathfrak{t}^k}$ by all atoms incident on the variables in \mathbf{q}_0 and then, to satisfy **(tw1)**, checks whether this extended set of atoms can be homomorphically mapped to a \mathcal{T} -tree. Since *EL* contains no inverse roles, any two tree witnesses \mathfrak{t}^i and \mathfrak{t}^j that share a common root must become part of the *interior* of a larger tree witness (if it exists at all) and therefore, the number of sets of tree witnesses that need consideration is bounded by the number of (existentially quantified) variables in \mathbf{q} . Therefore, this algorithm constructs all tree witnesses and runs in polynomial time in the size of \mathbf{q} and \mathcal{T} . (We note in passing that the presented algorithm resembles the construction of the equivalence relation \sim_q defined by Lutz et al. (2009) except that the equivalence relation does not take account of the TBox and the distinction between answer variables and existentially quantified variables.)

It follows that the number of tree witnesses for any $\mathbf{q}(\mathbf{x})$ and \mathcal{T} does not exceed the number of atoms in \mathbf{q} . Moreover, each pair, \mathfrak{t} and \mathfrak{t}' , of tree witnesses for $\mathbf{q}(\mathbf{x})$ and \mathcal{T} is compatible, that is,

$$\text{either } \mathbf{q}_{\mathfrak{t}} \cap \mathbf{q}_{\mathfrak{t}'} = \emptyset \text{ or } \mathbf{q}_{\mathfrak{t}} \subseteq \mathbf{q}_{\mathfrak{t}'} \text{ or } \mathbf{q}_{\mathfrak{t}} \supseteq \mathbf{q}_{\mathfrak{t}'}.$$

(In other words, the tree witnesses are partially ordered by the \subseteq relation on their sets of atoms.) It follows that we can use the ‘modified’ tree-witness rewriting $\mathbf{q}'_{\text{tw}}(\mathbf{x})$ over H-complete ABoxes defined at the end of Section 4.2, or rather its datalog representation. Let $Q_{S(\mathbf{z})}$, for each atom $S(\mathbf{z})$ in \mathbf{q} , and D be fresh k -ary predicate symbols, for $k = |\mathbf{x}| + |\mathbf{y}|$, and let D_1 be a fresh unary predicate symbol. Intuitively, the $Q_{S(\mathbf{z})}$ are the rewritings of individual atoms of the query, the interpretation of D_1 consists of individuals from the ABox that are relevant to the query and the interpretation of D of all tuples of such individuals. Formally, let $\Omega_{\mathcal{T}}^{\mathbf{q}}$ comprise the rule

$$Q(\mathbf{x}) \leftarrow \bigwedge_{S(\mathbf{z}) \in \mathbf{q}} Q_{S(\mathbf{z})}(\mathbf{x}, \mathbf{y}),$$

the following rules, for each atom $S(\mathbf{z})$ in \mathbf{q} :

$$\begin{aligned} Q_{S(\mathbf{z})}(\mathbf{x}, \mathbf{y}) &\leftarrow D(\mathbf{x}, \mathbf{y}) \wedge S(\mathbf{z}), \\ Q_{S(\mathbf{z})}(\mathbf{x}, \mathbf{y}) &\leftarrow D(\mathbf{x}, \mathbf{y}) \wedge \text{tw}_{\mathfrak{t}}(\mathfrak{t}_r), \end{aligned} \quad \text{for } \mathfrak{t} \in \Theta_{\mathcal{T}}^{\mathbf{q}} \text{ with } S(\mathbf{z}) \in \mathbf{q}_{\mathfrak{t}},$$

and the following rules defining D and D_1 :

$$\begin{aligned} D(z_1, \dots, z_k) &\leftarrow D_1(z_1) \wedge \dots \wedge D_1(z_k), \\ D_1(z) &\leftarrow A(z), && \text{for a concept name } A \text{ in } \mathbf{q} \text{ or } \mathcal{T}, \\ D_1(z) &\leftarrow P(z, y), && \text{for a role name } P \text{ in } \mathbf{q} \text{ or } \mathcal{T}, \\ D_1(z) &\leftarrow P(y, z), && \text{for a role name } P \text{ in } \mathbf{q} \text{ or } \mathcal{T}. \end{aligned}$$

(Strictly speaking, D_1 is not the same as the set of individuals in the ABox because concept and role names that occur in the ABox but do not occur in the query or the TBox are not included in the definition above: any individual that belongs only to such a concept or role is simply not visible to the query.) Thus we obtain a datalog rewriting $(\Omega_{\mathcal{T}}^q, Q(\mathbf{x}))$ of $q(\mathbf{x})$ and \mathcal{T} over H-complete ABoxes:

Theorem 41. *Let \mathcal{T} be an EL TBox and $q(\mathbf{x})$ a CQ. For any ABox \mathcal{A} that is H-complete with respect to \mathcal{T} and any tuple \mathbf{a} in $\text{ind}(\mathcal{A})$, we have*

$$\mathcal{C}_{\mathcal{T}, \mathcal{A}} \models q(\mathbf{a}) \quad \text{if and only if} \quad (\Omega_{\mathcal{T}}^q, \mathcal{I}_{\mathcal{A}}) \models Q(\mathbf{a}).$$

Finally, since the number of tree witnesses is linear in the size of q , this datalog rewriting is of polynomial size in the size of q and \mathcal{T} . We also note in passing that the datalog rewriting over H-complete ABoxes is nonrecursive (none of the predicates is defined, even indirectly, in terms of itself) and the only recursive component of the rewriting over arbitrary ABoxes is the rules ensuring H-completeness of the ABox.

5 OBDA with Ontop and Databases

In the final section of this chapter, we briefly describe the ontology-based data access (OBDA) system *Ontop*¹⁶ (Rodríguez-Muro et al., 2013) implemented at the Free University of Bozen-Bolzano and available as a plugin for the ontology editor Protégé 4, a SPARQL end-point and OWLAPI and Sesame libraries. *Ontop* is the first system to support the W3C recommendations OWL2 QL, R2RML, SPARQL and the OWL2 QL direct semantics entailment regime.

We illustrate how *Ontop* works using an example from (Rodríguez-Muro et al., 2013), which involves a (simplified) database IMDb¹⁷ (in a typical OBDA scenario data comes from a relational database rather than an ABox). From a logical point of view, a *database schema* (Abiteboul et al., 1995) contains predicate symbols (with their arity) for both stored database relations (also known as tables) and views (with their definitions in terms of stored relations) as well as a set Σ of *integrity constraints* (in the form of functional and inclusion dependencies; for example, primary and foreign keys): any instance I of the database schema must satisfy its integrity constraints Σ .

The schema of IMDb contains relations $\text{title}[m, t, y]$ with information about movies (ID, title, production year), and $\text{castinfo}[p, m, r]$ with information about movies' cast (person ID, movie ID, person role). Thus, a data instance I_{IMDb} of this schema may contain the tables

title			castinfo		
<i>m</i>	<i>t</i>	<i>y</i>	<i>p</i>	<i>m</i>	<i>r</i>
728	‘Django Unchained’	2012	n37	728	1
			n38	728	1

¹⁶ ontop.inf.unibz.it

¹⁷ www.imdb.com/interfaces

The integrity constraints Σ_{IMDb} of IMDb include the following foreign key (an inclusion dependency):

$$\forall m (\exists p, r \text{ castinfo}(p, m, r) \rightarrow \exists t, y \text{ title}(m, t, y)) \quad (21)$$

(‘each tuple in `castinfo` must refer to an existing title’) and the following primary key (a functional dependency with determinant m):

$$\forall m \forall t_1 \forall t_2 (\exists y \text{ title}(m, t_1, y) \wedge \exists y \text{ title}(m, t_2, y) \rightarrow (t_1 = t_2)), \quad (22)$$

$$\forall m \forall y_1 \forall y_2 (\exists t \text{ title}(m, t, y_1) \wedge \exists t \text{ title}(m, t, y_2) \rightarrow (y_1 = y_2)) \quad (23)$$

(‘each title is uniquely determined by its ID m ’).

In the framework of OBDA, the users are not supposed to know the structure of the database. Instead, they are given an ontology, e.g., MO^{18} , which describes the application domain in terms of concepts and roles. In our example we have concepts `mo:Movie` and `mo:Person`, and roles `mo:cast`, `mo:title` and `mo:year` related by OWL 2 QL TBox \mathcal{T}_{MO} containing, in particular, the following inclusions:

$$\begin{array}{ll} \text{mo:Movie} \sqsubseteq \exists \text{mo:title}, & \text{mo:Movie} \sqsubseteq \exists \text{mo:cast}, \\ \text{mo:Movie} \sqsubseteq \exists \text{mo:year}, & \exists \text{mo:title} \sqsubseteq \text{mo:Movie}, \\ \exists \text{mo:cast} \sqsubseteq \text{mo:Movie}, & \exists \text{mo:cast}^- \sqsubseteq \text{mo:Person}. \end{array}$$

The vocabularies of the database schema and the given OWL 2 QL TBox are linked together by means of mappings produced by a domain expert or extracted (semi)automatically. There are different known types of mappings: LAV (local-as-views), GAV (global-as-views), GLAV, etc.; consult, e.g., (Lenzerini, 2002) for an overview. Here we concentrate on GAV mappings because they guarantee low complexity of query answering (in what follows we call them simply mappings)—*Ontop* uses R2RML¹⁹ to specify them. A *mapping*, \mathcal{M} , is a set of rules of the form

$$S(\mathbf{x}) \leftarrow \varphi(\mathbf{x}, \mathbf{z}),$$

where S is a concept or role name in the ontology and $\varphi(\mathbf{x}, \mathbf{z})$ is a conjunction of atoms with database relations (both stored and views) and a *filter*, that is, a Boolean combination of built-in predicates such as $=$ and $<$. (Note that, by including views in the schema, we can express any SQL query in mappings, which is important from the practical point of view.) In our example, a mapping \mathcal{M}_{MO} that relates the terms of MO to the schema of IMDb contains the following rules:

$$\text{mo:Movie}(m) \leftarrow \text{title}(m, t, y), \quad (24)$$

$$\text{mo:title}(m, t) \leftarrow \text{title}(m, t, y), \quad (25)$$

$$\text{mo:year}(m, y) \leftarrow \text{title}(m, t, y), \quad (26)$$

$$\text{mo:cast}(m, p) \leftarrow \text{castinfo}(p, m, r), \quad (27)$$

$$\text{mo:Person}(p) \leftarrow \text{castinfo}(p, m, r). \quad (28)$$

¹⁸ www.movieontology.org

¹⁹ www.w3.org/TR/r2rml

Given a mapping \mathcal{M} from a database schema to an OWL2QL TBox \mathcal{T} and an instance \mathbf{I} of this schema, the ground atoms

$$S(\mathbf{a}), \quad \text{for } S(\mathbf{x}) \leftarrow \varphi(\mathbf{x}, \mathbf{z}) \text{ in } \mathcal{M} \text{ and } \mathbf{I} \models \exists \mathbf{z} \varphi(\mathbf{a}, \mathbf{z}),$$

comprise the ABox denoted by $\mathcal{A}_{\mathbf{I}, \mathcal{M}}$ and called the *virtual ABox* for \mathcal{M} over \mathbf{I} (Rodríguez-Muro and Calvanese, 2011). This ABox is just a convenient presentational tool and does not have to be materialised by the system. Then the virtual ABox for \mathcal{M}_{MO} over \mathbf{I}_{IMDB} consists of the ground atoms

$$\begin{aligned} &\text{mo:Movie}(728), \text{ mo:title}(728, \text{'Django Unchained'}), \text{ mo:year}(728, 2012), \\ &\text{mo:Person}(n37), \text{ mo:cast}(728, n37), \\ &\text{mo:Person}(n38), \text{ mo:cast}(728, n38). \end{aligned}$$

The *certain answers* to a CQ $\mathbf{q}(\mathbf{x})$ over \mathcal{T} linked by \mathcal{M} to \mathbf{I} are defined as the certain answers to $\mathbf{q}(\mathbf{x})$ over $(\mathcal{T}, \mathcal{A}_{\mathbf{I}, \mathcal{M}})$. In order to find certain answers, one could first construct a PE-rewriting $\mathbf{q}'(\mathbf{x})$ of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} over *arbitrary* ABoxes (in the sequel, it will be convenient to represent such a rewriting as a non-recursive datalog program). The rewriting $\mathbf{q}'(\mathbf{x})$ could then be *unfolded* into an SQL query using the so-called *partial evaluation* (Lloyd and Shepherdson, 1991), which exhaustively applies SLD-resolution to $\mathbf{q}'(\mathbf{x})$ and the mapping \mathcal{M} and returns those rules whose bodies contain only database atoms. Consider, for example, CQ $\mathbf{q}(x) = \text{mo:Movie}(x)$. An obvious rewriting of $\mathbf{q}(x)$ and the TBox \mathcal{T}_{MO} (over arbitrary ABoxes) contains the following three rules:

$$\mathbf{q}'(x) \leftarrow \text{mo:Movie}(x), \quad (29)$$

$$\mathbf{q}'(x) \leftarrow \text{mo:title}(x, y), \quad (30)$$

$$\mathbf{q}'(x) \leftarrow \text{mo:cast}(x, y). \quad (31)$$

The unfolding applies the SLD-resolution procedure to these three rules and the mapping \mathcal{M}_{MO} and produces the rules

$$\mathbf{q}'(x) \leftarrow \text{title}(x, t, y), \quad (29+24)$$

$$\mathbf{q}'(x) \leftarrow \text{title}(x, t, y), \quad (30+25)$$

$$\mathbf{q}'(x) \leftarrow \text{castinfo}(p, x, r). \quad (31+27)$$

The resulting union of SELECT-PROJECT-JOIN queries could then be forwarded for execution to a relational database management system (RDBMS).

The same result can be achieved by using the tree-witness rewriting $\mathbf{q}_{\text{tw}}(\mathbf{x})$ of $\mathbf{q}(\mathbf{x})$ and \mathcal{T} over H-complete ABoxes introduced in Section 4.2. An obvious way to construct H-complete ABoxes is to take the *composition* of \mathcal{M} and the inclusions in \mathcal{T} , that is, a mapping $\mathcal{M}^{\mathcal{T}}$ given by

$$\begin{aligned} A(x) \leftarrow \varphi(x, \mathbf{z}), & \quad \text{if } A'(x) \leftarrow \varphi(x, \mathbf{z}) \in \mathcal{M} \text{ and } \mathcal{T} \models A' \sqsubseteq A, \\ A(x) \leftarrow \varphi(x, y, \mathbf{z}), & \quad \text{if } R(x, y) \leftarrow \varphi(x, y, \mathbf{z}) \in \mathcal{M} \text{ and } \mathcal{T} \models \exists R \sqsubseteq A, \\ P(x, y) \leftarrow \varphi(x, y, \mathbf{z}), & \quad \text{if } R(x, y) \leftarrow \varphi(x, y, \mathbf{z}) \in \mathcal{M} \text{ and } \mathcal{T} \models R \sqsubseteq P. \end{aligned}$$

(Recall that we do not distinguish between $P^-(y, x)$ and $P(x, y)$.) Thus, for any I and any tuple \mathbf{a} of individuals in $\mathcal{A}_{I, \mathcal{M}}$, we have:

$$(\mathcal{T}, \mathcal{A}_{I, \mathcal{M}}) \models \mathbf{q}(\mathbf{a}) \quad \text{if and only if} \quad \mathcal{A}_{I, \mathcal{M}\tau} \models \mathbf{q}_{\text{tw}}(\mathbf{a}). \quad (32)$$

So, to compute the answers to $\mathbf{q}(\mathbf{x})$ over \mathcal{T} linked by \mathcal{M} to I , one can unfold the tree-witness rewriting $\mathbf{q}_{\text{tw}}(\mathbf{x})$ over H-complete ABoxes with the help of the composition $\mathcal{M}^{\mathcal{T}}$. However, the resulting query will produce duplicating answers if the ontology axioms express the same properties of the application domain as the integrity constraints of the database. For example, the IMDb schema Σ_{IMDb} contains foreign key (21): movie ID in `castinfo` references movie ID in `title`, and therefore the unfolded rewriting above will return the same movie many times—once from `title` and once for each of the cast members of the movie in `castinfo`. Such a duplication is clearly an undesirable feature of this straightforward approach.

For this reason, before applying $\mathcal{M}^{\mathcal{T}}$ to unfold the tree-witness rewriting, *Ontop* optimises the mapping using the database integrity constraints Σ . This allows *Ontop* to reduce redundancy in answers and substantially shorten the SQL queries, which makes the OBDA system more efficient.

5.1 \mathcal{T} -mappings

We say that a mapping \mathcal{M} is a *\mathcal{T} -mapping over dependencies Σ* if the ABox $\mathcal{A}_{I, \mathcal{M}}$ is H-complete with respect to \mathcal{T} , for any data instance I satisfying Σ . The composition $\mathcal{M}^{\mathcal{T}}$ defined above is trivially a \mathcal{T} -mapping over any Σ . *Ontop* starts with $\mathcal{M}^{\mathcal{T}}$ and then applies a series of optimisations to construct a simpler \mathcal{T} -mapping.

Inclusion Dependencies. Suppose $\mathcal{M} \cup \{S(\mathbf{x}) \leftarrow \psi_1(\mathbf{x}, \mathbf{z})\}$ is a \mathcal{T} -mapping over Σ . If there is a more specific rule than $S(\mathbf{x}) \leftarrow \psi_1(\mathbf{x}, \mathbf{z})$ in \mathcal{M} , then \mathcal{M} itself will also be a \mathcal{T} -mapping. To discover such ‘more specific’ rules, we run the standard query containment check (Abiteboul et al., 1995) taking account of the inclusion dependencies. For example, since $\mathcal{T}_{\text{MO}} \models \exists \text{mo:cast} \sqsubseteq \text{mo:Movie}$ in our running example, the composition of mapping \mathcal{M}_{MO} and ontology \mathcal{T}_{MO} contains the following rules for `mo:Movie`:

$$\begin{aligned} \text{mo:Movie}(m) &\leftarrow \text{title}(m, t, y), \\ \text{mo:Movie}(m) &\leftarrow \text{castinfo}(p, m, r). \end{aligned}$$

As we observed above, the latter rule is redundant because Σ_{IMDb} contains inclusion dependency (21), which is repeated here for reference:

$$\forall m (\exists p, r \text{ castinfo}(p, m, r) \rightarrow \exists t, y \text{ title}(m, t, y)).$$

Disjunctions in SQL. Another way to reduce the size of a \mathcal{T} -mapping is to identify pairs of rules whose bodies are equivalent up to *filters with respect to*

concept	index	interval
mo:Actor	1	[1,1]
mo:Artist	2	[1,2]
mo:Director	3	[3,3]
mo:Person	4	[1,4]

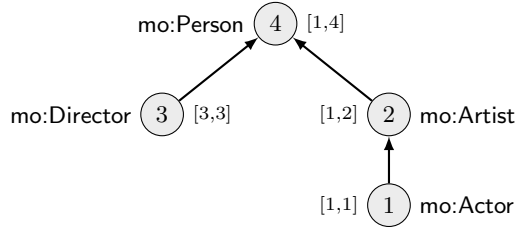


Fig. 12. Semantic Index example

constant values. This optimisation deals with the rules introduced due to the so-called type (discriminating) attributes (Elmasri and Navathe, 2010) in database schemas. For example, the mapping \mathcal{M}_{MO} contains six rules for sub-concepts of `mo:Person`:

$$\begin{aligned}
 \text{mo:Actor}(p) &\leftarrow \text{castinfo}(c, p, m, r), (r = 1), \\
 \text{mo:Actress}(p) &\leftarrow \text{castinfo}(c, p, m, r), (v = 2), \\
 &\dots \\
 \text{mo:Editor}(p) &\leftarrow \text{castinfo}(c, p, m, r), (r = 6).
 \end{aligned}$$

Thus, the composition of \mathcal{M}_{MO} and \mathcal{T}_{MO} contains six rules for `mo:Person` that differ only in the last condition ($r = k$), $1 \leq k \leq 6$. These can be reduced to a single rule:

$$\text{mo:Person}(p) \leftarrow \text{castinfo}(c, p, m, r), ((r = 1) \vee \dots \vee (r = 6)).$$

Note that such disjunctions lend themselves to efficient evaluation by RDBMSs.

Materialised ABoxes and Semantic Index. In addition to working with proper relational data sources, *Ontop* also supports ABox storage in the form of structureless *universal tables*: a binary relation $\text{CA}[id, \text{concept-id}]$ and a ternary relation $\text{RA}[id_1, id_2, \text{role-id}]$ represent concept and role assertions. The universal tables give rise to trivial mappings, and *Ontop* implements a technique, the *semantic index* (Rodríguez-Muro and Calvanese, 2011), that takes advantage of SQL features in \mathcal{T} -mappings for this scenario. The key observation is that, since the IDs in the universal tables CA and RA can be chosen by the system, each concept and role in the TBox \mathcal{T} can be assigned a numeric *index* and a set of numerical *intervals* in such a way that the resulting \mathcal{T} -mapping contains simple SQL queries with interval filter conditions. For example, in \mathcal{T}_{MO} , we have

$$\begin{aligned}
 \text{mo:Actor} &\sqsubseteq \text{mo:Artist}, \\
 \text{mo:Artist} &\sqsubseteq \text{mo:Person}, \\
 \text{mo:Director} &\sqsubseteq \text{mo:Person},
 \end{aligned}$$

so we can choose indexes and intervals for these concepts as in Fig. 12. It can be seen that these intervals respect the concept inclusions of the TBox: for instance,

[1,1] for `mo:Actor` is a subset of [1,2] for `mo:Artist`. This will generate a \mathcal{T} -mapping with

$$\begin{aligned} \text{mo:Actor}(p) &\leftarrow \text{CA}(p, \text{concept-id}), \quad (\text{concept-id} = 1), \\ \text{mo:Artist}(p) &\leftarrow \text{CA}(p, \text{concept-id}), \quad (1 \leq \text{concept-id} \leq 2), \\ \text{mo:Director}(p) &\leftarrow \text{CA}(p, \text{concept-id}), \quad (\text{concept-id} = 3), \\ \text{mo:Person}(p) &\leftarrow \text{CA}(p, \text{concept-id}), \quad (1 \leq \text{concept-id} \leq 4). \end{aligned}$$

Thus, by choosing appropriate concept and role IDs, we effectively construct H-complete ABoxes *without* the expensive forward chaining procedure (and the need to store large amounts of the derived assertions and the extra complications with updating the data). On the other hand, the semantic index \mathcal{T} -mappings are based on range expressions that can be evaluated efficiently by RDBMSs using standard B-tree indexes (Elmasri and Navathe, 2010).

5.2 Unfolding with Semantic Query Optimisation (SQO)

Ontop applies the Semantic Query Optimisation (Chakravarthy et al., 1986) to rules obtained at the intermediate steps of unfolding. In particular, this eliminates redundant JOIN operations caused by reification of database relations by means of concepts and roles. Consider, for example, the CQ

$$\mathbf{q}(t, y) \leftarrow \text{mo:Movie}(m), \text{ mo:title}(m, t), \text{ mo:year}(m, y), (y > 2010).$$

It has no tree witnesses, and so $\mathbf{q}_{\text{tw}}(t, y) = \mathbf{q}(t, y)$. By straightforwardly applying the unfolding to $\mathbf{q}_{\text{tw}}(t, y)$ and the \mathcal{T} -mapping \mathcal{M}_{MO} , we obtain the query

$$\mathbf{q}'_{\text{tw}}(t, y) \leftarrow \text{title}(m, t_0, y_0), \text{ title}(m, t, y_1), \text{ title}(m, t_2, y), (y > 2010),$$

which requires two (potentially) expensive JOIN operations. However, by using functional dependencies (22) and (23) for the primary key m of `title`, which are repeated below:

$$\begin{aligned} \forall m \forall t_1 \forall t_2 (\exists y \text{ title}(m, t_1, y) \wedge \exists y \text{ title}(m, t_2, y) \rightarrow (t_1 = t_2)), \\ \forall m \forall y_1 \forall y_2 (\exists t \text{ title}(m, t, y_1) \wedge \exists t \text{ title}(m, t, y_2) \rightarrow (y_1 = y_2)), \end{aligned}$$

we can reduce two JOIN operations in the first three atoms of $\mathbf{q}'_{\text{tw}}(t, y)$ to a single atom $\text{title}(m, t, y)$:

$$\mathbf{q}''_{\text{tw}}(t, y) \leftarrow \text{title}(m, t, y), (y > 2010).$$

Note that these two JOIN operations were introduced to reconstruct the ternary relation from its reification by means of the roles `mo:title` and `mo:year`.

The role of SQO in OBDA systems appears to be much more prominent than in conventional RDBMSs, where it was initially proposed to optimise SQL queries. While some of SQO techniques reached industrial RDBMSs, it never had a strong impact on the database community because it is costly compared to statistics- and heuristics-based methods, and because most SQL queries are written by highly-skilled experts (and so are nearly optimal anyway). In OBDA scenarios, in contrast, SQL queries are generated automatically, and so SQO becomes the only tool to avoid redundant and expensive JOIN operations.

References

- Abiteboul, S., Hull, R., Vianu, V.: *Foundations of Databases*. Addison-Wesley (1995)
- Arora, S., Barak, B.: *Computational Complexity: A Modern Approach*, 1st edn. Cambridge University Press, New York (2009)
- Artale, A., Calvanese, D., Kontchakov, R., Zakharyashev, M.: The DL-Lite family and relations. *Journal of Artificial Intelligence Research (JAIR)* 36, 1–69 (2009)
- Baader, F., Brandt, S., Lutz, C.: Pushing the EL envelope. In: Kaelbling, L.P., Saffiotti, A. (eds.) *Proceedings of the 19th Int. Joint Conf. on Artificial Intelligence, IJCAI-2005*, pp. 364–369. Professional Book Center (2005)
- Baader, F., Brandt, S., Lutz, C.: Pushing the EL envelope further. In: Clark, K., Patel-Schneider, P.F. (eds.) *Proceedings of the OWLED 2008 DC Workshop on OWL: Experiences and Directions* (2008)
- Baget, J.-F., Leclère, M., Mugnier, M.-L., Salvat, E.: On rules with existential variables: Walking the decidability line. *Artificial Intelligence* 175(9–10), 1620–1654 (2011)
- Brandt, S.: Polynomial time reasoning in a description logic with existential restrictions, GCI axioms, and—what else? In: *Proc. of the 16th European Conf. on Artificial Intelligence, ECAI-2004*, pp. 298–302. IOS Press (2004)
- Calì, A., Gottlob, G., Kifer, M.: Taming the infinite chase: Query answering under expressive relational constraints. *J. of Artificial Intelligence Research* 48, 115–174 (2013)
- Calì, A., Gottlob, G., Lukasiewicz, T.: A general datalog-based framework for tractable query answering over ontologies. *J. of Web Semantics* 14, 57–83 (2012)
- Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Rosati, R.: DL-Lite: Tractable description logics for ontologies. In: *Proc. of AAAI*, pp. 602–607. AAAI Press / The MIT Press (2005)
- Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Rosati, R.: 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*, pp. 260–270 (2006)
- Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Rosati, R.: Tractable reasoning and efficient query answering in description logics: The DL-Lite family. *J. of Automated Reasoning* 39(3), 385–429 (2007)
- Ceri, S., Gottlob, G., Tanca, L.: What you always wanted to know about datalog (and never dared to ask). *IEEE Trans. Knowl. Data Eng.* 1(1), 146–166 (1989)
- Chagrov, A., Zakharyashev, M.: *Modal Logic*. Oxford Logic Guides, vol. 35. Clarendon Press, Oxford (1997)
- Chakravarthy, U.S., Fishman, D.H., Minker, J.: *Semantic query optimization in expert systems and database systems*. Benjamin-Cummings Publishing Co., Inc. (1986)
- Chortaras, A., Trivela, D., Stamou, G.: Optimized query rewriting for OWL 2 QL. In: Bjørner, N., Sofronie-Stokkermans, V. (eds.) *CADE 2011*. LNCS, vol. 6803, pp. 192–206. Springer, Heidelberg (2011)
- Eiter, T., Lutz, C., Ortiz, M., Šimkus, M.: Query answering in description logics: The knots approach. In: Ono, H., Kanazawa, M., de Queiroz, R. (eds.) *WoLLIC 2009*. LNCS, vol. 5514, pp. 26–36. Springer, Heidelberg (2009)
- Eiter, T., Ortiz, M., Šimkus, M., Tran, T.-K., Xiao, G.: Query rewriting for Horn-SHIQ plus rules. In: *Proc. of AAAI 2012*. AAAI Press (2012)
- Elmasri, R., Navathe, S.: *Fundamentals of Database Systems*, 6th edn. Addison-Wesley (2010)
- Garey, M., Johnson, D.: *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman (1979)

- Gottlob, G., Kikot, S., Kontchakov, R., Podolskii, V.V., Schwentick, T., Zakharyashev, M.: The price of query rewriting in ontology-based data access. *Artificial Intelligence* 213, 42–59 (2014)
- Gottlob, G., Orsi, G., Pieris, A.: Ontological queries: Rewriting and optimization. In: *Proc. of ICDE 2011*, pp. 2–13. IEEE Computer Society Press (2011)
- Grosz, B.N., Horrocks, I., Volz, R., Decker, S.: Description logic programs: Combining logic programs with description logic. In: *Proc. of the 12th Int. World Wide Web Conference, WWW 2003*, pp. 48–57 (2003)
- Horrocks, I., Kutz, O., Sattler, U.: The even more irresistible SROIQ. In: *Proc. of the 10th Int. Conf. on Principles of Knowledge Representation and Reasoning, KR 2006*, pp. 57–67. AAAI Press (2006)
- Horrocks, I., Motik, B., Wang, Z.: The HermiT OWL reasoner. In: *Proc. of ORE, CEUR Workshop Proceedings*, vol. 858, CEUR-WS.org (2012)
- Johnson, D.S., Klug, A.C.: Testing containment of conjunctive queries under functional and inclusion dependencies. *J. Comput. Syst. Sci.* 28(1), 167–189 (1984)
- Kikot, S., Kontchakov, R., Podolskii, V., Zakharyashev, M.: Exponential lower bounds and separation for query rewriting. In: Czumaj, A., Mehlhorn, K., Pitts, A., Wattenhofer, R. (eds.) *ICALP 2012, Part II. LNCS*, vol. 7392, pp. 263–274. Springer, Heidelberg (2012a)
- Kikot, S., Kontchakov, R., Zakharyashev, M.: Conjunctive query answering with OWL 2 QL. In: *Proc. of KR 2012*, AAAI Press (2012b)
- König, M., Leclère, M., Mugnier, M.-L., Thomazo, M.: A sound and complete backward chaining algorithm for existential rules. In: Krötzsch, M., Straccia, U. (eds.) *RR 2012. LNCS*, vol. 7497, pp. 122–138. Springer, Heidelberg (2012)
- Kontchakov, R., Lutz, C., Toman, D., Wolter, F., Zakharyashev, M.: The combined approach to ontology-based data access. In: *Proceedings of the 20th Int. Joint Conf. on Artificial Intelligence, IJCAI-2011*, pp. 2656–2661. AAAI Press (2011)
- Kozen, D.: *Theory of Computation*. Springer (2006)
- Lenzerini, M.: Data integration: A theoretical perspective. In: *Proc. of the 21st ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS 2002*, pp. 233–246. ACM (2002)
- Libkin, L.: *Elements Of Finite Model Theory*. Springer (2004)
- Lloyd, J., Shepherdson, J.: Partial Evaluation in Logic Programming. *The Journal of Logic Programming* 11(3-4), 217–242 (1991)
- Lutz, C., Seylan, İ., Toman, D., Wolter, F.: The combined approach to OBDA: Taming role hierarchies using filters. In: Alani, H., et al. (eds.) *ISWC 2013, Part I. LNCS*, vol. 8218, pp. 314–330. Springer, Heidelberg (2013)
- Lutz, C., Toman, D., Wolter, F.: Conjunctive query answering in the description logic EL using a relational database system. In: *Proceedings of the 21st Int. Joint Conf. on Artificial Intelligence, IJCAI 2009*, pp. 2070–2075 (2009)
- Lutz, C., Wolter, F.: Conservative extensions of the lightweight description logic EL. In: Pfenning, F. (ed.) *CADE 2007. LNCS (LNAI)*, vol. 4603, pp. 84–99. Springer, Heidelberg (2007)
- Motik, B.: On the properties of metamodeling in OWL. *J. Log. Comput.* 17(4), 617–637 (2007)
- Papadimitriou, C.: *Computational Complexity*. Addison-Wesley (1994)
- Pérez-Urbina, H., Motik, B., Horrocks, I.: A comparison of query rewriting techniques for DL-Lite. In: *Proc. of DL 2009*, vol. 477, CEUR-WS (2009)
- Pérez-Urbina, H., Rodríguez-Díaz, E., Grove, M., Konstantinidis, G., Sirin, E.: Evaluation of query rewriting approaches for OWL 2. In: *Proc. of SSWS+HPCSW 2012*, vol. 943, CEUR-WS (2012)

- Poggi, A., Lembo, D., Calvanese, D., De Giacomo, G., Lenzerini, M., Rosati, R.: Linking data to ontologies. *J. on Data Semantics*, X:133–X:173 (2008)
- Rodríguez-Muro, M., Calvanese, D.: Dependencies: Making ontology based data access work. In: *Proc. of AMW 2011*, vol. 749, CEUR-WS.org (2011)
- Rodríguez-Muro, M., Kontchakov, R., Zakharyashev, M.: Ontology-based data access: Ontop of databases. In: Alani, H., et al. (eds.) *ISWC 2013, Part I*. LNCS, vol. 8218, pp. 558–573. Springer, Heidelberg (2013)
- Rosati, R.: Prexto: Query rewriting under extensional constraints in DL-Lite. In: Simperl, E., Cimiano, P., Polleres, A., Corcho, O., Presutti, V. (eds.) *ESWC 2012*. LNCS, vol. 7295, pp. 360–374. Springer, Heidelberg (2012)
- Rosati, R., Almatelli, A.: Improving query answering over DL-Lite ontologies. In: *Proc. of KR 2010*. AAAI Press (2010)
- Schaerf, A.: On the complexity of the instance checking problem in concept languages with existential quantification. *J. of Intelligent Information Systems* 2, 265–278 (1993)
- Sirin, E., Parsia, B., Cuenca Grau, B., Kalyanpur, A., Katz, Y.: Pellet: A practical OWL-DL reasoner. *J. of Web Semantics* 5(2), 51–53 (2007)
- Steigmiller, A., Liebig, T., Glimm, B.: Konclude: System description. *J. of Web Semantics* (2014)
- ter Horst, H.J.: Completeness, decidability and complexity of entailment for RDF Schema and a semantic extension involving the OWL vocabulary. *J. of Web Semantics* 3(2-3), 79–115 (2005)
- Tobies, S.: Complexity results and practical algorithms for logics in knowledge representation. PhD thesis, LuFG Theoretical Computer Science, RWTH-Aachen, Germany (2001)
- Tsarkov, D., Horrocks, I.: FaCT++ description logic reasoner: System description. In: Furbach, U., Shankar, N. (eds.) *IJCAR 2006*. LNCS (LNAI), vol. 4130, pp. 292–297. Springer, Heidelberg (2006)
- Vardi, M.: The complexity of relational query languages (extended abstract). In: *Proc. of the 14th ACM SIGACT Symp. on Theory of Computing, STOC 1982*, pp. 137–146 (1982)