

# Beyond OWL 2 QL in OBDA: Rewritings and Approximations

Elena Botoeva,<sup>1</sup> Diego Calvanese,<sup>1</sup> Valerio Santarelli,<sup>2</sup> Domenico F. Savo,<sup>2</sup>  
Alessandro Solimando,<sup>3</sup> and Guohui Xiao<sup>1</sup>

<sup>1</sup> KRDB Research Centre for Knowledge and Data, Free University of Bozen-Bolzano, Italy, lastname@inf.unibz.it

<sup>2</sup> Dip. di Ing. Informatica Automatica e Gestionale, Sapienza Università di Roma, Italy, lastname@dis.uniroma1.it

<sup>3</sup> DIBRIS, University of Genova, Italy, alessandro.solimando@unige.it

## Abstract

Ontology-based data access (OBDA) is a novel paradigm facilitating access to relational data, realized by linking data sources to an ontology by means of declarative mappings. *DL-Lite<sub>R</sub>*, which is the logic underpinning the W3C ontology language OWL 2 QL and the current language of choice for OBDA, has been designed with the goal of delegating query answering to the underlying database engine, and thus is restricted in expressive power. E.g., it does not allow one to express disjunctive information, and any form of recursion on the data. The aim of this paper is to overcome these limitations of *DL-Lite<sub>R</sub>*, and extend OBDA to more expressive ontology languages, while still leveraging the underlying relational technology for query answering. We achieve this by relying on two well-known mechanisms, namely conservative rewriting and approximation, but significantly extend their practical impact by bringing into the picture the mapping, an essential component of OBDA. Specifically, we develop techniques to rewrite OBDA specifications with an expressive ontology to “equivalent” ones with a *DL-Lite<sub>R</sub>* ontology, if possible, and to approximate them otherwise. We do so by exploiting the high expressive power of the mapping layer to capture part of the domain semantics of rich ontology languages. We have implemented our techniques in the prototype system ONTOPROX, making use of the state-of-the-art OBDA system ONTOP and the query answering system CLIPPER, and we have shown their feasibility and effectiveness with experiments on synthetic and real-world data.

## 1 Introduction

Ontology-Based Data Access (OBDA) is a popular paradigm that enables end users to access data sources through an ontology, abstracting away low-level details of the data sources themselves. The ontology provides a high-level description of the domain of interest, and is semantically linked to the data sources by means of a set of mapping assertions (Calvanese et al. 2009; Giese et al. 2015). Typically, the data sources are represented as relational data, the ontology is constituted by a set of logical axioms over concepts and roles, and each mapping assertion relates an SQL query over the database to a concept or role of the ontology.

As an example, consider a bank domain, where we can specify that a checking account in the name of a person is a

simple account by means of the axiom (expressed in description logic notation)  $\text{CAcc} \sqcap \exists \text{inNameOf}.\text{Person} \sqsubseteq \text{SAcc}$ . We assume that the information about the accounts and their owners is stored in a database  $\mathcal{D}$ , and that the ontology terms  $\text{CAcc}$ ,  $\text{inNameOf}$ , and  $\text{Person}$  are connected to  $\mathcal{D}$  respectively via the mapping assertions  $\text{sql}_1(x) \rightsquigarrow \text{CAcc}(x)$ ,  $\text{sql}_2(x, y) \rightsquigarrow \text{inNameOf}(x, y)$  and  $\text{sql}_3(x) \rightsquigarrow \text{Person}(x)$ , where each  $\text{sql}_i$  is a (possibly very complex) SQL query over  $\mathcal{D}$ . Suppose now that the user intends to extract all simple accounts from  $\mathcal{D}$ . Formulating such a query directly over  $\mathcal{D}$  would require to know precisely how  $\mathcal{D}$  is structured, and thus could be complicated. Instead, exploiting OBDA, the user can simply query the ontology with  $q(x) = \text{SAcc}(x)$ , and rely on the OBDA system to get the answers.

Making OBDA work efficiently over large amounts of data, requires that query answering over the ontology is *first-order (FO)-rewritable*<sup>1</sup> (Calvanese et al. 2007; Artale et al. 2009), which in turn limits the expressiveness of the ontology language, and the degree of detail with which the domain of interest can be captured. The current language of choice for OBDA is *DL-Lite<sub>R</sub>*, the logic underlying OWL 2 QL (Motik et al. 2009), which has been specifically designed to ensure FO-rewritability of query answering. Hence, it does not allow one to express disjunctive information, or any form of recursion on the data (e.g., as resulting from qualified existentials on the left-hand side of concept inclusions), since using such constructs in general causes the loss of FO-rewritability (Calvanese et al. 2013). For this reason, in many situations the expressive power of *DL-Lite<sub>R</sub>* is too restricted to capture real-world scenarios; e.g., the axiom in our example is not expressible in *DL-Lite<sub>R</sub>*.

The aim of this work is to overcome these limitations of *DL-Lite<sub>R</sub>* by allowing the use of additional constructs in the ontology. To be able to exploit the added value coming from OBDA in real-world settings, an important requirement is the efficiency of query answering, achieved through a rewriting-based approach. This is only possible for ontology languages that are FO-rewritable. Two general mechanisms that have been proposed to cope with computational complexity coming from high expressiveness of ontology languages, and that allow one to regain FO-rewritability, are conservative rewriting (Lutz, Piro, and

<sup>1</sup>Recall that FO queries constitute the core of SQL.

Wolter 2011) and approximation (Ren, Pan, and Zhao 2010; Console et al. 2014). Given an ontology in a powerful language, in the former approach it is rewritten, when possible, into an equivalent one in a restricted language, while in the latter it is approximated, thus losing part of its semantics.

In this work, we significantly extend the practical impact of both approaches by bringing into the picture *the mapping*, an essential component of OBDA that has been ignored so far. Indeed, it is a fairly expressive component of an OBDA system, since it allows one to make use of arbitrary SQL (hence FO) queries to relate the content of the data source to the elements of the ontology. Hence, a natural question is how one can use the mapping component to capture as much as possible additional domain semantics, resulting in better approximations or more cases where conservative rewritings are possible, while maintaining a  $DL\text{-}Lite_{\mathcal{R}}$  ontology.

We illustrate how this can be done on our running example, where the non- $DL\text{-}Lite_{\mathcal{R}}$  axiom can be encoded by adding the assertion  $sql_1(x) \bowtie sql_2(x, y) \bowtie sql_3(y) \rightsquigarrow SAcc(x)$  to the mapping. This assertion connects  $\mathcal{D}$  directly to the ontology term  $SAcc$  by making use of a join of the SQL queries in the original mapping. We observe that the resulting mapping, together with the ontology in which the non- $DL\text{-}Lite_{\mathcal{R}}$  axiom has been removed, constitutes a conservative rewriting of the original OBDA specification.

In this paper, we elaborate on this idea, by introducing a novel *framework for rewriting and approximation of OBDA specifications*. Specifically, we provide a notion of *rewriting based on query inseparability* of OBDA specifications (Bivenvenu and Rosati 2015). To deal with those cases where it is not possible to rewrite the OBDA specification into a query inseparable one whose ontology is in  $DL\text{-}Lite_{\mathcal{R}}$ , we give a notion of *approximation* that is sound for query answering. We develop techniques for rewriting and approximation of OBDA specifications based on compiling the extra expressiveness into the mappings. We target rather expressive ontology languages, and for Horn- $\mathcal{ALCHIQ}$ , a Horn fragment of OWL2, we study decidability of existence of OBDA rewritings, and techniques to compute them when they exist, and to approximate them, otherwise.

We have implemented our techniques in a prototype system called ONTOPROX, which exploits functionalities provided by the ONTOP (Rodríguez-Muro, Kontchakov, and Zakharyashev 2013) and CLIPPER systems (Eiter et al. 2012) to rewrite or approximate an OBDA specification expressed in Horn- $\mathcal{SHIQ}$  to one that can be directly processed by any OBDA system. We have evaluated ONTOPROX over synthetic and real OBDA instances against (i) the default ONTOP behavior, (ii) local semantic approximation (LSA), (iii) global semantic approximation (GSA), and (iv) CLIPPER over materialized ABoxes. We observe that using ONTOPROX, for a few queries we have been able to obtain more answers (in fact, complete answers, as confirmed by CLIPPER). However, for many queries ONTOPROX showed no difference with respect to the default ONTOP behavior. One reason for this is that in the considered real-world scenario, the mapping designers put significant effort to manually create complex mappings that overcome the limitations of  $DL\text{-}Lite_{\mathcal{R}}$ . Essentially they followed the principle of the technique pre-

sented here, and therefore produced an OBDA specification that was already “complete” by design.

The observations above immediately suggest a significant practical value of our approach, which can be used to facilitate the design of new OBDA specifications for existing expressive ontologies: instead of a manual compilation, which is cumbersome, error-prone, and difficult to maintain, mapping designers can write straightforward mappings, and the resulting OBDA specification can then be automatically transformed into a  $DL\text{-}Lite_{\mathcal{R}}$  OBDA specification with rich mappings.

The paper is structured as follows. In Section 2, we provide some preliminary notions, and in Section 3, we present our framework of OBDA rewriting and approximation. In Section 4, we illustrate a technique for computing the OBDA-rewriting of a given Horn- $\mathcal{ALCHIQ}$  specification. In Section 5, we address the problem of OBDA-rewritability, and show how to obtain an approximation when a rewriting does not exist. In Section 6, we discuss our prototype ONTOPROX and experiments. Finally, in Section 7, we conclude the paper. Omitted proofs can be found in the extended version of this paper (Botoeva et al. 2015).

## 2 Preliminaries

We give some basic notions about ontologies and OBDA.

### 2.1 Ontologies

We assume to have the following pairwise disjoint countably infinite alphabets:  $N_C$  of *concept names*,  $N_R$  of *role names*, and  $N_I$  of constants (also called *individuals*). We consider ontologies expressed in Description Logics (DLs). Here we present the logics Horn- $\mathcal{ALCHIQ}$ , the Horn fragment of  $\mathcal{SHIQ}$  without role transitivity, and  $DL\text{-}Lite_{\mathcal{R}}$ , for which we develop some of the technical results in the paper. However, the general approximation framework is applicable to any fragment of OWL 2.

A Horn- $\mathcal{ALCHIQ}$  TBox in normal form is a finite set of axioms: *concept inclusions* ( $CIs$ )  $\prod_i A_i \sqsubseteq C$ , *role inclusions* ( $RI$ s)  $R_1 \sqsubseteq R_2$  and *role disjointness* axioms  $R_1 \sqcap R_2 \sqsubseteq \perp$ , where  $A, A_i$  denote concept names,  $R, R_1, R_2$  denote role names  $P$  or their inverses  $P^-$ , and  $C$  denotes a concept of the form  $\perp, A, \exists R.A, \forall R.A$ , or  $\leq 1 R.A$  (Kazakov 2009). For an inverse role  $R = P^-$ , we use  $R^-$  to denote  $P$ .  $\perp$  denotes the empty concept/role. A  $DL\text{-}Lite_{\mathcal{R}}$  TBox is a finite set of axioms of the form  $B_1 \sqsubseteq B_2, B_1 \sqcap B_2 \sqsubseteq \perp, R_1 \sqsubseteq R_2$ , and  $R_1 \sqcap R_2 \sqsubseteq \perp$ , where  $B_i$  denotes a concept of the form  $A$  or  $\exists R.T$ . In what follows, for simplicity we write  $\exists R$  instead of  $\exists R.T$ , and we use  $N$  to denote either a concept or a role name. We also assume that all TBoxes are in normal form.

An *ABox* is a finite set of *membership assertions* of the form  $A(c)$  or  $P(c, c')$ , where  $c, c' \in N_I$ . For a DL  $\mathcal{L}$ , an  $\mathcal{L}$ -ontology is a pair  $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$ , where  $\mathcal{T}$  is an  $\mathcal{L}$ -TBox and  $\mathcal{A}$  is an ABox. A *signature*  $\Sigma$  is a finite set of concept and role names. An ontology  $\mathcal{O}$  is said to be defined over (or simply, *over*)  $\Sigma$  if all the concept and role names occurring in it belong to  $\Sigma$  (and likewise for TBoxes, ABoxes, concept inclusions, etc.). When  $\mathcal{T}$  is over  $\Sigma$ , we denote by  $\text{sig}(\mathcal{T})$  the subset of  $\Sigma$  actually occurring in  $\mathcal{T}$ . Moreover we denote with  $\text{Ind}(\mathcal{A})$ , the set of individuals appearing in  $\mathcal{A}$ .

The semantics, models, and the notions of satisfaction and consistency of ontologies are defined in the standard way. We only point out that we adopt the *Unique Name Assumption* (UNA), and for simplicity we also assume to have *standard names*, i.e., for every interpretation  $\mathcal{I}$  and every constant  $c \in \mathbb{N}_1$  interpreted by  $\mathcal{I}$ , we have that  $c^{\mathcal{I}} = c$ .

## 2.2 OBDA and Mappings

Let  $\mathcal{S}$  be a relational schema over a countably infinite set  $\mathbb{N}_S$  of database predicates. For simplicity, we assume to deal with plain relational schemas without constraints, and with database instances that directly store abstract objects (as opposed to values). In other words, a database instance  $\mathcal{D}$  of  $\mathcal{S}$  is a set of ground atoms over the predicates in  $\mathbb{N}_S$  and the constants in  $\mathbb{N}_1$ .<sup>2</sup> Queries over  $\mathcal{S}$  are expressed in SQL. We use  $\varphi(\vec{x})$  to denote that query  $\varphi$  has  $\vec{x} = x_1, \dots, x_n$  as free (i.e., answer) variables, where  $n$  is the arity of  $\varphi$ . Given a database instance  $\mathcal{D}$  of  $\mathcal{S}$  and a query  $\varphi$  over  $\mathcal{S}$ ,  $ans(\varphi, \mathcal{D})$  denotes the set of tuples of constants in  $\mathbb{N}_1$  computed by evaluating  $\varphi$  over  $\mathcal{D}$ .

In OBDA, one provides access to an (external) database through an ontology TBox, which is connected to the database by means of a mapping. Given a source schema  $\mathcal{S}$  and a TBox  $\mathcal{T}$ , a (GAV) *mapping assertion* between  $\mathcal{S}$  and  $\mathcal{T}$  has the form  $\varphi(x) \rightsquigarrow A(x)$  or  $\varphi'(x, x') \rightsquigarrow P(x, x')$ , where  $A$  and  $P$  are respectively concept and role names, and  $\varphi(x)$ ,  $\varphi'(x, x')$  are arbitrary (SQL) queries expressed over  $\mathcal{S}$ . Intuitively, given a database instance  $\mathcal{D}$  of  $\mathcal{S}$  and a mapping assertion  $m = \varphi(x) \rightsquigarrow A(x)$ , the instances of the concept  $A$  generated by  $m$  from  $\mathcal{D}$  is the set  $ans(\varphi, \mathcal{D})$ ; similarly for a mapping assertion  $\varphi(x, x') \rightsquigarrow P(x, x')$ .

An *OBDA specification* is a triple  $\mathcal{P} = \langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$ , where  $\mathcal{T}$  is a DL TBox,  $\mathcal{S}$  is a relational schema, and  $\mathcal{M}$  is a finite set of mapping assertions. Without loss of generality, we assume that all concept and role names appearing in  $\mathcal{M}$  are contained in  $\text{sig}(\mathcal{T})$ . An *OBDA instance* is a pair  $\langle \mathcal{P}, \mathcal{D} \rangle$ , where  $\mathcal{P}$  is an OBDA specification, and  $\mathcal{D}$  is a database instance of  $\mathcal{S}$ . The semantics of the OBDA instance  $\langle \mathcal{P}, \mathcal{D} \rangle$  is specified in terms of interpretations of the concepts and roles in  $\mathcal{T}$ . We define it by relying on the following (*virtual*<sup>3</sup>) ABox

$$\mathcal{A}_{\mathcal{M}, \mathcal{D}} = \{N(\vec{o}) \mid \vec{o} \in ans(\varphi, \mathcal{D}) \text{ and } \varphi(\vec{x}) \rightsquigarrow N(\vec{x}) \text{ in } \mathcal{M}\}$$

generated by  $\mathcal{M}$  from  $\mathcal{D}$ , where  $N$  is a concept or role name in  $\mathcal{T}$ . Then, a model of  $\langle \mathcal{P}, \mathcal{D} \rangle$  is simply a model of the ontology  $\langle \mathcal{T}, \mathcal{A}_{\mathcal{M}, \mathcal{D}} \rangle$ .

Following Di Pinto et al. (2013), we split each mapping assertion  $m = \varphi(\vec{x}) \rightsquigarrow N(\vec{x})$  in  $\mathcal{M}$  into two parts by introducing an intermediate view name  $V_m$  for the SQL query  $\varphi(\vec{x})$ . We obtain a *low-level* mapping assertion of the form  $\varphi(\vec{x}) \rightsquigarrow V_m(\vec{x})$ , and a *high-level* mapping assertion of the form  $V_m(\vec{x}) \rightsquigarrow N(\vec{x})$ . In our technical development, we deal only with the high-level mappings. Hence, we abstract away the low-level mapping part, and in the following we directly consider the intermediate views as our data sources.

<sup>2</sup>All our results easily extend to the case where objects are constructed from retrieved database values (Calvanese et al. 2009).

<sup>3</sup>We call such an ABox ‘virtual’, because we are not interested in actually materializing its facts.

## 2.3 Query Answering

We consider conjunctive queries, which are the basic and most important querying mechanism in relational database systems and ontologies. A *conjunctive query* (CQ)  $q(\vec{x})$  over a signature  $\Sigma$  is a formula  $\exists \vec{y}. \varphi(\vec{x}, \vec{y})$ , where  $\varphi$  is a conjunction of atoms  $N(\vec{z})$ , such that  $N$  is a concept or role name in  $\Sigma$ , and  $\vec{z}$  are variables from  $\vec{x}$  and  $\vec{y}$ . The set of *certain answers* to a CQ  $q(\vec{x})$  over an ontology  $\langle \mathcal{T}, \mathcal{A} \rangle$ , denoted  $cert(q, \langle \mathcal{T}, \mathcal{A} \rangle)$ , is the set of tuples  $\vec{c}$  of elements from  $\text{Ind}(\mathcal{A})$  of the same length as  $\vec{x}$ , such that  $q(\vec{c})$  (considered as a FO sentence) holds in every model of  $\langle \mathcal{T}, \mathcal{A} \rangle$ . We mention two more query classes. An *atomic query* (AQ) is a CQ consisting of exactly one atom whose variables are all free. A *CQ with inequalities* ( $CQ^\neq$ ) is a CQ that may contain inequality atoms between the variables of the predicate atoms.

Given a CQ  $q$ , an OBDA specification  $\mathcal{P} = \langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  and a database instance  $\mathcal{D}$  of  $\mathcal{S}$ , the answer to  $q$  over the OBDA instance  $\langle \mathcal{P}, \mathcal{D} \rangle$ , denoted  $cert(q, \mathcal{P}, \mathcal{D})$ , is defined as  $cert(q, \langle \mathcal{T}, \mathcal{A}_{\mathcal{M}, \mathcal{D}} \rangle)$ . Observe that, when  $\mathcal{D}$  is inconsistent with  $\mathcal{P}$  (i.e.,  $\langle \mathcal{P}, \mathcal{D} \rangle$  does not have a model), then  $cert(q, \mathcal{P}, \mathcal{D})$  is the set of all possible tuples of constants in  $\mathcal{A}_{\mathcal{M}, \mathcal{D}}$  (of the same arity as  $q$ ).

## 3 An OBDA Rewriting Framework

We extend the notion of query inseparability of ontologies (Botoeva et al. 2014) to OBDA specifications. We adopt the proposal by Bienvenu and Rosati (2015), but we do not enforce preservation of inconsistency.

**Definition 1.** *Let  $\Sigma$  be a signature. Two OBDA specifications  $\mathcal{P}_1 = \langle \mathcal{T}_1, \mathcal{M}_1, \mathcal{S} \rangle$  and  $\mathcal{P}_2 = \langle \mathcal{T}_2, \mathcal{M}_2, \mathcal{S} \rangle$  are  $\Sigma$ -CQ inseparable if  $cert(q, \mathcal{P}_1, \mathcal{D}) = cert(q, \mathcal{P}_2, \mathcal{D})$ , for every CQ  $q$  over  $\Sigma$  and every database instance  $\mathcal{D}$  of  $\mathcal{S}$ .*

In OBDA, one must deal with the trade-off between the computational complexity of query answering and the expressiveness of the ontology language. Suppose that for an OBDA specification  $\mathcal{P} = \langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$ ,  $\mathcal{T}$  is expressed in an ontology language  $\mathcal{L}$  that does not allow for efficient query answering. A possible solution is to exploit the expressive power of the mapping layer to compute a new OBDA specification  $\mathcal{P}' = \langle \mathcal{T}', \mathcal{M}', \mathcal{S} \rangle$  in which  $\mathcal{T}'$  is expressed in a language  $\mathcal{L}_t$  more suitable for query answering than  $\mathcal{L}$ . The aim is to encode in  $\mathcal{M}'$  not only  $\mathcal{M}$  but also part of the semantics of  $\mathcal{T}$ , so that  $\mathcal{P}'$  is query-inseparable from  $\mathcal{P}$ . This leads to the notion of rewriting of OBDA specifications.

**Definition 2.** *Let  $\mathcal{L}_t$  be an ontology language. The OBDA specification  $\mathcal{P}' = \langle \mathcal{T}', \mathcal{M}', \mathcal{S} \rangle$  is a CQ-rewriting in  $\mathcal{L}_t$  of the OBDA specification  $\mathcal{P} = \langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  if (i)  $\text{sig}(\mathcal{T}) \subseteq \text{sig}(\mathcal{T}')$ , (ii)  $\mathcal{T}'$  is an  $\mathcal{L}_t$ -TBox, and (iii)  $\mathcal{P}$  and  $\mathcal{P}'$  are  $\Sigma$ -CQ inseparable, for  $\Sigma = \text{sig}(\mathcal{T})$ . If such  $\mathcal{P}'$  exists, we say that  $\mathcal{P}$  is CQ-rewritable into  $\mathcal{L}_t$ .*

We observe that the new OBDA specification can be defined over a signature that is an extension of that of the original TBox. This is specified by condition (i). In condition (ii), we impose that the new ontology is specified in the target language  $\mathcal{L}_t$ . Finally, condition (iii) imposes that the OBDA specifications cannot be distinguished by CQs over the original TBox. Note that the definition allows for changing the

ontology and the mappings, but not the source schema, accounting for the fact that the data sources might not be under the control of the designer of the OBDA specification.

As expected, it is not always possible to obtain a CQ-rewriting of  $\mathcal{P}$  in an ontology language  $\mathcal{L}_t$  that allows for efficient query answering. Indeed, the combined expressiveness of  $\mathcal{L}_t$  with the new mappings might not be sufficient to simulate query answering over  $\mathcal{P}$  without loss. In these cases, we can resort to approximating query answers over  $\mathcal{P}$  in a *sound* way, which means that the answers to queries posed over the new specification are contained in those produced by querying  $\mathcal{P}$ . Hence, we say that the OBDA specification  $\mathcal{P}' = \langle \mathcal{T}', \mathcal{M}', \mathcal{S} \rangle$  is a *sound CQ-approximation* in  $\mathcal{L}_t$  of the OBDA specification  $\mathcal{P} = \langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  if  $\mathcal{P}'$  satisfies (i), (ii), and  $\text{cert}(q, \mathcal{P}', \mathcal{D}) \subseteq \text{cert}(q, \mathcal{P}, \mathcal{D})$ , for each CQ  $q$  over  $\text{sig}(\mathcal{T})$  and for each instance  $\mathcal{D}$  of  $\mathcal{S}$ .

Next, we study CQ-rewritability of OBDA specifications into *DL-Lite $\mathcal{R}$* , developing suitable techniques.

## 4 Rewriting OBDA Specifications

In this section, we develop our OBDA rewriting technique, which relies on Datalog rewritings of the TBox (and mappings). Recall that a *Datalog program* (with inequalities) is a finite set of definite *Horn* clauses *without* functions symbols, i.e., rules of the form  $\text{head} \leftarrow \varphi$ , where  $\varphi$  is a finite non-empty list of predicate atoms and guarded inequalities called the body of the rule, and  $\text{head}$  is an atom, called the head of the rule, all of whose variables occur in the body. The predicates that occur in rule heads are called *intensional* (IDB), the other predicates are called *extensional* (EDB).

### 4.1 ET-mappings

Now, we extend the notion of T-mappings introduced by Rodriguez-Muro, Kontchakov, and Zakharyashev (2013), and define the notion of an ET-mapping that results from compiling into the mapping the expressiveness of ontology languages that are Datalog rewritable, as introduced below.

We first introduce notation we need. Let  $\Pi$  be a Datalog program and  $N$  an IDB predicate. For a database  $\mathcal{D}$  over the EDB predicates of  $\Pi$ , let  $N_{\Pi}^i(\mathcal{D})$  denote the set of facts about  $N$  that can be deduced from  $\mathcal{D}$  by at most  $i \geq 1$  applications of the rules in  $\Pi$ , and let  $N_{\Pi}^{\infty}(\mathcal{D}) = \bigcup_{i \geq 1} N_{\Pi}^i(\mathcal{D})$ . It is known that the predicate  $N_{\Pi}^{\infty}(\cdot)$  defined by  $N$  in  $\Pi$  can be characterized by a possibly infinite union of CQ $\neq$ s (Cosmadakis et al. 1988), i.e., there exist CQ $\neq$ s  $\varphi_0^N, \varphi_1^N, \dots$  such that  $N_{\Pi}^{\infty}(\mathcal{D}) = \bigcup_{i \geq 0} \{N(\vec{a}) \mid \vec{a} \in \text{ans}(\varphi_i^N, \mathcal{D})\}$ , for every  $\mathcal{D}$ . The  $\varphi_i^N$ 's are called the *expansions* of  $N$  and can be described in terms of expansion trees; cf. (Botoeva et al. 2015, Appendix A). We denote by  $\Phi_{\Pi}(N)$  the set of expansion trees for  $N$  in  $\Pi$ , and abusing notation also the (possibly infinite) union of CQ $\neq$ s corresponding to it. Note that  $\Phi_{\Pi}(N)$  might be infinite due to the presence of IDB predicates that are *recursive*, i.e., either directly or indirectly refer to themselves.

We call a TBox  $\mathcal{T}$  *Datalog rewritable* if it admits a translation  $\Pi_{\mathcal{T}}$  to Datalog that preserves consistency and answers to AQs (see, e.g., the translations by Hustadt, Motik, and Sattler (2005), Eiter et al. (2012), and Trivela et al. (2015)

for Horn-*SHIQ*, and by Cuenca Grau et al. (2013) for *SHI*). We assume that  $\Pi_{\mathcal{T}}$  makes use of a special nullary predicate  $\perp$  that encodes inconsistency, i.e., for an ABox  $\mathcal{A}$ ,  $\langle \mathcal{T}, \mathcal{A} \rangle$  is consistent iff  $\perp_{\Pi_{\mathcal{T}}}^{\infty}(\mathcal{A})$  is empty.<sup>4</sup> We also assume that  $\Pi_{\mathcal{T}}$  includes the following auxiliary rules, which ensure that  $\Pi_{\mathcal{T}}$  derives all possible facts constructed over  $\text{sig}(\mathcal{T})$  and  $\text{Ind}(\mathcal{A})$  whenever  $\langle \mathcal{T}, \mathcal{A} \rangle$  is inconsistent:

$$\begin{aligned} \top_{\Delta}(x) &\leftarrow A(x); \top_{\Delta}(x) \leftarrow P(x, y); \top_{\Delta}(y) \leftarrow P(x, y); \\ A(x) &\leftarrow \perp, \top_{\Delta}(x); P(x, y) \leftarrow \perp, \top_{\Delta}(x), \top_{\Delta}(y); \end{aligned}$$

where  $A$  and  $P$  respectively range over concept and role names in  $\text{sig}(\mathcal{T})$ , and  $\top_{\Delta}$  is a fresh unary predicate denoting the set of all the individuals appearing in  $\mathcal{A}$ .

In the following, we denote with  $\Pi_{\mathcal{M}}$  the (high-level) mapping  $\mathcal{M}$  viewed as a Datalog program, and with  $\Pi_{\mathcal{T}, \mathcal{M}}$  the Datalog program  $\Pi_{\mathcal{T}} \cup \Pi_{\mathcal{M}}$  associated to a Datalog rewritable TBox  $\mathcal{T}$  and a mapping  $\mathcal{M}$ . From the properties of the translation  $\Pi_{\mathcal{T}}$  (and the simple structure of  $\Pi_{\mathcal{M}}$ ), we obtain that  $\Pi_{\mathcal{T}, \mathcal{M}}$  satisfies the following:

**Lemma 3.** *Let  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  be an OBDA specification where  $\mathcal{T}$  is Datalog rewritable. Then, for every database instance  $\mathcal{D}$  of  $\mathcal{S}$ , concept or role name  $N$  of  $\mathcal{T}$ , and  $\vec{a}$  in  $\text{Ind}(\mathcal{A}_{\mathcal{M}, \mathcal{D}})$ , we have that  $\langle \mathcal{T}, \mathcal{A}_{\mathcal{M}, \mathcal{D}} \rangle \models N(\vec{a})$  iff  $N(\vec{a}) \in N_{\Pi_{\mathcal{T}, \mathcal{M}}}^{\infty}(\mathcal{D})$ .*

For a predicate  $N$ , we say that an expansion  $\varphi^N \in \Phi_{\Pi_{\mathcal{T}, \mathcal{M}}}(N)$  is *DB-defined* if  $\varphi^N$  is defined over database predicates. Now we are ready to define ET-mappings.

**Definition 4.** *Let  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  be an OBDA specification where  $\mathcal{T}$  is Datalog rewritable. The ET-mapping for  $\mathcal{M}$  and  $\mathcal{T}$ , denoted  $\text{etm}_{\mathcal{T}}(\mathcal{M})$ , is defined as the set of assertions of the form  $\varphi^N(\vec{x}) \rightsquigarrow N(\vec{x})$  such that  $N$  is a concept or role name in  $\mathcal{T}$ , and  $\varphi^N \in \Phi_{\Pi_{\mathcal{T}, \mathcal{M}}}(N)$  is DB-defined.*

It is easy to show that, for  $\mathcal{M}' = \text{etm}_{\mathcal{T}}(\mathcal{M})$  and each database instance  $\mathcal{D}$ , the virtual ABox  $\mathcal{A}_{\mathcal{M}', \mathcal{D}}$  (which can be defined for ET-mappings as for ordinary mappings) contains all facts entailed by  $\langle \mathcal{T}, \mathcal{A}_{\mathcal{M}, \mathcal{D}} \rangle$ . In this sense, the ET-mapping  $\text{etm}_{\mathcal{T}}(\mathcal{M})$  plays for a Datalog rewritable TBox  $\mathcal{T}$  the same role as T-mappings play for (the simpler) *DL-Lite $\mathcal{R}$*  TBoxes. Note that, in general, an ET-mapping is not a mapping, as it may contain infinitely many assertions. However,  $\mathcal{A}_{\mathcal{M}', \mathcal{D}}$  is still finite, given that it is constructed over the finite number of constants appearing in  $\mathcal{D}$ .

### 4.2 Rewriting Horn-*ALCHIQ* OBDA Specifications to *DL-Lite $\mathcal{R}$*

Let  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  be an OBDA specification, where  $\mathcal{T}$  is a Horn-*ALCHIQ* TBox over a signature  $\Sigma$ . Figure 1 describes the algorithm  $\text{RewObda}(\mathcal{T}, \mathcal{M})$ , which constructs a *DL-Lite $\mathcal{R}$*  TBox  $\mathcal{T}_r$  and an ET-mapping  $\mathcal{M}_c$  such that  $\langle \mathcal{T}_r, \mathcal{M}_c, \mathcal{S} \rangle$  is  $\Sigma$ -CQ inseparable from  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$ .

In Step 2, the algorithm applies to  $\mathcal{T}_1$  the normalization procedure  $\text{norm}_{\exists}$ , which gets rid of concepts of the form  $\exists R.(\prod A_j)$  in the right-hand side of CIs. This is achieved by the following well-known substitution (Artale et al. 2009): every CI  $\prod_{i=1}^m A_i \sqsubseteq \exists R.(\prod_{j=1}^n A_j)$  in  $\mathcal{T}_1$  is replaced with  $\prod_{i=1}^m A_i \sqsubseteq \exists P_{\text{new}}. P_{\text{new}} \sqsubseteq R$ , and  $\top \sqsubseteq \forall P_{\text{new}}. A_j$ , for  $1 \leq$

<sup>4</sup>Here we simply consider  $\mathcal{A}$  as a database.

**Input:** Horn- $\mathcal{ALCHIQ}$  TBox  $\mathcal{T}$  and mapping  $\mathcal{M}$ .  
**Output:**  $DL\text{-Lite}_{\mathcal{R}}$  TBox  $\mathcal{T}_r$  and ET-mapping  $\mathcal{M}_c$ .  
**Step 1:**  $\mathcal{T}_1$  is obtained from  $\mathcal{T}$  by adding all CIs of the form  $\bigcap A_i \sqsubseteq \exists R.(\bigcap A'_j)$  entailed by  $\mathcal{T}$ , for concept names  $A_i, A'_j \in \text{sig}(\mathcal{T})$ .  
**Step 2:**  $\mathcal{T}_2 = \text{norm}_{\exists}(\mathcal{T}_1)$ .  
**Step 3:**  $\mathcal{T}_3 = \text{norm}_{\sqcap}(\mathcal{T}_2)$ .  
**Step 4:**  $\mathcal{M}_c$  is  $\text{etm}_{\mathcal{T}_3}(\mathcal{M})$ , and  $\mathcal{T}_r$  is the  $DL\text{-Lite}_{\mathcal{R}}$  TBox consisting of all  $DL\text{-Lite}_{\mathcal{R}}$  axioms over  $\text{sig}(\mathcal{T}_3)$  entailed by  $\mathcal{T}_3$  (including the trivial ones  $N \sqsubseteq N$ ).

Figure 1: OBDA specification rewriting algorithm RewObda.

$j \leq n$ , where  $P_{new}$  is a fresh role name. Notice that the latter two forms of inclusions introduced by  $\text{norm}_{\exists}$  are actually in  $DL\text{-Lite}_{\mathcal{R}}$ , as  $\top \sqsubseteq \forall P_{new}.A'_j$  is equivalent to  $\exists P_{new}^- \sqsubseteq A'_j$ . In Step 3, the algorithm applies to  $\mathcal{T}_2$  a further normalization procedure,  $\text{norm}_{\sqcap}$ , which introduces a fresh concept name  $A_{A_1 \sqcap \dots \sqcap A_n}$  for each concept conjunction  $A_1 \sqcap \dots \sqcap A_n$  appearing in  $\mathcal{T}_2$ , and adds  $A_1 \sqcap \dots \sqcap A_n \equiv A_{A_1 \sqcap \dots \sqcap A_n}$ <sup>5</sup> to the TBox. Note that  $\text{norm}_{\exists}(\mathcal{T}_1)$  and  $\text{norm}_{\sqcap}(\mathcal{T}_2)$  are model-conservative extensions of  $\mathcal{T}_1$  and  $\mathcal{T}_2$ , respectively (Lutz, Walther, and Wolter 2007), as one can easily show. We denote by  $\text{rew}(\mathcal{T})$  the resulting TBox  $\mathcal{T}_r$ , which in general is exponential in the size of  $\mathcal{T}$ , and by  $\text{comp}(\mathcal{T}, \mathcal{M})$  the resulting ET-mapping  $\mathcal{M}_c$ , which in general is infinite.

**Example 5.** Assume that the domain knowledge is represented by the axiom about bank accounts from Section 1. The normalization of this axiom is the TBox  $\mathcal{T}^b = \{\text{Person} \sqsubseteq \forall \text{inNameOf}^- . A_1, \text{CAcc} \sqcap A_1 \sqsubseteq \text{SAcc}\}$ . Assume that the database schema  $\mathcal{S}^b$  consists of the two relations  $\text{ENT}(\text{ID}, \text{TYPE}, \text{EMPID})$ ,  $\text{PROD}(\text{NUM}, \text{TYPE}, \text{CUSTID})$ , whose data are mapped to the ontology terms by means of the following mapping  $\mathcal{M}$ :

```

mP: SELECT ID AS X FROM ENT WHERE ENT.TYPE='P' ~> Person(X)
mN: SELECT NUM AS X, CUSTID AS Y FROM PROD ~> inNameOf(X, Y)
mC: SELECT NUM AS X FROM PROD P WHERE P.TYPE='B' ~> CAcc(X)

```

We will work with the corresponding high-level mapping  $\mathcal{M}^b$  consisting of the assertions:

$$\begin{aligned}
h_P &: \{x \mid V_{\text{Person}}(x)\} \rightsquigarrow \text{Person}(x) \\
h_N &: \{x, y \mid V_{\text{inNameOf}}(x, y)\} \rightsquigarrow \text{inNameOf}(x, y) \\
h_C &: \{x \mid V_{\text{CAcc}}(x)\} \rightsquigarrow \text{CAcc}(x)
\end{aligned}$$

Now, consider the OBDA specification  $\mathcal{P}^b = \langle \mathcal{T}^b, \mathcal{M}^b, \mathcal{S}^b \rangle$ . The RewObda algorithm invoked on  $(\mathcal{T}^b, \mathcal{M}^b)$  produces:

- The intermediate TBoxes  $\mathcal{T}_1^b$  and  $\mathcal{T}_2^b$  coinciding with  $\mathcal{T}^b$ , and  $\mathcal{T}_3^b$  extending  $\mathcal{T}^b$  with  $A_{\text{CAcc} \sqcap A_1} \equiv \text{CAcc} \sqcap A_1$ .
- The ET-mapping  $\mathcal{M}_c^b = \text{etm}_{\mathcal{T}_3^b}(\mathcal{M}^b)$ , which extends  $\mathcal{M}^b$  with the assertions  $\{x \mid V_{\text{inNameOf}}(x, y), V_{\text{Person}}(y)\} \rightsquigarrow A_1(x)$ ,  $\{x \mid V_{\text{CAcc}}(x), V_{\text{inNameOf}}(x, y), V_{\text{Person}}(y)\} \rightsquigarrow \text{SAcc}(x)$ , and  $\{x \mid V_{\text{CAcc}}(x), V_{\text{inNameOf}}(x, y), V_{\text{Person}}(y)\} \rightsquigarrow A_{\text{CAcc} \sqcap A_1}(x)$ .

<sup>5</sup>We use '≡' to abbreviate inclusion in both directions.

The algorithm returns the  $DL\text{-Lite}_{\mathcal{R}}$  TBox  $\mathcal{T}_r^b = \{A_{\text{CAcc} \sqcap A_1} \sqsubseteq \text{CAcc}, A_{\text{CAcc} \sqcap A_1} \sqsubseteq A_1, A_{\text{CAcc} \sqcap A_1} \sqsubseteq \text{SAcc}\}$  and the mapping  $\mathcal{M}_c^b$ . It is possible to show that  $\mathcal{P}_{DL\text{-Lite}_{\mathcal{R}}}^b = \langle \mathcal{T}_r^b, \mathcal{M}_c^b, \mathcal{S}^b \rangle$  is a CQ-rewriting of  $\mathcal{P}^b$  into  $DL\text{-Lite}_{\mathcal{R}}$ . ■

The TBox  $\mathcal{T}_3$  obtained as an intermediate result in Step 3 of RewObda( $\mathcal{T}, \mathcal{M}$ ), is a model-conservative extension of  $\mathcal{T}$  that is tailored towards capturing in  $DL\text{-Lite}_{\mathcal{R}}$  the answers to tree-shaped CQs. This is obtained by introducing in Step 2 sufficiently new role names, and in Step 3 new concept names, so as to capture entailed axioms that generate the tree-shaped parts of models. On the other hand, the ET-mapping  $\mathcal{M}_c = \text{comp}(\mathcal{T}, \mathcal{M})$  is such that it generates from a database instance a virtual ABox  $\mathcal{A}^v$  that is complete with respect to all ABox facts that might be involved in the generation of the tree-shaped parts of models of  $\mathcal{T}_r$  and  $\mathcal{A}^v$ . This allows us to prove the main result of this section.

**Theorem 6.** *Let  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  be an OBDA specification such that  $\mathcal{T}$  is a Horn- $\mathcal{ALCHIQ}$  TBox, and let  $\langle \mathcal{T}_r, \mathcal{M}_c \rangle = \text{RewObda}(\mathcal{T}, \mathcal{M})$ . Then  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  and  $\langle \mathcal{T}_r, \mathcal{M}_c, \mathcal{S} \rangle$  are  $\Sigma$ -CQ inseparable, for  $\Sigma = \text{sig}(\mathcal{T})$ .*

Clearly,  $\langle \mathcal{T}_r, \mathcal{M}_c, \mathcal{S} \rangle$  is a candidate for being a CQ-rewriting of  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  into  $DL\text{-Lite}_{\mathcal{R}}$ . However, since  $\mathcal{M}_c$  might be an infinite set,  $\langle \mathcal{T}_r, \mathcal{M}_c, \mathcal{S} \rangle$  might not be an OBDA specification and hence might not be effectively usable for query answering. Next we address this issue, and show that in some cases we obtain proper CQ-rewritings, while in others we have to resort to approximations.

## 5 Approximating OBDA Specifications

To obtain from an ET-mapping a proper mapping, we exploit the notion of predicate boundedness in Datalog, and use a bound on the depth of Datalog expansion trees.

An IDB predicate  $N$  is said to be *bounded* in a Datalog program  $\Pi$ , if there exists a constant  $k$  depending only on  $\Pi$  such that, for every database  $\mathcal{D}$ , we have  $N_{\Pi}^k(\mathcal{D}) = N_{\Pi}^{\infty}(\mathcal{D})$  (Cosmadakis et al. 1988). If  $N$  is bounded in  $\Pi$ , then there exists an equivalent Datalog program  $\Pi'$  such that  $\Phi_{\Pi'}(N)$  is *finite*, and thus represents a finite union of CQ $^{\neq}$ s. It is well known that predicate boundedness for Datalog is undecidable in general (Gaifman et al. 1987). We say that  $\Omega$  is a *boundedness oracle* if for a Datalog program  $\Pi$  and a predicate  $N$  it returns one of the three answers:  $N$  is bounded in  $\Pi$ ,  $N$  is not bounded in  $\Pi$ , or unknown. When  $N$  is bounded,  $\Omega$  returns also a *finite* union of CQ $^{\neq}$ s, denoted  $\Omega_{\Pi}(N)$ , defining  $N$ . Given a constant  $k$ ,  $\Phi_{\Pi}^k(N)$  denotes the set of trees (and the corresponding union of CQ $^{\neq}$ s) in  $\Phi_{\Pi}(N)$  of depth at most  $k$ , hence  $\Phi_{\Pi}^k(N)$  is always finite.

We introduce a *cutting operator*  $\text{cut}_k^{\Omega}$ , which is parametric with respect to the cutting depth  $k > 0$  and the boundedness oracle  $\Omega$ , which, when applied to a predicate  $N$  and a Datalog program  $\Pi$ , returns a finite union of CQ $^{\neq}$ s as follows:

$$\text{cut}_k^{\Omega}(N, \Pi) = \begin{cases} \Omega_{\Pi}(N), & \text{if } N \text{ is bounded in } \Pi \text{ w.r.t. } \Omega \\ \Phi_{\Pi}^k(N), & \text{otherwise.} \end{cases}$$

We apply cutting also to ET-mappings: given an ET-mapping  $\text{etm}_{\mathcal{T}}(\mathcal{M})$ , the *mapping*  $\text{cut}_k^{\Omega}(\text{etm}_{\mathcal{T}}(\mathcal{M}))$  is the (finite) set of mapping assertions  $\varphi^N(\vec{x}) \rightsquigarrow N(\vec{x})$  s.t.  $N$  is a concept or role name in  $\mathcal{T}$ , and  $\varphi^N \in \text{cut}_k^{\Omega}(N, \Pi_{\mathcal{T}, \mathcal{M}})$  is DB-defined.

The following theorem provides a sufficient condition for CQ-rewritability into  $DL\text{-}Lite_{\mathcal{R}}$  in terms of the well-known notion of first-order (FO)-rewritability, which we recall here: a query  $q$  is *FO-rewritable* with respect to a TBox  $\mathcal{T}$ , if there exists a FO query  $q'$  such that  $cert(q, \langle \mathcal{T}, \mathcal{A} \rangle) = ans(q', \mathcal{A})$ , for every ABox  $\mathcal{A}$  over  $sig(\mathcal{T})$  (viewed as a database). It uses the fact that if an AQ is FO-rewritable with respect to a Horn- $\mathcal{ALCHIQ}$  TBox  $\mathcal{T}$ , then it is actually rewritable into a union of CQ $\neq$ s, and the fact that if  $\mathcal{T}$  is FO-rewritable for AQs (i.e., every AQ is FO-rewritable with respect to  $\mathcal{T}$ ), then each concept and role name is bounded in  $\Pi_{\mathcal{T}}$  (Lutz and Wolter 2011; Biennu, Lutz, and Wolter 2013).

**Theorem 7.** *Let  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  be an OBDA specification such that  $\mathcal{T}$  is a Horn- $\mathcal{ALCHIQ}$  TBox. Further, let  $\mathcal{T}_r = rew(\mathcal{T})$  and  $\mathcal{M}' = cut_k^{\Omega}(\text{comp}(\mathcal{T}, \mathcal{M}))$ , for a boundedness oracle  $\Omega$  and some  $k > 0$ . If  $\mathcal{T}$  is FO-rewritable for AQs, then  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  is CQ-rewritable into  $DL\text{-}Lite_{\mathcal{R}}$ , and  $\langle \mathcal{T}_r, \mathcal{M}', \mathcal{S} \rangle$  is its CQ-rewriting. Otherwise,  $\langle \mathcal{T}_r, \mathcal{M}', \mathcal{S} \rangle$  is a sound CQ-approximation of  $\langle \mathcal{T}, \mathcal{M}, \mathcal{S} \rangle$  in  $DL\text{-}Lite_{\mathcal{R}}$ .*

The above result provides us with decidable conditions for rewritability of OBDA specifications in several significant cases. It is shown by Biennu, Lutz, and Wolter (2013) and Lutz and Wolter (2011) that FO-rewritability of AQs relative to Horn- $\mathcal{SHL}$ -TBoxes, Horn- $\mathcal{ALCF}$ -TBoxes, and Horn- $\mathcal{ALCF}$ -TBoxes of depth two is decidable. In fact, these FO-rewritability algorithms provide us with a boundedness oracle  $\Omega$ : for each concept and role name  $N$  in  $\mathcal{T}$ , they return a FO-rewriting of the AQ  $N(\vec{x})$  that combined with the mapping  $\mathcal{M}$  results in  $\Omega_{\Pi_{\mathcal{T}}, \mathcal{M}}(N)$ .

Unfortunately, a complete characterization of CQ-rewritability into  $DL\text{-}Lite_{\mathcal{R}}$  is not possible if arbitrary FO-queries are allowed in the (low-level) mapping.

**Theorem 8.** *The problem of checking whether an OBDA specification with an  $\mathcal{EL}$  ontology and FO source queries in the mapping is CQ-rewritable into  $DL\text{-}Lite_{\mathcal{R}}$  is undecidable.*

However, if we admit only unions of CQs in the (low-level) mapping, we can fully characterize CQ-rewritability.

**Theorem 9.** *The problem of checking whether an OBDA specification with a Horn- $\mathcal{ALCHI}$  ontology of depth one and unions of CQs as source queries in the mapping is CQ-rewritable into  $DL\text{-}Lite_{\mathcal{R}}$  is decidable.*

## 6 Implementation and Experiments

To demonstrate the feasibility of our OBDA specification rewriting technique, we have implemented a prototype system called ONTOPROX<sup>6</sup> and evaluated it over synthetic and real OBDA instances. Our system relies on the OBDA reasoner ONTOP<sup>7</sup> and the complete Horn- $\mathcal{SHIQ}$  CQ-answering system CLIPPER<sup>8</sup>, used as Java libraries. ONTOPROX also relies on a standard Prolog engine (SWI-PROLOG<sup>9</sup>) and on an OWL 2 reasoner (HERMIT<sup>10</sup>).

<sup>6</sup><https://github.com/ontop/ontoprox/>

<sup>7</sup><http://ontop.inf.unibz.it/>

<sup>8</sup><http://www.kr.tuwien.ac.at/research/systems/clipper/>

<sup>9</sup><http://www.swi-prolog.org/>

<sup>10</sup><http://hermit-reasoner.com/>

Essentially, ONTOPROX implements the rewriting and compiling procedure described in Figure 1, but instead of computing the (possibly infinite) ET-mapping  $\text{comp}(\mathcal{T}, \mathcal{M})$ , it computes its finite part  $\text{cut}_k(\text{comp}(\mathcal{T}, \mathcal{M}))$ . So, it gets as input an OWL 2 OBDA specification  $\langle \mathcal{T}_{\text{OWL2}}, \mathcal{M}, \mathcal{S} \rangle$  and a positive integer  $k$ , and produces a  $DL\text{-}Lite_{\mathcal{R}}$  OBDA specification that can be used with any OBDA system. Below we describe some of the implementation details:

- (1)  $\mathcal{T}_{\text{OWL2}}$  is first approximated to the Horn- $\mathcal{SHIQ}$  TBox  $\mathcal{T}$  by dropping the axioms outside this fragment.
- (2)  $\mathcal{T}$  is translated into a (possibly recursive) Datalog program  $\Pi$  and saturated with all CIs of the form  $\prod A_i \sqsubseteq \exists R.(\prod A'_j)$ , using functionalities provided by CLIPPER.
- (3) The expansions  $\text{cut}_k(\Phi_{\Pi}(X))$  are computed by an auxiliary Prolog program using Prolog meta-programming.
- (4) To produce actual mappings that can be used by an OBDA reasoner, the views in the high-level mapping  $\text{cut}_k(\text{comp}(\mathcal{T}, \mathcal{M}))$  are replaced with their original SQL definitions using functionalities of ONTOP.
- (5) The  $DL\text{-}Lite_{\mathcal{R}}$  closure is computed by relying on the OWL 2 reasoner for Horn- $\mathcal{SHIQ}$  TBox classification.

For the experiments, we have considered two scenarios:

**UOBM.** The university ontology benchmark (UOBM) (Ma et al. 2006) comes with a  $\mathcal{SHOIN}$  ontology (with 69 concepts, 35 roles, 9 attributes, and 204 TBox axioms), and an ABox generator. We have designed a database schema for the generated ABox, converted the ABox to a 10MB database instance for the schema, and manually created the mapping, consisting of 96 assertions<sup>11</sup>.

Among others, we have considered the following queries:

```

Q1u: SELECT DISTINCT ?X WHERE
      { ?X a ub:Person . }
Q2u: SELECT DISTINCT ?X WHERE
      { ?X a ub:Employee . }
Q3u: SELECT DISTINCT ?X ?Y WHERE
      { ?X rdf:type ub:ResearchGroup .
        ?X ub:subOrganizationOf ?Y . }
Q4u: SELECT DISTINCT ?X ?Y ?Z WHERE
      { ?X rdf:type ub:Chair .
        ?X ub:worksFor ?Y .
        ?Y rdf:type ub:Department .
        ?Y ub:subOrganizationOf ?Z . }

```

**Telecom benchmark.** The telecommunications ontology models a portion of the network of a leading telecommunications company, namely the portion connecting subscribers to the operating centers of their service providers. The current specification consists of an OWL 2 ontology with 152 concepts, 53 roles, 73 attributes, 458 TBox axioms, and of a mapping with 264 mapping assertions. The database instance contains 32GB of real-world data.

In the following, we only provide a description of some of the queries because the telecommunications ontology itself is bound by a confidentiality agreement.

<sup>11</sup><https://github.com/ontop/ontop-examples/tree/master/aaai-2016-ontoprox/uobm>

Table 1: Query evaluation with respect to 5 setups (number of answers / running time in seconds)

		ONTOP	LSA	GSA	ONTOPROX	CLIPPER
UOBM	$Q_1^u$	14,129 / 0.08	14,197 / 0.11	14,197 / 0.43	14,197 / 0.42	14,197 / 21.4
	$Q_2^u$	1,105 / 0.09	2,170 / 0.15	2,170 / 0.42	2,170 / 0.44	2,170 / 21.3
	$Q_3^u$	235 / 0.20	235 / 0.24	235 / 0.88	247 / 0.83	247 / 19.6
	$Q_4^u$	19 / 0.13	19 / 0.15	19 / 0.43	38 / 0.52	38 / 21.4
Telecom	$Q_1^t$	0 / 2.91	0 / 0.72	0 / 1.91	82,455 / 5.21	N/A
	$Q_2^t$	0 / 0.72	0 / 0.21	0 / 0.67	16,487 / 198	N/A
	$Q_3^t$	5,201,363 / 128	5,201,363 / 105	5,201,363 / 538	5,260,346 / 437	N/A

Table 2: ONTOPROX pre-computation time and output size

	UOBM	Telecom
Time (s)	8.47	8.72
Number of mapping assertions	441	907
Number of TBox axioms	294	620
Number of new concepts	26	60
Number of new roles	30	7

- Query  $Q_1^t$  asks, for each cable in the telecommunications network, the single segments of which the cable is composed, and the network line (between two devices) that the cable covers. For each cable, it also returns its bandwidth and its status (functioning, non-functioning, etc.).
- Query  $Q_2^t$  asks for each path in the network that runs on fiber-optic cable, to return the specific device from which the path originates, and also requires to provide the number of different channels available in the path.
- Query  $Q_3^t$  asks, for each cable in the telecommunications network, the port to which the cable is attached, the slot on the device in which the port is installed, and, for each such slot, its status and its type. For each cable, it also returns its status.

For each OBDA instance  $\langle\langle\mathcal{T}, \mathcal{M}, \mathcal{S}\rangle, \mathcal{D}\rangle$ , we have evaluated the number of query answers and the query answering time with respect to five different setups:

- (1) The default behavior of ONTOP v1.15, which simply ignores all non-*DL-Lite $\mathcal{R}$*  axioms in  $\mathcal{T}$ , i.e., using  $\langle\mathcal{T}^1, \mathcal{M}, \mathcal{S}\rangle$  where  $\mathcal{T}^1$  are all the *DL-Lite $\mathcal{R}$*  axioms in  $\mathcal{T}$ .
- (2) The local semantic approximation (LSA) of  $\mathcal{T}$  in *DL-Lite $\mathcal{R}$* , i.e., using  $\langle\mathcal{T}^2, \mathcal{M}, \mathcal{S}\rangle$  where  $\mathcal{T}^2$  is obtained as the union, for each axiom  $\alpha \in \mathcal{T}$ , of the set of *DL-Lite $\mathcal{R}$*  axioms  $\Gamma(\alpha)$  entailed by  $\alpha$  (Console et al. 2014).
- (3) The global semantic approximation (GSA) of  $\mathcal{T}$  in *DL-Lite $\mathcal{R}$* , i.e., using  $\langle\mathcal{T}^3, \mathcal{M}, \mathcal{S}\rangle$  where  $\mathcal{T}^3$  is the *DL-Lite $\mathcal{R}$*  closure of  $\mathcal{T}$  (Pan and Thomas 2007).
- (4) Result of ONTOPROX,  $\langle\text{rew}(\mathcal{T}), \text{cut}_5(\text{comp}(\mathcal{T}, \mathcal{M})), \mathcal{S}\rangle$ .
- (5) CLIPPER over the materialization of the virtual ABox.

In Table 1, we present details of the evaluation for some of the queries for which we obtained significant results. In Table 2, we provide statistics about the ONTOPROX pre-computations. The performed evaluation led to the following findings:

- For the considered set of queries LSA and GSA produce the same answers.
- Compared to the default ONTOP behavior, LSA/GSA produces more answers for 2 queries out of 4 for UOBM.
- ONTOPROX produces more answers than LSA/GSA for 2 queries out of 4 for UOBM, and for all Telecom queries. In particular, note that for  $Q_1^t$  and  $Q_2^t$ , LSA and GSA returned no answers at all.
- For UOBM, ONTOPROX answers are complete, as confirmed by the comparison with the results provided by CLIPPER. We cannot determine completeness for the Telecom queries, because the Telecom database was too large and its materialization in an ABox was not feasible.
- Query answering of ONTOPROX is  $\sim 3$ – $5$  times slower than ONTOP, when the result sets are of comparable size (note that for  $Q_2^t$  the result set is significantly larger).
- The size of the new *DL-Lite $\mathcal{R}$*  OBDA specifications is comparable with that of the original specifications.

## 7 Conclusions

We proposed a novel framework for rewriting and approximation of OBDA specifications in an expressive ontology language to specifications in a weaker language, in which the core idea is to exploit the mapping layer to encode part of the semantics of the original OBDA specification, and we developed techniques for *DL-Lite $\mathcal{R}$*  as the target language.

We plan to continue our work along the following directions: (i) extend our technique to Horn-*SHIQ*, and, more generally, to Datalog rewritable TBoxes (Cuenca Grau et al. 2013); (ii) deepen our understanding of the computational complexity of deciding CQ-rewritability of OBDA specifications into *DL-Lite $\mathcal{R}$* ; (iii) extend our technique to SPARQL queries under different OWL entailment regimes (Kontchakov et al. 2014); (iv) carry out more extensive experiments, considering queries that contain existentially quantified variables. This will allow us to verify the effectiveness of RewObda, which was designed specifically to deal with existentially implied objects.

**Acknowledgement.** This paper is supported by the EU under the large-scale integrating project (IP) Optique (*Scalable End-user Access to Big Data*), grant agreement n. FP7-318338. We thank Martin Rezk for insightful discussions, and Benjamin Cogrel and Elem Güzel for help with the experimentation.

## References

- Artale, A.; Calvanese, D.; Kontchakov, R.; and Zakharyashev, M. 2009. The *DL-Lite* family and relations. *J. of Artificial Intelligence Research* 36:1–69.
- Bienvenu, M., and Rosati, R. 2015. Query-based comparison of OBDA specifications. In *Proc. of the 28th Int. Workshop on Description Logic (DL)*, volume 1350 of *CEUR Electronic Workshop Proceedings*.
- Bienvenu, M.; Lutz, C.; and Wolter, F. 2013. First-order rewritability of atomic queries in Horn description logics. In *Proc. of the 23rd Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 754–760.
- Botoeva, E.; Kontchakov, R.; Ryzhikov, V.; Wolter, F.; and Zakharyashev, M. 2014. Query inseparability for description logic knowledge bases. In *Proc. of the 14th Int. Conf. on the Principles of Knowledge Representation and Reasoning (KR)*, 238–247. AAAI Press.
- Botoeva, E.; Calvanese, D.; Santarelli, V.; Savo, D. F.; Solimando, A.; and Xiao, G. 2015. Beyond OWL 2 QL in OBDA: Rewritings and approximations (Extended version). CoRR Technical Report abs/1511.08412, arXiv.org e-Print archive. Available at <http://arxiv.org/abs/1511.08412>.
- Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; and Rosati, R. 2007. Tractable reasoning and efficient query answering in description logics: The *DL-Lite* family. *J. of Automated Reasoning* 39(3):385–429.
- Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; Poggi, A.; Rodriguez-Muro, M.; and Rosati, R. 2009. Ontologies and databases: The *DL-Lite* approach. In *5th Reasoning Web Int. Summer School Tutorial Lectures (RW)*, volume 5689 of *LNCS*. Springer. 255–356.
- Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; and Rosati, R. 2013. Data complexity of query answering in description logics. *Artificial Intelligence* 195:335–360.
- Console, M.; Mora, J.; Rosati, R.; Santarelli, V.; and Savo, D. F. 2014. Effective computation of maximal sound approximations of description logic ontologies. In *Proc. of the 13th Int. Semantic Web Conf. (ISWC)*, volume 8797 of *LNCS*, 164–179. Springer.
- Cosmadakis, S. S.; Gaifman, H.; Kanellakis, P. C.; and Vardi, M. Y. 1988. Decidable optimization problems for database logic programs. In *Proc. of the 20th ACM SIGACT Symp. on Theory of Computing (STOC)*, 477–490.
- Cuenca Grau, B.; Motik, B.; Stoilos, G.; and Horrocks, I. 2013. Computing datalog rewritings beyond Horn ontologies. In *Proc. of the 23rd Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 832–838.
- Di Pinto, F.; Lembo, D.; Lenzerini, M.; Mancini, R.; Poggi, A.; Rosati, R.; Ruzzi, M.; and Savo, D. F. 2013. Optimizing query rewriting in ontology-based data access. In *Proc. of the 16th Int. Conf. on Extending Database Technology (EDBT)*, 561–572. ACM Press.
- Eiter, T.; Ortiz, M.; Simkus, M.; Tran, T.-K.; and Xiao, G. 2012. Query rewriting for Horn-SHIQ plus rules. In *Proc. of the 26th AAAI Conf. on Artificial Intelligence (AAAI)*, 726–733. AAAI Press.
- Gaifman, H.; Mairson, H. G.; Sagiv, Y.; and Vardi, M. Y. 1987. Undecidable optimization problems for database logic programs. In *Proc. of the 2nd IEEE Symp. on Logic in Computer Science (LICS)*, 106–115.
- Giese, M.; Soyulu, A.; Vega-Gorgojo, G.; Waaler, A.; Haase, P.; Jiménez-Ruiz, E.; Lanti, D.; Rezk, M.; Xiao, G.; Özçep, Ö. L.; and Rosati, R. 2015. Optique: Zooming in on Big Data. *IEEE Computer* 48(3):60–67.
- Hustadt, U.; Motik, B.; and Sattler, U. 2005. Data complexity of reasoning in very expressive description logics. In *Proc. of the 19th Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 466–471.
- Kazakov, Y. 2009. Consequence-driven reasoning for Horn-SHIQ ontologies. In *Proc. of the 21st Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 2040–2045.
- Kontchakov, R.; Rezk, M.; Rodriguez-Muro, M.; Xiao, G.; and Zakharyashev, M. 2014. Answering SPARQL queries over databases under OWL 2 QL entailment regime. In *Proc. of International Semantic Web Conference (ISWC 2014)*, Lecture Notes in Computer Science. Springer.
- Lutz, C., and Wolter, F. 2011. Non-uniform data complexity of query answering in description logics. In *Proc. of the 24th Int. Workshop on Description Logic (DL)*, volume 745 of *CEUR Electronic Workshop Proceedings*.
- Lutz, C.; Piro, R.; and Wolter, F. 2011. Description logic TBoxes: Model-theoretic characterizations and rewritability. In *Proc. of the 22nd Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 983–988.
- Lutz, C.; Walther, D.; and Wolter, F. 2007. Conservative extensions in expressive description logics. In *Proc. of the 20th Int. Joint Conf. on Artificial Intelligence (IJCAI)*, 453–458.
- Ma, L.; Yang, Y.; Qiu, Z.; Xie, G.; Pan, Y.; and Liu, S. 2006. Towards a complete OWL ontology benchmark. In *Proc. of the 3rd European Semantic Web Conf. (ESWC)*, volume 4011 of *LNCS*, 125–139. Springer.
- Motik, B.; Fokoue, A.; Horrocks, I.; Wu, Z.; Lutz, C.; and Cuenca Grau, B. 2009. OWL Web Ontology Language profiles. W3C Recommendation, World Wide Web Consortium. Available at <http://www.w3.org/TR/owl-profiles/>.
- Pan, J. Z., and Thomas, E. 2007. Approximating OWL-DL ontologies. In *Proc. of the 21st AAAI Conf. on Artificial Intelligence (AAAI)*, 1434–1439.
- Ren, Y.; Pan, J. Z.; and Zhao, Y. 2010. Soundness preserving approximation for TBox reasoning. In *Proc. of the 24th AAAI Conf. on Artificial Intelligence (AAAI)*.
- Rodriguez-Muro, M.; Kontchakov, R.; and Zakharyashev, M. 2013. Ontology-based data access: Ontop of databases. In *Proc. of the 12th Int. Semantic Web Conf. (ISWC)*, volume 8218 of *LNCS*, 558–573. Springer.
- Trivela, D.; Stoilos, G.; Chortaras, A.; and Stamou, G. 2015. Optimising resolution-based rewriting algorithms for OWL ontologies. *J. of Web Semantics* 30–49.