

# Monotone Abstractions in Ontology-Based Data Management

Gianluca Cima, Marco Console, Maurizio Lenzerini, Antonella Poggi

# ► To cite this version:

Gianluca Cima, Marco Console, Maurizio Lenzerini, Antonella Poggi. Monotone Abstractions in Ontology-Based Data Management. The Thirty-Sixth AAAI Conference on Artificial Intelligence, Feb 2022, Virtual Event, Canada. hal-03686667

# HAL Id: hal-03686667 https://hal.science/hal-03686667

Submitted on 2 Jun2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Monotone Abstractions in Ontology-Based Data Management

Gianluca Cima<sup>1</sup>, Marco Console<sup>2</sup>, Maurizio Lenzerini<sup>2</sup>, Antonella Poggi<sup>2</sup>

<sup>1</sup>CNRS & University of Bordeaux <sup>2</sup>Sapienza University of Rome gianluca.cima@u-bordeaux.fr, {console, lenzerini, poggi}@diag.uniroma1.it

#### Abstract

In Ontology-Based Data Management (OBDM), an abstraction of a source query q is a query over the ontology capturing the semantics of q in terms of the concepts and the relations available in the ontology. Since a perfect characterization of a source query may not exist, the notions of best sound and complete approximations of an abstraction have been introduced and studied in the typical OBDM context, i.e., in the case where the ontology is expressed in DL-Lite, and source queries are expressed as unions of conjunctive queries (UCQs). Interestingly, if we restrict our attention to abstractions expressed as UCQs, even best approximations of abstractions are not guaranteed to exist. Thus, a natural question to ask is whether such limitations affect even larger classes of queries. In this paper, we answer this fundamental question for an essential class of queries, namely the class of monotone queries. We define a monotone query language based on disjunctive Datalog enriched with an epistemic operator, and show that its expressive power suffices for expressing the best approximations of monotone abstractions of UCQs.

## 1 Introduction

In Ontology-Based Data Management (OBDM) (Poggi et al. 2008), an ontology, i.e., a formal, logic-based representation of a domain of interest, is used to provide a high-level conceptual tool for accessing and managing the data sources of an information system. Suitable mappings declaratively specify the relationship between the data at the sources and the elements in the ontology, and this enables the user to carry out many relevant tasks on data through the lens of the ontology (Lenzerini 2018; De Giacomo et al. 2018).

In the last decade, there have been extensive investigations on the problem of answering queries expressed over the ontology by means of algorithms that, taking into account both the ontology and the mapping, produce suitable queries that extract the correct data from the sources. The problem is complicated by the need of reasoning about all components of the system, so as not to miss any answer logically implied by the axioms expressing both the ontology and the mapping (see the surveys in (Bienvenu 2016; Xiao et al. 2018; Ortiz 2018)).

Recent papers (Cima 2017; Lutz, Marti, and Sabellek 2018; Cima, Lenzerini, and Poggi 2019, 2020b; Cima et al. 2021) address a different issue in OBDM: starting from a query  $q_S$  expressed over the sources, the goal is to find a so-called *abstraction* of  $q_S$  (Cima 2022), i.e., an ontologybased characterization of  $q_S$  expressed in terms of the ontology elements, and whose answers coincide with the answer to the original query, modulo the ontology and the mapping. We encountered the need of abstraction during a joint project with a public statistical research institute. The institute's departments must publish subsets of the data they gather in the form of semantically described linked open data. To compute the content of the datasets the departments execute suitable queries over the data sources mapped to a shared ontology. Notably, when the dataset is published, it must be documented through a SPARQL query expressed in terms of the ontology. This task is currently done manually. The notion of abstraction perfectly captures this scenario and provides the formal tool for automating the process: given the query over the sources computing the content of the dataset, the abstraction of such query with respect to the mapping and the ontology is exactly the SPARQL query to be associated to the open dataset. Besides the above use case, abstraction can be the appropriate tool for checking the quality of mapping in capturing relevant source queries (Lutz, Marti, and Sabellek 2018), or deriving ontology-mediated queries from examples (Ortiz 2019; Cima, Croce, and Lenzerini 2021).

The first investigations on abstraction appear in (Cima 2017; Lutz, Marti, and Sabellek 2018). Both papers point out that the "perfect" abstraction of a union of conjunctive queries (UCQ) expressed over the data source not always exists, and present algorithms for computing such abstraction in the case where it both exists, and can be expressed as a UCQ over the ontology. In (Cima 2017; Cima, Lenzerini, and Poggi 2019) the notion of (sound and complete) approximations of the perfect abstraction is introduced, exactly to cope with situations in which perfectness cannot be achieved. Moreover, both papers make it clear that, for a given class of queries C, one is probably interested in two specific forms of approximations, called C-minimally complete and C-maximally sound abstractions. Based on these notions, (Cima, Lenzerini, and Poggi 2019) presents a thorough analysis of the verification problem (check whether a given query is a complete or sound abstraction) and the com-

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

putation problem of UCQ-minimally complete and UCQmaximally sound abstractions of UCQ source queries in OBDM systems based on *DL-Lite*. In (Cima, Lenzerini, and Poggi 2020b) the computation problem is studied in the context of a specific class of non-monotone queries for expressing abstractions, and it is shown that this class can provide abstractions that are better than the one in the UCQ class.

Thus, with the exception of (Cima, Lenzerini, and Poggi 2020b), all the results on abstractions have been obtained under the assumptions that abstractions are expressed as UCQs over the ontology, and many of them originate from the observations that best approximations in the UCQ class are not guaranteed to exist. Thus, a natural issue to investigate is whether such limitations affect even larger classes of queries for expressing abstractions. The main goal of this paper is to address the following question: do approximations of perfect abstraction that are best in a fundamental class of queries, namely the class of monotone queries, always exist? Obviously, a related goal is to derive algorithms for computing approximations of abstractions that are best in the class of monotone queries, if they exist. The class of monotone queries includes queries expressible in First-Order Logic (FOL) and is therefore extremely important.

In this paper we answer positively to the above-mentioned fundamental question. More specifically, we provide the following contributions.

- We present a general framework for abstraction in OBDM, based on the definition of queries as functions from the logical models of OBDM systems to sets of tuples. This is important for the goal of making our investigation, and in particular the notion of monotone queries, to be independent from any particular syntactic form. The usual concept of queries whose semantics is based on certain answers is a special case of our general notion.
- 2. The framework includes a new monotone query language, called  $Datalog_{\mathbf{K}}^{\vee}$ , based on disjunctive Datalog enriched with inequalities and an epistemic operator.
- 3. We consider a scenario where the OBDM specification J is based on *DL-Lite<sub>RDFS</sub>* (i.e., the fragment of RDFS expressible in Description Logic), and present algorithms that, given a source query  $q_S$  expressed as UCQ, compute the best (sound or complete) approximations of the J-abstraction of  $q_S$  in the class of monotone queries, expressed in Datalog $_K^{\vee}$ .
- 4. The above algorithms provide the proof that, in the considered scenario, for any UCQ  $q_S$ , the best (sound or complete) approximations of the *J*-abstraction of  $q_S$  in the class of monotone queries always exists. As a consequence, if the perfect abstraction exists and is in the class of monotone queries, then it can be expressed in Datalog<sup>V</sup><sub>K</sub>.

The paper is organized as follows. In Section 2 we recall some preliminary notions. In Section 3 we describe the formal framework for abstraction, including  $Datalog_{\mathbf{K}}^{\vee}$ . Starting from Section 4, we focus on the above-mentioned scenario, and present the algorithms for computing the best complete (Section 4) and the best sound (Section 5) abstractions of source UCQs. In Section 6 we discuss perfect mono-

tone abstractions. Finally, Section 7 concludes the paper by mentioning future works in the context of abstraction.

#### 2 Preliminaries

A relational schema R is a finite set of predicate symbols (each one with a fixed arity) and an R-database D is a finite set of facts using predicate symbols in R and constants from a countable infinite set Const. A query q of arity n over a schema R is a function mapping each R-database to n-tuples of constants occurring in Const. We denote by  $q^D$  the set of tuples obtained by applying q to the R-database D. For two queries  $q_1$  and  $q_2$  over a schema S, we write  $q_1 \sqsubseteq_R q_2$  if  $q_1^D \subseteq q_2^D$  for each R-database D.

Let Var be a countable infinite set of symbols, called variables, disjoint from Const. A relational atom over R is an expression constituted by a predicate symbol in R applied to as many arguments (constants or variables) as its arity. We also consider inequality atoms, where the predicate symbol is inequality  $(\neq)$ . For a set of atoms C, we let  $Dom(C) = Var(C) \cup Const(C)$ , where Var(C) and Const(C) are the set of all variables and all constants, respectively, occurring in C.

A conjunctive query (CQ) q over a schema R is a query expressible as  $\{\bar{x} \mid \exists \bar{y}.\phi(\bar{x},\bar{y})\}$ , where  $\bar{x}$  is a tuple of variables, called the *distinguished variables* of q,  $\bar{y}$  is a tuple of variables, called the *existential variables* of q, and  $\phi(\bar{x},\bar{y})$ is the *body* of q, i.e., a finite conjunction of relational atoms over R containing all the distinguished and existential variables of q. The arity of q is the arity of  $\bar{x}$ . A *union of conjunctive queries (UCQ)* is a finite union of CQs with same arity, called its *disjuncts*. A (U)CQ with inequality ((U)CQ<sup> $\neq$ </sup>) is simply a (U)CQ allowing inequality atoms in its body.

Let  $C_1$  and  $C_2$  be two sets of atoms. A homomorphism h from  $C_1$  to  $C_2$  (written  $h : C_1 \to C_2$ ) is a mapping from  $Dom(C_1)$  to  $Dom(C_2)$  such that: (i)  $h(\alpha) \in C_2$ , for each relational atom  $\alpha \in C_1$ ; (ii) h(c) = c, for each constant  $c \in Const(C_1)$ ; and (iii)  $h(x) \neq h(y)$ , for each inequality atom  $x \neq y$  occurring in  $C_1$ .

Given a signature of unary and binary predicates, a *Description Logic (DL)* ontology is simply a TBox (set of axioms) expressed in a specific DL (Baader et al. 2003). We assume a special unary predicate  $\top$  whose extension in every interpretation coincides with the entire domain. In the technical sections, we are interested in ontologies expressed in *DL-Lite<sub>RDFS</sub>* (Cuenca Grau 2004; Rosati 2007; Cima, Lenzerini, and Poggi 2020a). Such language corresponds to the fragment of RDFS (Brickley and Guha 2014) expressible in *DL-Lite*. Specifically, a *DL-Lite<sub>RDFS</sub>* ontology is a finite set of assertions of the form  $B \sqsubseteq A$  or  $R_1 \sqsubseteq R_2$ , where *B* denotes a basic concept, i.e., an expression of the form  $A, \exists P$ , or  $\exists P^-$ , with A and P denoting an atomic concept and an atomic role, respectively, and  $R_1$  and  $R_2$  denote basic roles, i.e., expressions of the form P, or  $P^-$ .

An *OBDM specification J* is a triple  $\langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ , where  $\mathcal{O}$  is a DL ontology,  $\mathcal{S}$  is a relational schema, called the source schema of *J*, and  $\mathcal{M}$  is a *Global-and-Local-As-View (GLAV)* mapping, i.e., a set of GLAV mapping assertions of the form:

$$\forall \bar{x}. \forall \bar{y}. \phi(\bar{x}, \bar{y}) \to \exists \bar{z}. \psi(\bar{x}, \bar{z}),$$

where  $\bar{x}$  are called the *frontier variables* of m,  $\phi(\bar{x}, \bar{y})$ , called the *body* of m, is a conjunction of relational atoms over S, and  $\psi(\bar{x}, \bar{z})$ , called the *head* of m, is a conjunction of relational atoms over  $\mathcal{O}$ .

An *OBDM system* is a pair  $\Sigma = \langle J, D \rangle$ , where  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$  is an OBDM specification, and D is an *S*-database. For a mapping  $\mathcal{M}$  relating a source schema  $\mathcal{S}$  to an ontology  $\mathcal{O}$  and an *S*-database D, we denote by  $\mathcal{M}(D)$  the set of atoms over  $\mathcal{O}$  obtained by *chasing* D with respect to  $\mathcal{M}$  as described in (Fagin et al. 2005).

The semantics of an OBDM system is given in terms of interpretations for  $\mathcal{O}$  following the so-called *standard names assumption* (Calvanese et al. 2007a). This means that the domain of every interpretation is *Const*, with every constant in the signature interpreted as itself. Note that this also implies the *unique name assumption*. We say that an interpretation  $\mathcal{I}$  for  $\mathcal{O}$  is a model of an OBDM system  $\Sigma = \langle \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle, D \rangle$  if (*i*)  $\mathcal{I} \models \mathcal{O}$ , and (*ii*)  $(D, \mathcal{I}) \models \mathcal{M}$ , where  $(D, \mathcal{I})$  is the interpretation obtained from  $\mathcal{I}$  by adding the set of facts corresponding to D. We denote by  $mod(\langle J, D \rangle)$  the set of models of an OBDM system  $\langle J, D \rangle$ , and when  $mod(\langle J, D \rangle) \neq \emptyset$  we say that D is *consistent* with J. We next provide an example of OBDM specification.

**Example 1.** Let  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$  be an OBDM specification such that  $\mathcal{O} = \emptyset$ ,  $\mathcal{S} = \{s_1/2, s_2/1, s_3/1\}$  and  $\mathcal{M} = \{m_1, m_2, m_3, m_4\}$ , where:

$m_1$ :	$\forall x_1, x_2.s_1(x_1, x_2)$	$\rightarrow$	${\sf E}(x_1, x_2)$
$m_2$ :	$\forall x.s_1(x,x)$	$\rightarrow$	SN(x)
<b>m</b> 3:	$\forall x.s_2(x)$	$\rightarrow$	$\exists z. E(x, z)$
$m_4$ :	$\forall x.s_3(x)$	$\rightarrow$	E(x, x)

The mapping  $\mathcal{M}$  of J establishes how predicates in the schema S relate to the ontology predicates  $\mathsf{E}$  (which stands for edge) and  $\mathsf{SN}$  (which stands for special node).  $\bigtriangleup$ 

A query q of arity n over the signature of an ontology  $\mathcal{O}$ is a function whose domain is the power set of the set of all the interpretations for the ontology, and whose range is the power set of the set of all *n*-tuples of constants in Const. A query q over  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$  is a query over the signature of the ontology  $\mathcal{O}$  that, for every  $\mathcal{S}$ -database D, associates to  $mod(\langle J, D \rangle)$  a set of tuples of constants. The answers to q with respect to the OBDM system  $\langle J, D \rangle$ , denoted as  $q^{J,D}$ . is obtained by applying the function q to the set of models  $mod(\langle J, D \rangle)$ . Although our framework is general enough to accept any kind of function as query, a very common method for defining  $q^{J,D}$  when q is expressed as a FOL query is by sanctioning that the function q applied to  $mod(\langle J, D \rangle)$  returns the set of *certain answers*. A certain answer is a tuple of constants in Dom(J, D, q) that, once assigned to the distinguished variables of q, satisfies of q in every model in  $mod(\langle J, D \rangle)$ . Here, Dom(J, D, q) is the set of all constants occurring in J, D, and q.

Given a set of constants  $\Lambda \subseteq Const$ , we denote by  $q_{|\Lambda}^{J,D}$ the answers to q w.r.t.  $\langle J, D \rangle$  restricted to the constants of  $\Lambda$ , i.e.,  $q_{|\Lambda}^{J,D} = q^{J,D} \cap \Lambda^n$ .

The class of monotone queries is particularly important in this paper. In the context of an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ , a query q over J is monotone if, for every pairs of S-databases  $D, D', mod(\langle J, D \rangle) \subseteq mod(\langle J, D' \rangle)$ implies  $q_{\Lambda(D,D')}^{J,D'} \subseteq q_{\Lambda(D,D')}^{J,D}$ , where  $\Lambda(D,D') = Dom(J, D, q) \cap Dom(J, D', q)$ . This captures the usual intuition of monotonicity of queries in logic: when the knowledge possessed by the formal system increases (i.e., the set of models decreases), the answers to a query increase. Note that the restriction on the constants occurring in both the *active domains* of D and D' (other than in J and q) is imposed for ensuring that the property of monotonicity looks only for those constants occurring in both databases. We use  $\mathfrak{M}^J$  to denote the class of all monotone queries in the context of J, and when J is understood, we simply use  $\mathfrak{M}$ . Note that, for each OBDM specification J,  $\mathfrak{M}^J$  includes the whole class of FOL queries evaluated under the certain answer semantics, that is arguably the most studied class in both the KR and the DB literature.

For two queries  $q_1$  and  $q_2$  over J, we write  $q_1 \sqsubseteq_J q_2$  if  $q_1^{J,D} \subseteq q_2^{J,D}$  for each S-database D, and write  $q_1 \sqsubset_J q_2$  if  $q_1 \sqsubseteq_J q_2$  and  $q_1^{J,D} \subset q_2^{J,D}$  for at least one S-database D. We say that  $q_1$  and  $q_2$  are J-equivalent if both  $q_1 \sqsubseteq_J q_2$  and  $q_2 \sqsubseteq_J q_1$  hold.

#### **3** Framework

In what follows, we implicitly refer to an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ , and when we denote a query by  $q_{\mathcal{O}}$ (resp.,  $q_{\mathcal{S}}$ ) we mean that the query is a query for J (resp., a source query), i.e., is over the signature of the ontology  $\mathcal{O}$  (resp., the schema  $\mathcal{S}$ ). We follow (Cima, Lenzerini, and Poggi 2019) for the basic definitions related to abstraction.

We say that  $q_{\mathcal{O}}$  is a *perfect J-abstraction* of  $q_{\mathcal{S}}$  if  $q_{\mathcal{O}}^{J,D} = q_{\mathcal{S}}^{D}$ , for each  $\mathcal{S}$ -database D consistent with J.

As the condition for an ontology query to be a perfect abstraction of a source query is a strong one, it might be very well the case that a perfect abstraction of a source query does not exist. It is then reasonable to consider weaker notions, such as sound or complete approximations, of perfectness.

We say that  $q_{\mathcal{O}}$  is a *complete (resp. sound) J-abstraction* of  $q_{\mathcal{S}}$  if  $q_{\mathcal{S}}^D \subseteq q_{\mathcal{O}}^{J,D}$  (resp.  $q_{\mathcal{O}}^{J,D} \subseteq q_{\mathcal{S}}^D$ ), for each  $\mathcal{S}$ -database D consistent with J.

Obviously, we might be interested in complete or sound abstractions that approximate  $q_S$  at best, at least in the context of a specific class of queries. If  $\mathcal{L}_{\mathcal{O}}$  is a class of queries, we say that a query  $q_{\mathcal{O}} \in \mathcal{L}_{\mathcal{O}}$  is an  $\mathcal{L}_{\mathcal{O}}$ -minimally complete (resp.,  $\mathcal{L}_{\mathcal{O}}$ -maximally sound) J-abstraction of  $q_S$  if  $q_{\mathcal{O}}$  is a complete (resp., sound) J-abstraction of  $q_S$  and there is no query  $q'_{\mathcal{O}} \in \mathcal{L}_{\mathcal{O}}$  such that  $q'_{\mathcal{O}}$  is a complete (resp., sound) J-abstraction of  $q_S$  and  $q'_{\mathcal{O}} \sqsubset_J q_{\mathcal{O}}$  (resp.,  $q_{\mathcal{O}} \sqsubset_J q'_{\mathcal{O}}$ ).

We now introduce a new language, called  $Datalog_{\mathbf{K}}^{\vee}$ , for expressing queries over OBDM specifications. The language is based on disjunctive Datalog, and is used in this paper for expressing abstractions. The basic component of a  $Datalog_{\mathbf{K}}^{\vee}$  query is a rule. Assume two disjoint and countably infinite sets of predicates  $\mathcal{E}$  and  $\mathcal{P}$ , called extensional and intensional, respectively. A  $Datalog_{\mathbf{K}}^{\vee}$  rule has one of the following forms:

• The typical form of disjunctive Datalog, i.e.,

$$\gamma(\bar{x}) \to \alpha_1(\bar{x}_1) \lor \ldots \lor \alpha_n(\bar{x}_n) \tag{1}$$

where  $\gamma(\bar{x})$  is a conjunction of relational atoms on the predicates of  $\mathcal{P}$  with  $\bar{x}$  as variables, and for i = 1, ..., n,  $\alpha_i(\bar{x}_i)$  is a single relational atom whose predicate is in  $\mathcal{P}$  such that  $\bar{x}_i \subseteq \bar{x}$ ,

• A new form specified as follows

$$\mathbf{K}(\exists \bar{z}.\phi(\bar{x},\bar{z}) \land \xi(\bar{x})) \to \psi_1(\bar{x}_1) \lor \ldots \lor \psi_n(\bar{x}_n) \quad (2)$$

where  $\phi(\bar{x}, \bar{z})$  is a conjunction of relational atoms on the predicates of  $\mathcal{E}, \xi(\bar{x})$  is a conjunction of inequality atoms involving only variables from  $\bar{x}$  and for  $j = 1, \ldots, n$ ,  $\psi_j(\bar{x}_j) = \exists \bar{y}_j. \gamma_j(\bar{x}_j, \bar{y}_j)$ , where  $\gamma_j(\bar{x}_j, \bar{y}_j)$  is a conjunction of relational atoms on  $\mathcal{P}$  such that  $\bar{x}_j \subseteq \bar{x}$ .

An *n*-ary Datalog<sup>\vee</sup><sub>K</sub> query  $q_{\mathcal{O}}$  over an OBDM specification J is a finite set of Datalog<sup>\vee</sup><sub>K</sub> rules whose extensional predicates coincide with the alphabet of  $\mathcal{O}$ , and whose intensional predicates include a special *n*-ary predicate Ans. The semantics of  $q_{\mathcal{O}}$  is provided relative to an OBDM system. Given an OBDM system  $\langle J, D \rangle$ , an interpretation for  $q_{\mathcal{O}}$  is a pair  $I = (mod(\langle J, D \rangle), f)$ , where f is a first-order interpretation (with domain Const) for the predicates in  $\mathcal{P}$ . As usual, we may also see f as the set of facts  $\{p(\bar{c}) \mid \bar{c} \in p^f\}$ . We now define when I satisfies a Datalog<sup>\vee</sup><sub>K</sub> rule.

- I satisfies a rule of the form (1) if the first-order formula  $\forall \bar{x}.\gamma(\bar{x}) \rightarrow \alpha_1(\bar{x}_1) \lor \ldots \lor \alpha_n(\bar{x}_n)$  is true in f,
- I satisfies a rule of the form (2) if for all tuples c̄ of constants, the fact that c̄ is a certain answers to the FOL query ∃z̄.φ(x̄, z̄) ∧ ξ(x̄) w.r.t. ⟨J, D⟩ implies that f satisfies the formula ∃ȳ<sub>j</sub>.γ<sub>j</sub>(c̄<sub>j</sub>, ȳ<sub>j</sub>), for some j = 1,...,n.

An interpretation I for  $q_{\mathcal{O}}$  is called a *model* of  $q_{\mathcal{O}}$ if all the rules of  $q_{\mathcal{O}}$  are satisfied by I. Finally, we define the notion of answers to an *n*-ary Datalog<sup> $\vee$ </sup><sub>**K**</sub> query  $q_{\mathcal{O}}$  w.r.t. an OBDM system  $\langle J, D \rangle$ , denoted by  $q_{\mathcal{O}}^{J,D}$ , as follows: { $\bar{c} \in Dom(J, D, q_{\mathcal{O}})^n \mid \bar{c} \in Ans^f$  for each model  $(mod(\langle J, D \rangle), f)$  of  $q_{\mathcal{O}}$  }.

**Example 2.** Consider the OBDM specification J illustrated in Example 1. The following  $Datalog_{\mathbf{K}}^{\vee}$  query  $q_{\mathcal{O}}$  over J returns all the pairs  $(v_1, v_2)$  of special nodes that are known to be distinct and such that there is a path from  $v_1$  to  $v_2$ passing only for nodes known to be special:

$$\begin{array}{rcl} \mathbf{K}(\mathsf{E}(x_1, x_2) \land \mathsf{SN}(x_1) \land \mathsf{SN}(x_2)) & \to & T_1(x_1, x_2) \\ \mathbf{K}(\mathsf{SN}(x_1) \land \mathsf{SN}(x_2) \land x_1 \neq x_2) & \to & T_2(x_1, x_2) \\ T_1(x_1, y) \land T_1(y, x_2) & \to & T_1(x_1, x_2) \\ T_1(x_1, x_2) \land T_2(x_1, x_2) & \to & Ans(x_1, x_2) & \bigtriangleup \end{array}$$

The following proposition shows that  $\mathsf{Datalog}_{\mathbf{K}}^{\vee}$  is a monotone query language.

**Proposition 1.** Every  $Datalog_{\mathbf{K}}^{\vee}$  query over J is in  $\mathfrak{M}^{J}$ , for every OBDM specification J.

The semantics should make it clear that **K** is the knowledge operator in the S5 epistemic logic: the formula **K**A should be read as "A is known (i.e., logically implied) by the system" (Levesque and Lakemeyer 2001). Therefore, when accessing the information modeled by  $\langle J, D \rangle$ , a Datalog<sup>V</sup> query extracts what is known by the system, and this characteristic is crucial for not falling into undecidability resulting from using Datalog rules jointly with Description Logics (see (Levy and Rousset 1998; Calvanese and Rosati 2003)), as stated in the following proposition.

**Proposition 2.** Let  $\Sigma$  be an OBDM system. Answering  $Datalog_{\mathbf{K}}^{\vee}$  queries w.r.t.  $\Sigma$  is decidable if and only if answering CQs w.r.t.  $\Sigma$  is decidable<sup>1</sup>.

Although our framework is general enough to consider any DL for expressing ontologies and any query language for expressing source queries, in the rest of this paper we will carry out our investigation in the following setting: (i) ontologies are expressed in *DL-Lite<sub>RDFS</sub>*, and (ii) source queries are expressed as UCQs.

At this point, one may wonder whether  $Datalog_{\mathbf{K}}^{\vee}$  is the right language to express monotone abstractions in this setting. While a thorough analysis of the language is outside the scope of the present paper, the following proposition provides a positive answer to this question, at least from the computational point of view.

**Proposition 3.** In our setting, (i) answering  $Datalog_{\mathbf{K}}^{\vee}$  queries is in coNP in data complexity, and (ii) there exists an OBDM specification J and a CQ  $q_S$  such that, given an S-database D, answering the  $\mathfrak{M}$ -maximally sound J-abstraction of  $q_S$  is coNP-hard in data complexity.

To ease the presentation, from now on we assume that mappings and source queries do not mention constants. However, all our results can be straightforwardly adapted to the case where constants are allowed.

#### **4** Minimally Complete Abstractions

We now prove that  $\mathfrak{M}$ -minimally complete *J*-abstractions of UCQs always exist and can be expressed in Datalog<sup> $\vee$ </sup><sub>K</sub>. Actually, a fragment of Datalog<sup> $\vee$ </sup><sub>K</sub> suffices for this purpose, namely the one using only rules of form (2) without disjunction. We start with some useful subroutines.

Given a CQ  $q_S = \{\bar{x} \mid \exists \bar{y}.\phi(\bar{x},\bar{y})\}$ , the subroutine SaturateQ( $q_S$ ) computes a UCQ<sup>\neq </sup> in the following way: for each possible unifier  $\mu$  on the variables in  $\bar{x} \cup \bar{y}$  such that  $\mu(x) \in \bar{x}$  for each  $x \in \bar{x}$ , SaturateQ( $q_S$ ) contains a query obtained from  $\mu(q_S)$  by adding the inequality atom  $t_1 \neq t_2$  for each pair of distinct variables  $t_1, t_2$  occurring in  $\mu(q_S)$ . For a UCQ  $q_S$ , we denote by SaturateQ( $q_S$ ) the UCQ<sup>\neq </sup> obtained by applying SaturateQ(q) to each disjunct q of  $q_S$ . We write each CQ<sup>\neq </sup> q generated by SaturateQ( $q_S$ ) as  $q = \{\bar{x} \mid \exists \bar{y}.\phi(\bar{x},\bar{y}) \land \xi(\bar{x},\bar{y})\}$ , where  $\phi(\bar{x},\bar{y})$  and  $\xi(\bar{x},\bar{y})$ are the conjunctions of the relational atoms over S and of the inequality atoms, respectively, occurring in the body of q.

Moreover, for an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ and a  $CQ^{\neq} q = \{ \bar{x} \mid \exists \bar{y}.\phi(\bar{x},\bar{y}) \land \xi(\bar{x},\bar{y}) \}$  over  $\mathcal{S}$ , we denote by  $r_q$  the following Datalog<sub>K</sub> rule of form (2):

$$r_q = \mathbf{K}(\exists \bar{z}. \top(\bar{x}) \land \mathcal{M}(q) \land \xi(\bar{x}, \bar{y})) \to Ans(\bar{x}),$$

where (i)  $\top(\bar{x})$  denotes  $\top(x_1) \land \ldots \land \top(x_n)$  for  $\bar{x} = (x_1, \ldots, x_n)$ , (ii)  $\mathcal{M}(q)$  is computed by simply ignoring the inequality atoms and chasing the set of relational atoms occurring in the body of q; (iii)  $\bar{y} \subseteq \bar{y}$  is the subset of the

<sup>&</sup>lt;sup>1</sup>With answering we implicitly refer to the associated *recognition problem*, i.e., check whether a tuple is in the answer to a query.

existential variables of q occurring in  $\mathcal{M}(q)$ ;  $(iv) \bar{z}$  are the fresh variables introduced when computing  $\mathcal{M}(q)$ ; and (v) $\xi(\bar{x}, \bar{y})$  is the conjunction of the inequality atoms obtained from  $\xi(\bar{x}, \bar{y})$  by removing all those atoms of the form  $y \neq t$ and  $t \neq y$  in which y is an existential variable occurring in  $\bar{y}$  but not in  $\bar{y}$  (i.e., not in  $\mathcal{M}(q)$ ) and t is any other possible variable. Observe that the epistemic operator is exploited to bind the existential variables coming from q. This is achieved by pushing the subset  $\bar{y}$  of the existential variables  $\bar{y}$  of q occurring in  $\mathcal{M}(q)$  inside the **K** operator.

We are now ready to present the algorithm  $\mathfrak{M}$ -MinComplete for computing the  $\mathfrak{M}$ -minimally complete J-abstractions. Given an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$  and a UCQ  $q_{\mathcal{S}}$  over  $\mathcal{S}$  such that SaturateQ $(q_{\mathcal{S}}) = q_1 \cup \ldots \cup q_n$ ,  $\mathfrak{M}$ -MinComplete $(J, q_{\mathcal{S}})$  outputs the Datalog<sup>K</sup> query  $q_{\mathcal{O}} = \{r_{q_1}, \ldots, r_{q_n}\}$  over J.

**Example 3.** Consider the OBDM specification J illustrated in Example 1 and the CQ  $q_S = \{(x) \mid \exists y.s_1(x,y)\}$  over S. One can verify that  $\mathfrak{M}$ -MinComplete $(J, q_S)$  returns the following Datalog<sup>\neq</sup> query  $q_O$  over J asking for all those nodes v such that either v is connected to a node v' known to be different from v or v is a special node with a self-loop:

$\mathbf{K}(E(x,y) \land x \neq y)$	$\rightarrow$	Ans(x)
$\mathbf{K}(E(x,x) \wedge SN(x))$	$\rightarrow$	Ans(x)

Note that  $q_{\mathcal{O}}$  is a better complete approximation than the query  $\{(x) \mid \exists y. \mathsf{E}(x, y)\}$ , which is the UCQ-minimally complete J-abstraction of  $q_{\mathcal{S}}$  (Cima, Lenzerini, and Poggi 2019).

**Theorem 1.**  $\mathfrak{M}$ -MinComplete $(J, q_S)$  terminates and returns the unique (up to J-equivalence)  $\mathfrak{M}$ -minimally complete J-abstraction of  $q_S$ .

Furthermore, we observe that the result of  $\mathfrak{M}$ -MinComplete $(J, q_S)$  is independent from the assertions occurring in the ontology of the OBDM specification J. Similar results can be obtained for OBDM specifications based on more expressive Horn DL ontologies.

Before concluding this section, we point out that answering Datalog<sup> $\vee$ </sup> queries returned by  $\mathfrak{M}$ -MinComplete $(J, q_S)$ w.r.t. OBDM systems of our setting is FOL-rewritable, as the following proposition shows.

**Proposition 4.** Let  $q_{\mathcal{O}}$  be a  $\mathsf{Datalog}_{\mathbf{K}}^{\vee}$  with only rules of the form  $\mathbf{K}(\exists \bar{z}.\phi(\bar{x},\bar{z}) \land \xi(\bar{x})) \to Ans(\bar{x}_a)$  over an OBDM specification J. It is possible to compute a  $UCQ^{\neq} q_r$  over S such that  $q_{\mathcal{O}}^{J,D} = q_r^D$ , for each S-database D.

### 5 Maximally Sound Abstractions

We present the algorithm  $\mathfrak{M}$ -MaxSound for the computation of  $\mathfrak{M}$ -maximally sound *J*-abstractions of UCQs. The output of  $\mathfrak{M}$ -MaxSound is a Datalog<sup>V</sup><sub>K</sub> query, thus proving that such abstractions always exist and can be expressed in Datalog<sup>V</sup><sub>K</sub>. To present our algorithm, we first need to introduce the notion of *inverse mappings*, which extends the one introduced in (Duschka and Genesereth 1998) for the problem of rewriting queries using *disjunctive views*<sup>2</sup>.

To ease the presentation of our techniques, we now introduce a technical tool to encode a *DL-Lite<sub>RDFS</sub>* ontology into a set of GLAV mapping assertions. Given a conjunction of atoms  $\psi$  over a *DL-Lite<sub>RDFS</sub>* ontology  $\mathcal{O}$ , we denote by  $\mathcal{O}(\psi)$  the conjunction of atoms over  $\mathcal{O}$  obtained by chasing the atoms in  $\psi$  w.r.t.  $\mathcal{O}$  as described in (Calvanese et al. 2007b). Note that, since O is a *DL-Lite<sub>RDFS</sub>* ontology,  $\mathcal{O}(\psi)$  is always finite and do not introduce further existential variables. Starting from an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ , we can generate an equivalent OBDM specification  $J' = \langle \emptyset, \mathcal{S}, \mathcal{M}' \rangle$  as follows. For each mapping assertion  $m \in \mathcal{M}$  of the form  $\forall \bar{x}.\forall \bar{y}.\phi(\bar{x},\bar{y}) \rightarrow$  $\exists \bar{z}.\psi(\bar{x},\bar{z})$ , the mapping  $\mathcal{M}'$  contains the mapping assertion  $\forall \bar{x}.\forall \bar{y}.\phi(\bar{x},\bar{y}) \rightarrow \exists \bar{z}.\mathcal{O}(\psi(\bar{x},\bar{z}))$ . From now on, without loss of generality, we assume OBDM specifications with an empty set of ontological assertions. Furthermore, in each mapping  $\mathcal{M}$  relating a source schema  $\mathcal{S}$  to an ontology  $\mathcal{O}$ , again without loss of generality we assume that  $\mathcal{M}$  contains the assertion  $\forall x_1, \ldots, x_n.s(x_1, x_2, \ldots, x_n) \rightarrow$  $\top(x_1) \land \top(x_2) \land \ldots \land \top(x_n)$  for each *n*-ary predicate  $s \in S$ .

#### 5.1 Inverse Mapping

Our algorithm for the computation of  $\mathfrak{M}$ -maximally sound abstractions is based on the notion of inverse mappings. Intuitively, given a GLAV mapping  $\mathcal{M}$ , the inverse of  $m \in \mathcal{M}$ is a Datalog<sup> $\vee$ </sup> rule that characterizes those that may give rise to an homomorphic image of the head of m when chased with  $\mathcal{M}$ . To compute inverse mappings, we will define the algorithm lnvMap. In order to define such algorithm, we need to introduce several additional notions.

**Splitting.** Given a set of relational atoms A and a set of variables U occurring in A, we define the U-graph  $G_A^U$  of A as follows: (i) each atom  $\alpha$  in A is a node of  $G_A^U$ , and (ii) there is an edge between nodes  $\alpha_1, \alpha_2 \in A$  if and only if there exists a variable  $z \in U$  such that z occurs in both  $\alpha_1$  and  $\alpha_2$ . The *splitting of* A around U (written Split(A, U)) is the set of all those  $C \subseteq A$  such that the atoms in C form a connected component of  $G_A^U$ , i.e., C is a maximal subset of A s.t. there exists a path from  $\alpha$  to  $\alpha'$  in  $G_A^U$ , for each  $\alpha, \alpha' \in C$ . Whenever necessary, we will assume that sets in Split(A, U) are presented as a tuple.

Saturate. Given a GLAV mapping  $\mathcal{M}$ , the algorithm SaturateM( $\mathcal{M}$ ) produces an equivalent set of rules with explicit inequalities on the frontier variables. Given a GLAV mapping assertion  $m = \forall \bar{x} . \forall \bar{y} . \phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} . \psi(\bar{x}, \bar{z})$ , the algorithm Saturate M(m) computes a set of rules in the following way: for each possible unifier  $\mu$  between the variables in  $\bar{x}$ , SaturateM(m) contains a rule obtained from  $\mu(m)$  by adding the inequality atom  $x_1 \neq x_2$  in the body of  $\mu(m)$ , for each pair of distinct variables  $x_1, x_2 \in \mu(\bar{x})$ . For a GLAV mapping  $\mathcal{M}$ , we denote by SaturateM( $\mathcal{M}$ ) the set of rules obtained by applying SaturateM( $\mathcal{M}$ ) to each mapping assertion m of  $\mathcal{M}$ . It is easy to see that  $(D, \mathcal{I}) \models \mathcal{M}$  if and only if  $(D, \mathcal{I}) \models \mathsf{SaturateM}(\mathcal{M})$ . We call each assertion in SaturateM( $\mathcal{M}$ ) a saturated GLAV mapping assertion and SaturateM( $\mathcal{M}$ ) a saturated GLAV mapping. We write each rule  $m \in \mathsf{SaturateM}(\mathcal{M})$  as  $\phi(\bar{x}, \bar{y}) \land \xi(\bar{x}) \to \exists \bar{z}. \psi(\bar{x}, \bar{z}),$ where  $\phi(\bar{x}, \bar{y})$  and  $\psi(\bar{x}, \bar{z})$  are conjunctions of relational

<sup>&</sup>lt;sup>2</sup>We refer to (Cima et al. 2021) for the relation between abstraction and the problem of rewriting queries using disjunctive views.

atoms over S and over O, respectively, whereas  $\xi(\bar{x})$  is a conjunction of inequality atoms. We write r-body(m) for the set of relational atoms in  $\phi(\bar{x}, \bar{y})$ , and, as in GLAV mappings, head(m) for the set of relational atoms in  $\psi(\bar{x}, \bar{z})$ .

**Generators and Supports.** Assume a saturated GLAV mapping  $\mathcal{M}$  from a source schema  $\mathcal{S}$  to an ontology  $\mathcal{O}$ . Without loss of generality, we assume that mapping assertions in  $\mathcal{M}$  use distinct variables. Given a set of atoms A over  $\mathcal{O}$ , a generator of A w.r.t.  $\mathcal{M}$  is a pair  $g = \langle h, m \rangle$ , where h is a homomorphism from A to the head of m and  $m \in \mathcal{M}$ . A generator g is V-distinguished, for some  $V \subseteq Var(A)$ , if h(v) is a frontier variable of m, for each  $v \in V$ , and h(u) is not a frontier variable of m, for each  $u \notin V$ . We write  $Gens(A, \mathcal{M}, V)$  for the set of V-distinguished generators of A w.r.t.  $\mathcal{M}$  Intuitively, a generator  $g = \langle h, m \rangle$  of A w.r.t.  $\mathcal{M}$  represents a possible way to obtain a homomorphic copy  $B \subseteq \mathcal{M}(D)$  of A by applying m to an S-database D. If such generator is V-distinguished, then the only constants occurring in B are images of the variables in V.

**Example 4.** Let  $\mathcal{M} = \{m_A, m_B, m_C, m_D\}$ , where:

$$\begin{split} \mathbf{m}_{A} &: s_{1}(x_{1}, x_{2}) \land x_{1} \neq x_{2} \to \exists Z. P_{1}(x_{1}, Z) \land P_{2}(Z, x_{2}), \\ \mathbf{m}_{B} &: s(x_{3}, x_{4}) \land x_{3} \neq x_{4} \to \exists Z'. P_{1}(x_{3}, Z') \land P_{2}(Z', x_{4}), \\ \mathbf{m}_{C} &: s_{3}(x_{5}, y_{1}) \land s_{4}(y_{1}, y_{2}) \to P_{1}(x_{5}, x_{5}), \\ \mathbf{m}_{D} &: s_{5}(x_{6}, x_{7}) \land x_{6} \neq x_{7} \to P_{1}(x_{6}, x_{7}) \land P_{2}(x_{7}, x_{6}). \end{split}$$

For the set of atoms  $A_1 = \{P_1(x_1, Z), P_2(Z, x_2)\}$ , we have that  $Gens(A_1, \mathcal{M}, \{x_1, x_2\}) = \{g_1, g_2\}$  with  $g_1 = \langle h_1, \mathbf{m}_A \rangle$  and  $g_2 = \langle h_2, \mathbf{m}_B \rangle$  where  $h_1$  is the identity function, whereas  $h_2(x_1) = x_3$ ,  $h_2(x_2) = x_4$ , and  $h_2(Z) = Z'$ .

For the set of atoms  $A_2 = \{P_1(x_1, Z)\}$ , we have that  $Gens(A_2, \mathcal{M}, \{x_1, Z\}) = \{g_3, g_4\}$  with  $g_3 = \langle h_3, \mathbf{m}_C \rangle$ and  $g_4 = \langle h_4, \mathbf{m}_D \rangle$  where  $h_3(x_1) = x_5$  and  $h_3(Z) = x_5$ , whereas  $h_4(x_1) = x_6$  and  $h_4(Z) = x_7$ .

For the set of atoms  $A_3 = \{P_2(Z, x_2)\}$ , we have that  $Gens(A_3, \mathcal{M}, \{Z, x_2\}) = \{g_5\}$  with  $g_5 = \langle h_5, \mathbf{m}_D \rangle$  where  $h_5(Z) = x_7$  and  $h_5(x_2) = x_6$ .

We now extend the notion of generators to the case of multiple applications of mapping assertions. Given a tuple of sets of atoms  $\mathcal{A} = \langle A_1, \ldots, A_n \rangle$  and a set of variables  $V \subseteq Var(A_1) \cup \ldots \cup Var(A_n)$ , the V-distinguished generators of  $\mathcal{A}$  w.r.t.  $\mathcal{M}$  are all the tuples  $G = \langle g_1, \ldots, g_n \rangle$  such that  $g_i \in Gens(A_i, \mathcal{M}, V)$ , for each  $i = 1, \ldots, n$ . With a slight abuse of notation, we use  $Gens(\mathcal{A}, \mathcal{M}, V)$  to denote the set of all tuples of V-distinguished generators of elements of  $\mathcal{A}$  w.r.t.  $\mathcal{M}$  and we call generator of  $\mathcal{A}$  every such tuple. Intuitively, each  $G \in Gens(\mathcal{A}, \mathcal{M}, V)$  represents a possible way to obtain a homomorphic copy of  $\mathcal{A}$  using a single assertion of  $\mathcal{M}$  for each  $A \in \mathcal{A}$ .

**Example 5.** Refer to Example 4. For  $\mathcal{A} = \langle A_1 \rangle$  and  $\mathcal{A}' = \langle A_2, A_3 \rangle$ , we have  $Gens(\mathcal{A}, \mathcal{M}, \{x_1, x_2\}) = \{G_1, G_2\}$  and  $Gens(\mathcal{A}', \mathcal{M}, \{x_1, x_2, Z\}) = \{G_3, G_4\}$ , resp., where  $G_1 = \langle g_1 \rangle$ ,  $G_2 = \langle g_2 \rangle$ ,  $G_3 = \langle g_3, g_5 \rangle$ , and  $G_4 = \langle g_4, g_5 \rangle$ .

We now characterize those databases D for which  $\mathcal{M}(D)$ contains an homomorphic copy of  $\mathcal{A}$  generated using the the mappings constituting some  $G \in Gens(\mathcal{A}, \mathcal{M}, V)$ . To this end, we define a formula  $\sigma_G$  that we call the *support* of  $G \in Gens(\mathcal{A}, \mathcal{M}, V)$ . Given  $G = \langle g_1, \ldots, g_n \rangle$  in  $Gens(\mathcal{A}, \mathcal{M}, V)$  with  $g_i = \langle h_i, m_i \rangle$  for each i = 1, ..., n, let  $\langle \rho_1, ..., \rho_n \rangle$  be a tuple of renaming for the variables of  $\mathcal{M}$  whose images are pairwise disjoint. For i = 1, ..., n, we define  $\beta_i = \rho_i(r \cdot body(m_i))$  and  $\sigma_G = \bigwedge_i \rho_i(r \cdot body(m_i))$ . Intuitively, each  $\beta_i$  represents a set of relational source atoms that is needed to generate an homomorphic image of  $h_i(A_i)$  using  $m_i$ .

**Example 6.** Refer to Example 5. We have  $\sigma_{G_1} = s_1(x'_1, x'_2)$ ,  $\sigma_{G_2} = s(x'_3, x'_4)$ ,  $\sigma_{G_3} = s_3(x'_5, y'_1) \land s_4(y'_1, y'_2) \land s_5(x''_6, x''_7)$ , and  $\sigma_{G_4} = s_5(x'_6, x'_7) \land s_5(x''_6, x''_7)$ .

In order to connect together the different  $h_i(A_i)$  to form an image of  $\mathcal{A}$ , we need to impose further restrictions on the variables of  $\sigma_G$ . These restrictions will take the form of equalities. Let  $\eta_G$  denote binary relation over  $Var(\sigma_G)$ defined as follows: for each  $i, j = 1, \ldots, n$ , and  $v \in V$ ,  $(\rho_i(h_i(v)), \rho_j(h_j(v))) \in \eta_G$ . Moreover, given a set  $\bar{x} \subseteq V$ , let  $\eta^{[\bar{x}]}$  denote the binary relation over  $Var(\sigma_G) \cup \bar{x}$  defined as follows: for each  $i = 1, \ldots, n$ , and  $x \in \bar{x}$ ,  $(x, \rho_i(h_i(x))) \in \eta^{[\bar{x}]}$ . We use  $\eta_G^{[\bar{x}]}$  to denote the reflexive and transitive closure of  $\eta_G \cup \eta^{[\bar{x}]}$ .

**Example 7.** Refer to Example 6. We have that  $\eta_{G_1}^{[\{x_1,x_2\}]}$ ,  $\eta_{G_2}^{[\{x_1,x_2\}]}$ ,  $\eta_{G_3}^{[\{x_1,x_2\}]}$ , and  $\eta_{G_4}^{[\{x_1,x_2\}]}$  are the reflexive and transitive closure of  $\{(x_1, x_1'), (x_2, x_2')\}$ ,  $\{(x_1, x_3'), (x_2, x_4')\}$ ,  $\{(x_5', x_7''), (x_1, x_5'), (x_2, x_6'')\}$ ,  $\{(x_7', x_7''), (x_1, x_6'), (x_2, x_6'')\}$ , respectively.

As customary, for each equivalence class E of  $\eta_G^{[\bar{x}]}$ , we select an element  $s(E) \in E$  as representative; if E contains elements of  $\bar{x}$ , we require  $s(E) \in \bar{x}$ . Moreover, given  $v \in Var(\sigma_G)$ , we use  $v_{\sim G,\bar{x}}$  to denote the equivalence class of  $\eta_G^{[\bar{x}]}$  that contains v. Finally, we define the  $\bar{x}$ -support  $\sigma_G^{\bar{x}}$ of G as the formula  $\exists \bar{w}. \eta_G^{[\bar{x}]}(\sigma_G)$ , where  $\eta_G^{[\bar{x}]}(\sigma_G)$  is obtained from  $\sigma_G$  by replacing each variable  $v \in Var(\sigma_G)$ with  $s(v_{\sim G,\bar{x}})$  and  $\bar{w}$  are the variables of  $\eta_G^{[\bar{x}]}(\sigma_G)$  not occurring in  $\bar{x}$ .

**Example 8.** Refer to the previous example. We have that  $\sigma_{G_1}^{\{x_1,x_2\}} = s_1(x_1,x_2), \ \sigma_{G_2}^{\{x_1,x_2\}} = s(x_1,x_2), \ \sigma_{G_3}^{\{x_1,x_2\}} = \exists y'_1, y'_2.s_3(x_1,y'_1) \land s_4(y'_1,y'_2) \land s_5(x_2,x_1), \ and \ \sigma_{G_4}^{\{x_1,x_2\}} = \exists x'_7.s_5(x_1,x'_7) \land s_5(x_2,x'_7).$ 

The lnvMap Algorithm. We are now ready to formally define the lnvMap algorithm. In the algorithm, we use  $Partitions(\bar{z})$ , where  $\bar{z}$  is a set of variables, to denote the set of all 2-tuples of pairwise disjoint sets  $\langle U, V \rangle$  such that  $U \cup V = \bar{z}$ . Intuitively, given an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ , lnvMap(J) returns a set of Datalog<sup>V</sup><sub>K</sub> rules that represents a *sound inverse of*  $\mathcal{M}$  w.r.t.  $\mathcal{O}$ , i.e., a set of formulae that, given an S-database D, characterize those databases D' such that  $mod(J, D) \supseteq mod(J, D')$ . Equivalently, given an S-database D, lnvMap(J) characterizes all those databases D' such that  $\mathcal{M}(D')$  is an homomorphic copy of  $\mathcal{M}(D)$ .

In order to define a sound inverse of  $\mathcal{M}$  w.r.t.  $\mathcal{O}$ , InvMap(J) relies on the notions of supports and generators. More specifically, for every saturated mapping assertion  $\phi(\bar{x}, \bar{y}) \wedge \xi(\bar{x}) \rightarrow \exists \bar{z}.\psi(\bar{x}, \bar{z})$ , InvMap(J) contains Algorithm 1: InvMap

**Input:** OBDM specification  $\langle \emptyset, S, \mathcal{M} \rangle$ **Output:** Set  $\mathcal{R}$  of Datalog<sup> $\vee$ </sup> rules 1:  $\mathcal{M}' = \mathsf{SaturateM}(\mathcal{M})$ 2:  $\mathcal{R} = \emptyset$ 3: for all  $m: \phi(\bar{x}, \bar{y}) \land \xi(\bar{x}) \to \exists \bar{z}. \psi(\bar{x}, \bar{z}) \in \mathcal{M}'$  do 4:  $\Gamma = \emptyset$ 5: for all  $\langle U, V \rangle \in Partitions(\bar{z})$  do 6: Let  $S_{\psi} = Split(\psi, U)$ for all  $G \in Gens(S_{\psi}, \mathcal{M}', V \cup \bar{x})$  do 7: if  $(x, x') \notin \eta_G^{[\bar{x}]}$  for each  $x, x' \in \bar{x}$  s.t.  $x \neq x'$  then Let  $\sigma_G^{\bar{x}}$  be the  $\bar{x}$ -support of G8: 9:  $\Gamma = \Gamma \cup \{\sigma_G^{\bar{x}}\}$ 10: end if 11: end for 12: 13: end for  $inv(m) = \mathbf{K}(\exists \bar{z}.\psi(\bar{x},\bar{z}) \land \xi(\bar{x})) \rightarrow \bigvee_{\gamma \in \Gamma} \gamma$ 14:  $\mathcal{R} = \mathcal{R} \cup \{inv(m)\}$ 15: 16: end for 17: return  $\mathcal{R}$ 

a Datalog<sup>V</sup><sub>K</sub> rule of the form  $\mathbf{K}(\exists \bar{z}.\psi(\bar{x},\bar{z}) \land \xi(\bar{x})) \rightarrow \bigvee_{\gamma \in \Gamma} \gamma$ , where  $\Gamma$  is the set of all  $\bar{x}$ -supports of the generators of all possible splittings of  $\psi$ . Intuitively, the latter guarantees that, if  $(mod(\langle J, D \rangle), D')$  satisfies  $\operatorname{InvMap}(J)$ , for every atom  $\alpha \in \mathcal{M}(D)$ , then there exists an assertion  $m \in \mathcal{M}$  such that  $h(body(m)) \subseteq D'$  and  $\ell(\alpha) \in h(head(m))$ , where h and  $\ell$  are homomorphisms.

**Example 9.** Consider the OBDM specification J illustrated in Example 1. One can verify that InvMap(J) returns the following set  $\mathcal{R}$  of  $Datalog_{\mathbf{K}}^{\vee}$  rules:

In the reminder of this section, we introduce two crucial properties of the InvMap algorithm. Assume an OBDM specification J and an S-database D. One can show that  $(mod(\langle J, D \rangle), D)$  satisfies InvMap(J).

**Lemma 1.** For every S-database D,  $(mod(\langle J, D \rangle), D) \models$ InvMap(J).

As it will turn out, Lemma 1 guarantees that our approximation is sound. Moreover, one can show that, whenever  $(mod(\langle J, D \rangle), D') \models InvMap(J)$ , we have that  $\mathcal{M}(D)$  can be homomorphically mapped into  $\mathcal{M}(D')$ . Informally, this property guarantees that our approximation is maximal.

**Lemma 2.** For every pair of S-databases D, D' such that  $(mod(\langle J, D \rangle), D') \models InvMap(J)$ , there exists a homomorphism from  $\mathcal{M}(D)$  to  $\mathcal{M}(D')$ .

#### 5.2 Maximally Sound Abstraction Algorithm

With the InvMap algorithm and its properties at hand, we are now ready to focus on the problem of computing the  $\mathfrak{M}$ -maximally sound *J*-abstractions. First, let us introduce

some additional notations. For an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$ , the *standard predicate renaming of* J, denoted by ren<sub>J</sub>, is the renaming such that ren<sub>J</sub>(s) = s', for every  $s \in \mathcal{S}$ , where s' is an IDB predicate with ar(s) = ar(s'). Furthermore, with a slight abuse of notation, given a set of formulae  $\Psi$  over  $\mathcal{S} \cup \mathcal{O}$ , we write ren<sub>J</sub>( $\Psi$ ) for the set of formulae obtained from  $\Psi$  by replacing each occurrence of s in  $\Psi$  with ren<sub>J</sub>(s), for each  $s \in \mathcal{S}$ .

We can now present the algorithm  $\mathfrak{M}$ -MaxSound for computing the  $\mathfrak{M}$ -maximally sound J-abstractions. Given an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$  and a UCQ  $q_{\mathcal{S}} = q_{\mathcal{S}}^1 \cup \ldots \cup q_{\mathcal{S}}^n$  over  $\mathcal{S}$  such that  $q_{\mathcal{S}}^i = \{\bar{x}_i \mid \exists \bar{y}_i.\phi_i(\bar{x}_i, \bar{y}_i)\}$  for each  $i = 1, \ldots, n$ ,  $\mathfrak{M}$ -MaxSound $(J, q_{\mathcal{S}})$  simply outputs the Datalog<sup> $\checkmark$ </sup> query  $q_{\mathcal{O}} = \operatorname{ren}_J(\Psi)$ , where  $\Psi = \operatorname{InvMap}(J) \cup \{\phi_1(\bar{x}_1, \bar{y}_1) \to Ans(\bar{x}_1)\} \cup \ldots \cup \{\phi_n(\bar{x}_n, \bar{y}_n) \to Ans(\bar{x}_n)\}$ .

We conclude this section by establishing termination and correctness of the  $\mathfrak{M}$ -MaxSound algorithm.

**Theorem 2.**  $\mathfrak{M}$ -MaxSound $(J, q_S)$  terminates and returns the unique (up to J-equivalence)  $\mathfrak{M}$ -maximally sound Jabstraction of  $q_S$ .

#### 6 Perfect Abstractions

It follows from Theorem 1 that either the perfect abstraction of a source UCQ can be expressed in Datalog<sup> $\lor$ </sup><sub>K</sub>, or it cannot be expressed as a monotone query (if it exists at all). We now present an algorithm that, given an OBDM specification J and a source UCQ  $q_S$ , returns the perfect J-abstraction of  $q_S$ , if and only if it exists and is in  $\mathfrak{M}$ . To this aim, we make use of Proposition 4, and refer to the  $q_r$  defined in that proposition as the *rewriting of*  $q_O$  w.r.t. J.

Our algorithm, that we call  $\mathfrak{M}$ -Perfect, goes as follows. Given an OBDM specification  $J = \langle \mathcal{O}, \mathcal{S}, \mathcal{M} \rangle$  and a UCQ  $q_{\mathcal{S}}$  over  $\mathcal{S}$  as input: if  $q_{\mathcal{O}} = \mathfrak{M}$ -MinComplete $(J, q_{\mathcal{S}})$  is such that  $q_r \sqsubseteq_{\mathcal{S}} q_{\mathcal{S}}$ ; then **return**  $q_{\mathcal{O}}$ ; otherwise, **report** "no perfect J-abstraction of  $q_{\mathcal{S}}$  is in  $\mathfrak{M}$ ".

**Example 10.** In Example 3,  $\mathfrak{M}$ -Perfect $(J, q_S)$  returns  $q_O$ , which is the perfect J-abstraction of  $q_S$ .

We conclude this section by establishing termination and correctness of the  $\mathfrak{M}$ -Perfect algorithm.

**Theorem 3.**  $\mathfrak{M}$ -Perfect $(J, q_S)$  terminates and returns the unique (up to J-equivalence) perfect J-abstraction of  $q_S$  if and only if such an abstraction can be expressed in  $\mathfrak{M}$ .

#### 7 Conclusion

We presented a thorough study of monotone abstractions of UCQs in OBDM systems. We proved that best approximations of such abstractions always exist and introduced a query language,  $Datalog_{\mathbf{K}}^{\vee}$ , that captures them. Directions for future work are many. In the context of monotone abstractions, we would like to investigate the case of more expressive ontology languages, e.g., DL-Lite<sub>R</sub>, as well as more expressive source query languages, e.g., unions of conjunctive queries with inequalities and disjunctive Datalog. Finally, the problem of checking whether given best approximations are expressible in simpler and more user friendly languages remains open.

## Acknowledgements

This work has been partially supported by the ANR AI Chair INTENDED (ANR-19-CHIA-0014), by MIUR under the PRIN 2017 project "HOPE" (prot. 2017MMJJRE), and by the EU under the H2020-EU.2.1.1 project TAILOR, grant id. 952215.

#### References

Baader, F.; Calvanese, D.; McGuinness, D.; Nardi, D.; and Patel-Schneider, P. F., eds. 2003. *The Description Logic Handbook: Theory, Implementation and Applications*. Cambridge University Press.

Bienvenu, M. 2016. Ontology-Mediated Query Answering: Harnessing Knowledge to Get More from Data. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI 2016)*, 4058–4061.

Brickley, D.; and Guha, R. V. 2014. RDF Schema 1.1. W3C Recommendation, World Wide Web Consortium. Available at https://www.w3.org/TR/2014/REC-rdf-schema-20140225/.

Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; and Rosati, R. 2007a. EQL-Lite: Effective First-Order Query Processing in Description Logics. In *Proceedings* of the Twentieth International Joint Conference on Artificial Intelligence (IJCAI 2007), 274–279.

Calvanese, D.; De Giacomo, G.; Lembo, D.; Lenzerini, M.; and Rosati, R. 2007b. Tractable Reasoning and Efficient Query Answering in Description Logics: The *DL-Lite* Family. *Journal of Automated Reasoning*, 39(3): 385–429.

Calvanese, D.; and Rosati, R. 2003. Answering recursive queries under keys and foreign keys is undecidable. In *Proceedings of the Tenth International Workshop on Knowledge Representation meets Databases (KRDB 2003)*, volume 79 of *CEUR Electronic Workshop Proceedings*.

Cima, G. 2017. Preliminary Results on Ontology-based Open Data Publishing. In *Proceedings of the Thirtieth International Workshop on Description Logics (DL 2017)*, volume 1879 of *CEUR Electronic Workshop Proceedings*.

Cima, G. 2022. Abstraction in Ontology-based Data Management, volume 348 of Frontiers in Artificial Intelligence and Applications. IOS Press.

Cima, G.; Console, M.; Lenzerini, M.; and Poggi, A. 2021. Abstraction in Data Integration. In *Proceedings of the Thirty-Sixth Annual ACM/IEEE Symposium on Logic in Computer Science (LICS 2021)*, 1–11.

Cima, G.; Croce, F.; and Lenzerini, M. 2021. Query Definability and Its Approximations in Ontology-based Data Management. In *Proceedings of the Thirtieth International Conference on Information and Knowledge Management* (CIKM 2021), 271–280.

Cima, G.; Lenzerini, M.; and Poggi, A. 2019. Semantic Characterization of Data Services through Ontologies. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI 2019)*, 1647–1653. Cima, G.; Lenzerini, M.; and Poggi, A. 2020a. Answering Conjunctive Queries with Inequalities in *DL-Lite*<sub> $\mathcal{R}$ </sub>. In *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI 2020)*, 2782–2789.

Cima, G.; Lenzerini, M.; and Poggi, A. 2020b. Non-Monotonic Ontology-based Abstractions of Data Services. In *Proceedings of the Seventeenth International Conference on Principles of Knowledge Representation and Reasoning* (*KR* 2020), 243–252.

Cuenca Grau, B. 2004. A possible simplification of the semantic web architecture. In *Proceedings of the Thirteenth International World Wide Web Conference (WWW 2004)*, 704–713.

De Giacomo, G.; Lembo, D.; Lenzerini, M.; Poggi, A.; and Rosati, R. 2018. Using Ontologies for Semantic Data Integration. In *A Comprehensive Guide Through the Italian Database Research Over the Last 25 Years*, 187–202.

Duschka, O. M.; and Genesereth, M. R. 1998. Query Planning with Disjunctive Sources. In *Proceedings of the AAAI-98 Workshop on AI and Information Integration*.

Fagin, R.; Kolaitis, P. G.; Miller, R. J.; and Popa, L. 2005. Data Exchange: Semantics and Query Answering. *Theoretical Computer Science*, 336(1): 89–124.

Lenzerini, M. 2018. Managing Data through the Lens of an Ontology. *AI Magazine*, 39(2): 65–74.

Levesque, H. J.; and Lakemeyer, G. 2001. *The Logic of Knowledge Bases*. The MIT Press.

Levy, A. Y.; and Rousset, M.-C. 1998. Combining Horn Rules and Description Logics in CARIN. *Artificial Intelligence*, 104(1–2): 165–209.

Lutz, C.; Marti, J.; and Sabellek, L. 2018. Query Expressibility and Verification in Ontology-based Data Access. In *Proceedings of the Sixteenth International Conference on the Principles of Knowledge Representation and Reasoning* (*KR* 2018), 389–398.

Ortiz, M. 2018. Improving Data Management using Domain Knowledge. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJ-CAI 2018)*, 5709–5713.

Ortiz, M. 2019. Ontology-Mediated Queries from Examples: a Glimpse at the DL-Lite Case. In *Proceedings* of the Fifth Global Conference on Artificial Intelligence (GCAI 2019), volume 65 of EPiC Series in Computing, 1–14.

Poggi, A.; Lembo, D.; Calvanese, D.; De Giacomo, G.; Lenzerini, M.; and Rosati, R. 2008. Linking Data to Ontologies. *Journal on Data Semantics*, X: 133–173.

Rosati, R. 2007. The Limits of Querying Ontologies. In *Proceedings of the Eleventh International Conference on Database Theory (ICDT 2007)*, volume 4353 of *Lecture Notes in Computer Science*, 164–178. Springer.

Xiao, G.; Calvanese, D.; Kontchakov, R.; Lembo, D.; Poggi, A.; Rosati, R.; and Zakharyaschev, M. 2018. Ontology-Based Data Access: A Survey. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI 2018)*, 5511–5519.