

Ontology-Mediated Queries for Probabilistic Databases

Stefan Borgwardt, İsmail İlkan Ceylan
 Faculty of Computer Science
 Technische Universität Dresden, Germany
firstname.lastname@tu-dresden.de

Thomas Lukasiewicz
 Department of Computer Science
 University of Oxford, UK
thomas.lukasiewicz@cs.ox.ac.uk

Abstract

Probabilistic databases (PDBs) are usually incomplete, e.g., containing only the facts that have been extracted from the Web with high confidence. However, missing facts are often treated as being false, which leads to unintuitive results when querying PDBs. Recently, open-world probabilistic databases (OpenPDBs) were proposed to address this issue by allowing probabilities of unknown facts to take any value from a fixed probability interval. In this paper, we extend OpenPDBs by Datalog[±] ontologies, under which both upper and lower probabilities of queries become even more informative, enabling us to distinguish queries that were indistinguishable before. We show that the dichotomy between P and PP in (Open)PDBs can be lifted to the case of first-order rewritable positive programs (without negative constraints); and that the problem can become NP^{PP}-complete, once negative constraints are allowed. We also propose an approximating semantics that circumvents the increase in complexity caused by negative constraints.

1 Introduction

The effort for building *large-scale knowledge bases* from data in an automated manner has resulted in a number of systems including NELL (Mitchell et al. 2015), Yago (Hofmann et al. 2013), DeepDive (Shin et al. 2015), Microsoft’s Probase (Wu et al. 2012), and Google’s Knowledge Vault (Dong et al. 2014). They combine methods from information extraction, natural language processing, relational learning, and databases to process large volumes of uncertain data. The state of the art to store and process such data is founded on *probabilistic databases* (PDBs) (Imielinski and Lipski 1984; Fuhr and Rölleke 1997; Suciu et al. 2011).

Each of the above systems encodes only a portion of the real world, and this description is necessarily *incomplete*. Thus, a meaningful querying semantics must provide a way to deal with missing information. Recently, an effort in this direction was made by introducing *open-world probabilistic databases* (OpenPDBs) (Ceylan, Darwiche, and Van den Broeck 2016), which generalize PDBs to be able to deal with incompleteness. More precisely, in OpenPDBs the probabilities of facts that are not in the database, called *open tuples*, are relaxed to a default probability interval, which is very

different from the *closed-world assumption* of PDBs, which requires the probabilities of such facts to be zero. In the resulting framework of OpenPDBs, query probabilities are given in terms of *upper* and *lower* probability values, which is more in line with an incomplete view of the world.

While forming a natural and flexible basis for querying incomplete data sources, OpenPDBs are limited in the following sense: All open tuples can take on probability values from a single *fixed* interval $[0, \lambda]$, which results in the *same* upper and lower probabilities for many queries. Consider, for instance the PDB containing the probabilistic tuples $\langle \text{Author}(a) : 0.8 \rangle$, $\langle \text{Pub}(a, b) : 0.6 \rangle$, $\langle \text{Pub}(c, d) : 0.9 \rangle$, $\langle \text{Novel}(d) : 1 \rangle$. In OpenPDBs, $\text{Author}(c)$ and $\text{Author}(d)$ evaluate to the *same lower and upper probabilities* (0 and λ , respectively), since both tuples are open. Intuition, however, tells us that c is more likely to be an author, as we already know (with high confidence) that c has published a novel. On the other hand, $\text{Author}(d)$ is unlikely to hold, since we know (almost surely) that d is a novel. Essentially, we lack the common-sense knowledge that

- (i) anyone who has published a novel is an author, and
- (ii) authors and novels are disjoint entities,

which helps us to distinguish such queries. Observe that (i) is a positive axiom and would lead to higher probabilities, whereas (ii) is a negative (constraining) axiom and would entail lower probabilities for some queries.

This problem has been widely studied in the context of classical databases under the name of *ontology-based data access* (OBDA) (Poggi et al. 2008), a popular paradigm that encodes the domain knowledge through an ontology, thus being able to deduce facts not explicitly specified in the database. Following this, we encode the domain knowledge using a Datalog[±] ontology (Cali, Gottlob, and Lukasiewicz 2012), which helps to break down the symmetries between open tuples, letting us distinguish more queries by comparing their upper and lower probability values.

We study the semantic and computational properties of OpenPDBs under Datalog[±] programs. The main distinction between a PDB and an OpenPDB is that the latter represents a set of probability distributions instead of a single one, and introduces the difficulty of choosing the distribution that will maximize (or minimize) the probability of a query. It is known that the data complexity of probabilistic UCQ evalu-

ation in OpenPDBs exhibits the same dichotomy between P and PP as in PDBs for unions of conjunctive queries (Dalvi and Suciu 2012; Ceylan, Darwiche, and Van den Broeck 2016). We lift this dichotomy to first-order rewritable (positive) Datalog[±] programs using standard techniques. We then show that, once negative constraints are allowed, reasoning can become NP^{PP}-hard. This result demonstrates the difference between OpenPDBs and PDBs, as in the latter reasoning with ontologies remains in PP.

We also propose an approximating semantics that circumvents the increase in complexity caused by negative constraints, and lift the dichotomy to general first-order rewritable programs under this semantics. We conclude with complexity results beyond the data complexity for ontology-mediated query evaluation relative to (tuple-independent) PDBs and OpenPDBs. All proofs can be found in the extended version of this paper (see <https://lat.inf.tu-dresden.de/research/papers.html>).

2 Background and Motivation

We briefly recall the basics of tuple-independent PDBs and their open-world variant OpenPDBs. We then highlight the advantages of accessing probabilistic data through a logical theory and provide an overview of Datalog[±] programs.

We consider a relational vocabulary γ consisting of *finite* sets \mathbf{R} of *predicates*, \mathbf{C} of *constants*, and \mathbf{V} of *variables*. A γ -*term* is a constant or a variable. A γ -*atom* is of the form $P(s_1, \dots, s_n)$, where P is an n -ary predicate, and s_1, \dots, s_n are γ -terms. A γ -*tuple* is a γ -atom without variables.

Queries and Databases. A *conjunctive query* (CQ) over γ is an existentially quantified formula $\exists x \phi$, where ϕ is a conjunction of γ -atoms, written as a comma-separated list. A *union of conjunctive queries* (UCQ) is a disjunction of CQs. A query is *Boolean* if it has no free variables. A database \mathcal{D} over γ is a finite set of γ -tuples. The central problem studied for databases is *query evaluation*: Finding all *answers* to a query Q over a database \mathcal{D} , which are assignments of the free variables in Q to constants such that the resulting first-order formula is satisfied in \mathcal{D} in the usual sense, i.e., there is a homomorphism from the atoms in Q to the tuples in \mathcal{D} . In the following, we consider only Boolean queries Q , and focus on the associated decision problem, i.e., deciding whether Q is satisfied in \mathcal{D} , denoted as usual by $\mathcal{D} \models Q$.

Example 1. Consider the database $\mathcal{D}_{ex} := \{\text{Author}(a), \text{Pub}(a, b), \text{Pub}(c, d), \text{Novel}(d)\}$ and the Boolean query $Q_1 := \exists x_1, x_2 \text{ Author}(x_1), \text{Pub}(x_1, x_2)$.¹ Then, $\mathcal{D}_{ex} \models Q_1$, since $\{\text{Author}(a), \text{Pub}(a, b)\} \models Q_1$.

Probabilistic Databases. The most elementary probabilistic database model is based on the tuple-independence assumption. We adopt this model and refer to (Suciu et al. 2011) for details on this model and alternatives. A probabilistic database induces a set of classical databases (called *worlds*), each of which is associated with a probability value.

Formally, a *probabilistic database* (PDB) \mathcal{P} over γ is a finite set of (*probabilistic*) *tuples* of the form $\langle t : p \rangle$, where t

is a γ -tuple and $p \in [0, 1]$, and, whenever $\langle t : p \rangle, \langle t : q \rangle \in \mathcal{P}$, then $p = q$. A PDB \mathcal{P} assigns, to every γ -tuple t , the probability p , if $\langle t : p \rangle \in \mathcal{P}$, and the probability 0, otherwise.

Under the *tuple-independence* assumption, any such probability assignment \mathcal{P} induces the following *unique joint probability distribution* over classical databases \mathcal{D} :

$$P(\mathcal{D}) := \prod_{t \in \mathcal{D}} P(t) \prod_{t \notin \mathcal{D}} (1 - P(t)).$$

Accordingly, query evaluation is enriched to also consider the probabilistic information. More formally, the *probability of a Boolean query* Q w.r.t. \mathcal{P} is $P(Q) := \sum_{\mathcal{D} \models Q} P(\mathcal{D})$. Here, we do not need to consider worlds with probability 0; e.g., if $P(t) = 0$, then the worlds containing t do not affect $P(Q)$.

Example 2. Consider the PDB \mathcal{P}_{ex} from the introduction and Q_1 from Example 1. The probability of Q_1 on \mathcal{P}_{ex} is obtained by summing the probabilities of the worlds that satisfy Q_1 , i.e., all worlds containing the first two tuples, resulting in the probability 0.48. In contrast, the natural query

$$Q_2 := \exists x_1, x_2 \text{ Author}(x_1), \text{Pub}(x_1, x_2), \text{Novel}(x_2)$$

evaluates to 0 on \mathcal{P}_{ex} , since all worlds that satisfy this query have probability 0.

Open-World Probabilistic Databases. An *open-world probabilistic database* (OpenPDB) over γ is a pair $\mathcal{G} = (\mathcal{P}, \lambda)$, where $\lambda \in [0, 1]$ and \mathcal{P} is a PDB. A λ -*completion* of \mathcal{G} is a PDB that is obtained by introducing, for each γ -tuple t that does not occur in \mathcal{P} (called an *open tuple*), a probabilistic tuple $\langle t : p \rangle$ with $p \in [0, \lambda]$. For a fixed value $\alpha \in [0, \lambda]$, we define a special λ -completion, denoted \mathcal{P}_α , in which the probabilities of all open tuples are equal to α . Note that \mathcal{P}_0 is equivalent to \mathcal{P} .

Example 3. Consider the OpenPDB $\mathcal{G}_{ex} := (\mathcal{P}_{ex}, 0.5)$. The set $\mathcal{P}_{ex} \cup \{\langle \text{Novel}(b) : 0.2 \rangle\}$ is a λ -completion of \mathcal{G}_{ex} (tuples with probability 0 are omitted).

An OpenPDB $\mathcal{G} = (\mathcal{P}, \lambda)$ defines the set $K_{\mathcal{G}}$ of all probability distributions P induced by the λ -completions of \mathcal{G} . $K_{\mathcal{G}}$ constitutes a so-called *credal set*, which means that it is closed, convex, and has a finite number of extremal points (Cozman 2000). The range of probabilities of a query under such a set can be expressed as a probability interval. Formally, the *probability interval* of a Boolean query Q w.r.t. \mathcal{G} is $K_{\mathcal{G}}(Q) := [\underline{P}_{\mathcal{G}}(Q), \overline{P}_{\mathcal{G}}(Q)]$, where

$$\underline{P}_{\mathcal{G}}(Q) := \min_{P \in K_{\mathcal{G}}} P(Q) \quad \text{and} \quad \overline{P}_{\mathcal{G}}(Q) := \max_{P \in K_{\mathcal{G}}} P(Q).$$

Example 4. Consider again the OpenPDB \mathcal{G}_{ex} . While the lower probability $\underline{P}_{\mathcal{G}}(Q_2)$ remains 0, the upper probability evaluates to $\overline{P}_{\mathcal{G}}(Q_2) > 0$ due to the λ -completion

$\mathcal{P}_{0.5} = \mathcal{P}_{ex} \cup \{\langle \text{Author}(b) : 0.5 \rangle, \langle \text{Author}(c) : 0.5 \rangle, \dots\}$, which contains all open tuples with probability $\lambda = 0.5$.

This example shows that OpenPDBs improve our view of the domain compared to PDBs. However, we have already illustrated in the introduction that OpenPDBs can further benefit from an axiomatic encoding of the domain knowledge, since many queries involving open tuples will yield

¹For ease of presentation, we assume that γ consists of the symbols appearing in the database and query (and later in the program).

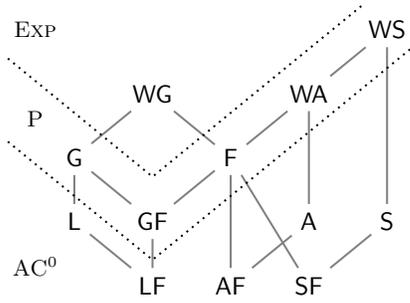


Figure 1: Inclusion relations and data complexity of UCQ entailment for Datalog[±] languages (Lukasiewicz et al. 2015)

the same lower and upper probabilities, although according to common-sense knowledge, they should differ. This motivates our introduction of a logical theory, in the form of Datalog[±] rules, to formalize such knowledge.

Datalog[±] Programs. We now extend the vocabulary γ by a (potentially infinite) set \mathbf{N} of *nulls*. An *instance* I over γ is a (possibly infinite) set of γ -tuples that may additionally contain nulls.

A *tuple-generating dependency* (TGD) σ is a first-order formula $\forall \mathbf{x} \varphi(\mathbf{x}) \rightarrow \exists \mathbf{y} \mathsf{P}(\mathbf{x}, \mathbf{y})$, where $\varphi(\mathbf{x})$ is a conjunction of γ -atoms, called the *body* of σ , and $\mathsf{P}(\mathbf{x}, \mathbf{y})$ is a γ -atom, called the *head* of σ . A *negative constraint* (NC) ν is a first-order formula $\forall \mathbf{x} \varphi(\mathbf{x}) \rightarrow \perp$, where $\varphi(\mathbf{x})$ is a conjunction of γ -atoms, called the *body* of ν , and \perp is the truth constant *false*. A (Datalog[±]) *program* Σ is a finite set of TGDs and NCs.² An *ontology-mediated query* (OMQ) is a pair (Q, Σ) , where Σ is a program, and Q is a Boolean query.

An instance I satisfies a TGD or NC σ , if $I \models \sigma$, where \models denotes the standard first-order entailment relation. I satisfies a program Σ , written $I \models \Sigma$, if I satisfies each formula in Σ . The set of *models* of a program Σ relative to a database \mathcal{D} , denoted $\text{mods}(\mathcal{D}, \Sigma)$, is $\{I \mid I \supseteq \mathcal{D} \text{ and } I \models \Sigma\}$. \mathcal{D} is *consistent* w.r.t. Σ , if $\text{mods}(\mathcal{D}, \Sigma)$ is non-empty. The OMQ (Q, Σ) is *entailed* by \mathcal{D} , denoted $\mathcal{D} \models (Q, \Sigma)$, if $I \models Q$ holds for all $I \in \text{mods}(\mathcal{D}, \Sigma)$.

In general, the entailment problem is undecidable (Beeri and Vardi 1981). For this reason, many different restrictions on the TGDs have been proposed. We consider here *guarded* (G), *linear* (L), *sticky* (S), *acyclic* (A), *weakly guarded* (WG), *weakly sticky* (WS), and *weakly acyclic* (WA) sets of TGDs (Calì, Gottlob, and Kifer 2013; Calì, Gottlob, and Pieris 2012). Other important classes are given by *full* TGDs (F), *full and guarded* TGDs (GF), and similarly for LF, SF, and AF. Figure 1 illustrates the inclusion relations between these classes; for a more detailed description, see the extended version of this paper. We extend all these notions to programs Σ in the obvious way; for instance, Σ is guarded if all the TGDs in Σ are guarded. In the following, we use \mathcal{L} to denote the set of Datalog[±] languages introduced above.

²For brevity, we omit the universal quantifiers in TGDs and NCs, and use commas (instead of \wedge) for conjoining atoms. For clarity, we consider single-atom-head TGDs; however, our results can be easily extended to TGDs with conjunctions of atoms in the head.

A key paradigm in OBDA is the *FO-rewritability* of queries; an OMQ (Q, Σ) is *FO-rewritable*, if there exists a Boolean UCQ Q_Σ such that, for all databases \mathcal{D} that are consistent w.r.t. Σ , we have $\mathcal{D} \models (Q, \Sigma)$ iff $\mathcal{D} \models Q_\Sigma$. In this case, Q_Σ is called a *FO-rewriting* of (Q, Σ) . A class of programs X is *FO-rewritable*, if it admits an FO-rewriting for any UCQ and program in X ; these classes are characterized by a data complexity of AC⁰ (see Figure 1).

3 Ontology-Mediated Queries for OpenPDBs

We now introduce the basics of OMQ evaluation relative to OpenPDBs. In the following, we assume that the input PDB \mathcal{P} induces a consistent distribution w.r.t. the program. Formally, a probability distribution P is *consistent* w.r.t. Σ , if the database $\{t \mid P(t) > 0\}$ is consistent w.r.t. Σ . Note that this assumption does not change the nature of the problem. The semantics of OMQs is again based on λ -completions. The difference appears in the deductive power provided by the Datalog[±] program, which is taken into consideration in the query semantics.

Definition 5 (Semantics). The *probability of an OMQ* (Q, Σ) *relative to a probability distribution* P is

$$P(Q, \Sigma) = \sum_{\mathcal{D} \models (Q, \Sigma)} P(\mathcal{D}),$$

where \mathcal{D} ranges over all databases over γ . The *probability interval of* (Q, Σ) *relative to an OpenPDB* \mathcal{G} is then given by $K_{\mathcal{G}}(Q, \Sigma) := [\underline{P}_{\mathcal{G}}(Q, \Sigma), \overline{P}_{\mathcal{G}}(Q, \Sigma)]$, where

$$\underline{P}_{\mathcal{G}}(Q, \Sigma) := \min_{P \in K_{\mathcal{G}}} \{P(Q, \Sigma) \mid P \text{ is consistent w.r.t. } \Sigma\},$$

$$\overline{P}_{\mathcal{G}}(Q, \Sigma) := \max_{P \in K_{\mathcal{G}}} \{P(Q, \Sigma) \mid P \text{ is consistent w.r.t. } \Sigma\}.$$

The special case of $\lambda = 0$ corresponds to having a single (closed-world) PDB \mathcal{P} . In this case, we simply speak of the *probability of* (Q, Σ) *relative to a PDB* \mathcal{P} .

This semantics defers the decision of whether a world satisfies a query to an entailment test. However, we maximize only over consistent λ -completions, i.e., the ones that induce consistent distributions, which is the most important aspect of this semantics.

3.1 Semantic Considerations

In the following, we evaluate our semantics w.r.t. the goals identified in the motivation of this paper, and discuss our choice of restricting to the consistent λ -completions.

Distinguishing Queries. We argued that OpenPDBs can benefit from an axiomatic encoding of the knowledge of the domain. Consider again our running example, which is now enriched with a program.

Example 6. Consider the OpenPDB \mathcal{G}_{ex} given before and the program $\Sigma_{ex} := \{\text{Author}(x), \text{Novel}(x) \rightarrow \perp, \text{Pub}(x, y), \text{Novel}(y) \rightarrow \text{Author}(x)\}$ which states that authors and novels are disjoint entities, and that anyone who has published a novel is an author. The lower probability of $\text{Author}(d)$ remains 0, while the upper probability is now reduced to 0 with the help of the program Σ_{ex} . In contrast,

the lower probability of $\text{Author}(c)$ increases to 0.9, while the upper probability increases to 0.95. These intervals are much more informative than the default interval $[0, 0.5]$.

Restricting to Consistent Distributions. The most subtle aspect of choosing the *best* distribution is the question of how to deal with inconsistent worlds. Ignoring inconsistencies (and optimizing over *all* completions) leads to a drowning effect: since inconsistent worlds entail everything, this semantics would be biased towards choosing inconsistent λ -completions. This does not satisfy our goals, as even an unsatisfiable query could evaluate to a positive probability.

An alternative approach, which is standard for (closed-world) PDBs, and is quite intuitive at first glance, would be to choose the distribution which maximizes the conditional probability $P((Q, \Sigma) \mid (\mathcal{D}, \Sigma) \neq \perp)$, i.e., the probability of the query on the set of all consistent worlds. A careful inspection, however, shows that this semantics also favors inconsistent distributions over consistent ones. To illustrate this, consider our running example, and suppose that we want to compute the upper probability of Q_2 (mediated by Σ_{ex}). The semantics based on the conditional probability would favor the λ -completion $\mathcal{P}_{0.5}$, even though this PDB is highly inconsistent. This is mainly due to the normalization process internal to the computation. As part of this normalization, the probability mass of inconsistent worlds is distributed to consistent worlds. As a consequence, it is often possible to increase the query probability by simply increasing the probability of inconsistent worlds. This is not a desired effect, since we are interested in finding the most suitable λ -completion from the open world, and not the one that increases the query probability by increasing the probability mass of inconsistent worlds.

To avoid such drowning effects, our proposal considers only consistent distributions. That is, we do not want to introduce inconsistencies when completing our knowledge over the domain by choosing a λ -completion. One drawback of our approach is the fact that inconsistencies are not tolerated even if the inconsistency degree is very small. However, it would be easy to introduce a threshold value, say 0.1, to tolerate the inconsistent completions where the probability of the inconsistent worlds does not exceed this threshold.

4 Data Complexity Results

We now formulate the task of probabilistic query evaluation as a decision problem.

Definition 7 (Decision Problems). Let (Q, Σ) be an OMQ, \mathcal{G} an OpenPDB and $p \in [0, 1]$. The problem of *upper* (resp., *lower*) *probabilistic query entailment* is to decide whether $\overline{P}_{\mathcal{G}}(Q, \Sigma) > p$ (resp., $\underline{P}_{\mathcal{G}}(Q, \Sigma) < p$) holds. *Probabilistic query entailment relative to PDBs* is a special case, where $\lambda = 0$.

Note that this definition is rather general, but in the scope of this paper, we are concerned with UCQs, and thus we use the term *probabilistic UCQ entailment* instead. Moreover, we are mainly concerned with the *data complexity*, which is calculated based on the size of the OpenPDB; i.e., the schema \mathbf{R} , the query Q , and the program Σ are assumed to

be fixed (Vardi 1982). The relevant data complexity results for UCQ entailment in Datalog $^{\pm}$ are summarized in Figure 1.

Most of our complexity results are related to the complexity class PP (Gill 1977), which comprises the languages recognized by a polynomial-time non-deterministic Turing machine that accepts an input if and only if *more than half* of the computation paths are accepting (Torán 1991). Intuitively, PP is the decision counterpart of #P (Valiant 1979). For details on the complexity classes used in our results, and the types of reductions, we refer to the extended version of this paper. It has been shown in (Dalvi and Suciu 2012) that probabilistic UCQ entailment for PDBs exhibits a dichotomy between P and PP. Queries that admit a P algorithm are called *safe* and the remaining ones *unsafe*. This result has been lifted to OpenPDBs in (Ceylan, Darwiche, and Van den Broeck 2016). For detailed insights on the class of safe queries, we refer to the original papers. The CQ $\exists x, y C(x) \wedge L(x, y) \wedge S(y)$ is the prototypical example of an unsafe query; it is connected and can not be decomposed into independent queries in an efficient manner (applying certain rules from (Dalvi and Suciu 2012)). However, removing any of the atoms from this query makes it safe.

We borrow this notion, and say that an OMQ (Q, Σ) is *safe*, if there exist polynomial-time algorithms for lower and upper probabilistic entailment of (Q, Σ) relative to any OpenPDB (resp., PDB).

4.1 Positive Programs

We first consider *positive* Datalog $^{\pm}$ programs, which do not contain NCs. Under this restriction, there are no inconsistent distributions, and Definition 5 simplifies. We later show that this distinction is important, since the complexity increases in the presence of NCs. This is surprising, as in the classical case NCs are usually not problematic.

Recall that OpenPDBs induce an infinite set of probability distributions that form a credal set, which has the following useful property (Cozman 2000): To determine the upper or lower probability of an event, it suffices to consider the *extremal* probability distributions, which are obtained by setting the probability values of all elementary events to one of the extreme points. In the context of OpenPDBs, this means that each of the open tuples may have probability λ or 0, but no intermediate choices need to be examined. For UCQs, this implies an even stronger result.

Lemma 8. *Let (Q, Σ) be an OMQ, where Q is a UCQ and Σ is a positive Datalog $^{\pm}$ program. Then, it holds that $K_{\mathcal{G}}(Q, \Sigma) = [P_{\mathcal{P}_0}(Q, \Sigma), P_{\mathcal{P}_\lambda}(Q, \Sigma)]$.*

Thus, it suffices to consider a single λ -completion (either \mathcal{P}_0 or \mathcal{P}_λ) and the particular distribution it induces. As a result, probabilistic UCQ entailment can be solved by standard methods; i.e., summing up the probabilities of all worlds that pass the entailment test. This naïve approach yields tight complexity bounds for the considered problems.

Theorem 9. *Probabilistic UCQ entailment is PP-complete for the languages in $\mathcal{L} \setminus \{\text{WG}\}$; it is EXP-complete in WG.*

This result is of no surprise given the PP-hardness of inference in OpenPDBs. However, all our PP-hardness results

are based on the result of (Dalvi and Suciu 2012), and hence are valid only with respect to Turing reductions. All other complexity results in this paper also hold under standard many-one reductions. It is an open problem to find a UCQ for which probabilistic entailment is PP-hard w.r.t. many-one reductions. The striving question is now whether it is possible to lift the dichotomy result from OpenPDBs. For this purpose, we elaborate on query rewritability.

Lemma 10. *Let (Q, Σ) be an OMQ, P be a tuple-independent probability distribution over worlds such that $P(\mathcal{D}) = 0$ whenever \mathcal{D} is inconsistent w.r.t. Σ , and Q_Σ be an FO-rewriting of (Q, Σ) . Then, we have $P(Q, \Sigma) = P(Q_\Sigma)$.*

Since all worlds are consistent under positive programs, Lemmas 8 and 10 imply that we can reduce probabilistic UCQ entailment under positive programs to the case of OpenPDBs via query rewriting.

Corollary 11. *Let (Q, Σ) be an OMQ, where Q is a UCQ, and Σ is a positive program, and Q_Σ be an FO-rewriting of (Q, Σ) . Then, for any OpenPDB \mathcal{G} , it holds that $\overline{P}_{\mathcal{G}}(Q, \Sigma) = \overline{P}_{\mathcal{G}}(Q_\Sigma)$ and $\underline{P}_{\mathcal{G}}(Q, \Sigma) = \underline{P}_{\mathcal{G}}(Q_\Sigma)$.*

We now obtain a dichotomy from the results in (Dalvi and Suciu 2012; Ceylan, Darwiche, and Van den Broeck 2016).

Theorem 12. *Let (Q, Σ) be an OMQ, where Q is a UCQ, and Σ is a positive program, and Q_Σ be a rewriting of (Q, Σ) . Then, (Q, Σ) is safe iff Q_Σ is safe (over OpenPDBs). If (Q, Σ) is not safe, then it is PP-hard.*

In particular, either all rewritings of (Q, Σ) are safe, or none of them are. Hence, in FO-rewritable languages, we can take an *arbitrary* rewriting and check safety using the characterization of (Dalvi and Suciu 2012). Such a rewriting can be obtained by well-known algorithms, e.g., using backward chaining of TGDs (Gottlob, Orsi, and Pieris 2011).

To conclude this section, we illustrate some effects that simple positive programs can have on the complexity of probabilistic query entailment.

Example 13. The query $\exists x, y C(x) \wedge M(x, y)$ is safe for OpenPDBs. It becomes unsafe under the TGD $R(x, y), T(y) \rightarrow M(x, y)$, since then it rewrites to the query $(\exists x, y C(x), M(x, y)) \vee (\exists x, y C(x), R(x, y), T(y))$. Conversely, the CQ $\exists x, y C(x) \wedge L(x, y) \wedge S(y)$ is not safe for OpenPDBs, but becomes safe under $L(x, y) \rightarrow S(y)$, as it rewrites to $\exists x, y C(x) \wedge L(x, y)$. Note that these are very simple TGDs, which are full, acyclic, guarded, and sticky.

4.2 Programs with Negative Constraints

In the presence of NCs, it still suffices to consider the extremal λ -completions. In fact, once the correct completion is known, the probabilistic UCQ entailment problem can still be reduced to probabilistic inference (in FO-rewritable languages). The key difference in the presence of NCs is that we have to make sure that this completion is consistent. That is, choosing the completion \mathcal{P}_λ that sets all open tuples to λ (as in Lemma 8) is not feasible, as this will very likely lead to inconsistencies. However, observe that the *lower* probability can still be obtained from the completion \mathcal{P}_0 (which we assumed to be consistent), and hence the previous results still hold for lower probabilistic UCQ entailment with NCs.

A naïve way of solving the upper probabilistic UCQ entailment problem is to *guess* a λ -completion and then check whether it is consistent and compare the resulting probability to the threshold. This yields an NP^{PP} upper bound for our decision problem. Our next result shows a matching lower bound for the class GF, and so for all considered Datalog $^\pm$ languages with data complexity above AC^0 (see Figure 1).

Theorem 14. *Upper probabilistic UCQ entailment is NP^{PP} -complete in full, guarded programs. It is PP-complete for all languages with polynomial data complexity once restricted to PDBs.*

This result is by reduction from the NP^{PP} -complete problem of finding a partial assignment for designated variables of a propositional formula in CNF, for which the number of satisfying assignments extending this partial assignment is above some threshold (Wagner 1986). On the one hand, this result is surprising, as NCs are not problematic for PDBs, even with normalization semantics; on the other hand, this is not so surprising, as non-monotonicity is also a source of additional hardness in OpenPDBs: query evaluation becomes NP^{PP} -complete in OpenPDBs if negated atoms are allowed in UCQs (Ceylan, Darwiche, and Van den Broeck 2016). In contrast, our result applies to UCQs without negated atoms, and thus it is much more involved. The proof encodes the non-determinism into the NCs, which are not as powerful as non-monotone queries, and uses TGDs to check the satisfaction condition of the clauses in the CNF.

Before concluding this section, we illustrate the effects of NCs on some examples, which also show the difficulties in lifting the dichotomy of Theorem 12 to NCs.

Example 15. Consider the query $(\exists x, y C(x) \wedge S(y)) \vee (\exists x, y C(x) \wedge L(x, y))$, which is not safe for OpenPDBs, but becomes safe relative to the NC $S(y), L(x, y) \rightarrow \perp$. The reason is that the algorithm of (Dalvi and Suciu 2012) that decides safety will produce the unsafe query $\exists x, y C(x) \wedge S(y) \wedge L(x, y)$ through a sequence of reduction rules; however, this query automatically has probability 0 under the given NC, and hence becomes trivially safe.

Approximations for Programs with NCs. Motivated by the high complexity of reasoning in programs with NCs, we propose an alternative semantics, which approximates the semantics of Definition 5. Observe that the upper probability $\overline{P}_{\mathcal{G}}(Q, \Sigma)$ will always be obtained at a λ -completion that adds as many open tuples as possible to the original \mathcal{P} without causing an inconsistency. This is related to the notion of a database *repair*, which is a maximal consistent subset of an inconsistent database (Arenas, Bertossi, and Chomicki 1999). Instead of considering all possible repairs, an easier alternative is to compute the intersection of all repairs and use this for query answering (Lembo et al. 2010). In our setting, however, we are not actually repairing an inconsistent initial database \mathcal{P} , but rather assume that all tuples in \mathcal{P} are correct and consistent, and hence need to take care that no such tuples are removed in this intersection. Formally, given an OMQ (Q, Σ) and an OpenPDB $\mathcal{G} = (\mathcal{P}, \lambda)$, we consider the special λ -completion \mathcal{P}_\cap that is constructed as the intersection of all \subseteq -maximal consistent subsets of \mathcal{P}_λ that con-

Datalog [±] Languages	PDBs		OpenPDBs	
	<i>fs-c.</i>	<i>fp-c.</i>	<i>fs-c.</i>	<i>fp-c.</i>
L, LF, AF	PP ^{NP}	PP ^{NP}	NP ^{PP^{NP}}	NP ^{PP^{NP}}
G	EXP	PP ^{NP}	EXP	NP ^{PP^{NP}}
WG	EXP	EXP	EXP	EXP
S, F, SF, GF	PP ^{NP}	PP ^{NP}	NP ^{PP^{NP}}	NP ^{PP^{NP}}
A	NEXP	PP ^{NP}	in P ^{NE}	NP ^{PP^{NP}}
WS, WA	2EXP	PP ^{NP}	2EXP	NP ^{PP^{NP}}

Table 1: (fs/fp)-combined complexity of probabilistic UCQ entailment relative to OpenPDBs and PDBs.

tain \mathcal{P} (all tuples not in this intersection have probability 0).

Definition 16 (Intersection Semantics). The *probability interval* of (Q, Σ) relative to an OpenPDB $\mathcal{G} = (\mathcal{P}, \lambda)$ under the *intersection semantics* is defined as $K_{\mathcal{P}}^{\cap}(Q, \Sigma) := [P_{\mathcal{P}_0}(Q, \Sigma), P_{\mathcal{P}_n}(Q, \Sigma)]$.

As with positive programs (cf. Lemma 8), probabilistic UCQ entailment under this semantics is PP-complete in all Datalog[±] languages where classical UCQ entailment is in P. More interestingly, we can also show a dichotomy for FO-rewritable queries with the help of Lemma 10.

Theorem 17. *Let (Q, Σ) be an OMQ, where Q is a UCQ, and Σ is a program, and Q_{Σ} be a rewriting of Q relative to Σ . Then, (Q, Σ) is safe under intersection semantics iff Q_{Σ} is safe (over OpenPDBs). If (Q, Σ) is not safe under intersection semantics, then it is PP-hard.*

5 Beyond Data Complexity

For the sake of completeness, we also provide results beyond the data complexity. We consider *fixed-program combined (fp-combined) complexity*, which is calculated in the size of the database and the query, while the program and schema remain fixed. Additionally, we remove the assumption that the program is fixed, and study *fixed-schema combined (fs-combined) complexity*. Our results are summarized in Table 1; all results except one are completeness results. The results are given relative to both PDBs and OpenPDBs to emphasize the computational differences.

Theorem 18. *Let X be a class of programs, and UCQ entailment in X be \mathbf{C} -complete in (fs/fp)-combined complexity. Then, probabilistic UCQ entailment in X is \mathbf{C} -hard and in $\text{PSPACE}^{\mathbf{C}}$ in (fs/fp)-combined complexity. If $\mathbf{C} = \text{NEXP}$, it is in P^{NE} , and NEXP-complete when restricted to PDBs.*

Hence, if $\mathbf{C} = \text{EXP}$ or $\mathbf{C} = 2\text{EXP}$, the complexity is not affected by adding OpenPDBs, since the complexity of UCQ entailment dominates the problem. We now consider the special case of NP-complete classes.

Theorem 19. *Let X be a class of programs. If UCQ entailment in X is NP-complete in (fs/fp)-combined complexity, then probabilistic UCQ entailment in X is complete for $\text{NP}^{\text{PP}^{\text{NP}}}$ in (fs/fp)-combined complexity; it is complete for PP^{NP} when restricted to a PDB.*

The hardness proof uses no TGDs and only one NC. This implies that the additional hardness in probabilistic UCQ entailment relative to OpenPDBs is caused solely by the interaction between NCs and the open-world semantics. This provides more evidence that OpenPDBs with NCs are more powerful than PDBs with NCs.

6 Related Work

Our work builds on the research on probabilistic databases, which has a long tradition (Imieliski and Lipski 1984; Fuhr and Rölleke 1997; Suciu et al. 2011). We focus on tuple-independent probabilistic databases, with an emphasis on the dichotomy result of Dalvi and Suciu (2012). The most closely related work is by Jung and Lutz (2012), where the authors lift the dichotomy result of PDBs to the lightweight description logics \mathcal{EL} and $DL\text{-Lite}$ over PDBs; they even describe the case of an ontology language that is not FO-rewritable and causes all CQs of a certain form to become #P-hard. In contrast, we consider the more expressive languages of the Datalog[±] family and provide results both relative to PDBs and OpenPDBs. We show that the semantic differences between these formalisms lead to different results (even in the data complexity).

Most of the recent work on probabilistic query answering using ontologies is based on lightweight ontology languages. Some (D’Amato, Fanizzi, and Lukasiewicz 2008; Ceylan and Peñaloza 2015; Gottlob et al. 2013) result from a combination of ontologies with probabilistic graphical models such as Bayesian networks (Pearl 1988) or Markov logic networks (Richardson and Domingos 2006). Both the semantics and the assumptions used in these works are very different than ours. More closely related is the work by Ceylan, Peñaloza, and Lukasiewicz (2016), where the computational complexity of query answering in probabilistic Datalog[±] under the possible world semantics is investigated. Note, however, that the authors consider PDBs, and thus a unique probability distribution. Moreover, even for PDBs, the results are not comparable as they allow conditional dependencies and hence the hardness results do not apply to the special case of tuple-independent PDBs.

Possible world semantics is common in probabilistic logic programming and relational probabilistic models (Renkens et al. 2012; Kwiatkowska, Norman, and Parker 2002; Poole 1997). OpenPDBs extend this semantics to a (finite) open universe, and allow imprecise probabilities (Levi 1980) for tuples in this universe. The latter can be seen as analogous to extending Bayesian networks (Pearl 1988) to credal networks (Cozman 2000; De Campos and Cozman 2005). Our framework enriches OpenPDBs further by mediating the query with an ontology, where the query evaluation problem over a database is replaced with a logical entailment problem, allowing us to deduce implicitly encoded facts.

7 Summary and Outlook

We introduced a refinement of the recently proposed OpenPDBs, using Datalog[±] ontologies to express additional background knowledge, and lifted the dichotomy from (Dalvi and Suciu 2012; Ceylan, Darwiche, and Van den Broeck

2016) to all FO-rewritable languages for positive programs. We showed that NCs can increase the worst-case complexity, and proposed an approximating semantics circumventing the increase in the complexity. Additionally, we provided complexity results beyond the data complexity.

In future work, we want to determine whether it is possible to obtain a dichotomy result for programs with NCs for FO-rewritable Datalog[±] languages. Similarly, the question whether the P-complete languages admit a dichotomy when restricting to positive programs is left as future work. Note also that we assume a finite set of constants (as in OpenPDBs), but allow infinitely many unknown individuals (nulls). Dealing with distributions over infinitely many objects as in BLOG (Milch et al. 2005) is an important task, and a crucial part of future work.

Acknowledgments

This work is supported by the German Research Foundation (DFG) within the Collaborative Research Center SFB 912 HAEC and the Graduiertenkolleg RoSI (GRK 1907), and by the UK EPSRC grants EP/J008346/1, EP/L012138/1, EP/M025268/1, and EP/N510129/1.

References

- Arenas, M.; Bertossi, L.; and Chomicki, J. 1999. Consistent query answers in inconsistent databases. In *Proc. of PODS*, 68–79. ACM.
- Beeri, C., and Vardi, M. Y. 1981. The implication problem for data dependencies. In *Proc. of ICALP*, 73–85. Springer.
- Calì, A.; Gottlob, G.; and Kifer, M. 2013. Taming the infinite chase: Query answering under expressive relational constraints. *JAIR* 48:115–174.
- Calì, A.; Gottlob, G.; and Lukasiewicz, T. 2012. A general Datalog-based framework for tractable query answering over ontologies. *J. Web Sem.* 14:57–83.
- Calì, A.; Gottlob, G.; and Pieris, A. 2012. Towards more expressive ontology languages: The query answering problem. *AIJ* 193:87–128.
- Ceylan, İ. İ., and Peñaloza, R. 2015. Probabilistic query answering in the Bayesian description logic BEL. In *Proc. of SUM*, 21–35.
- Ceylan, İ. İ.; Darwiche, A.; and Van den Broeck, G. 2016. Open-world probabilistic databases. In *Proc. of KR*. AAAI Press.
- Ceylan, İ. İ.; Peñaloza, R.; and Lukasiewicz, T. 2016. Complexity results for probabilistic Datalog[±]. In *Proc. of ECAI*. IOS Press.
- Cozman, F. G. 2000. Credal networks. *AIJ* 120(2):199–233.
- Dalvi, N., and Suciu, D. 2012. The dichotomy of probabilistic inference for unions of conjunctive queries. *J. ACM* 59(6):1–87.
- D’Amato, C.; Fanizzi, N.; and Lukasiewicz, T. 2008. Tractable reasoning with Bayesian description logics. In *Proc. of SUM*, 146–159. Springer.
- De Campos, C. P., and Cozman, F. G. 2005. The inferential complexity of Bayesian and credal networks. In *Proc. of IJCAI*, 1313–1318. AAAI Press.
- Dong, X. L.; Gabrielovich, E.; Heitz, G.; Horn, W.; Lao, N.; Murphy, K. P.; Strohmann, T.; Sun, S.; and Zhang, W. 2014. Knowledge Vault: A web-scale approach to probabilistic knowledge fusion. In *Proc. of SIGKDD*, 601–610. ACM.
- Fuhr, N., and Rölleke, T. 1997. A probabilistic relational algebra for the integration of information retrieval and database systems. *ACM Trans. Inf. Systems* 15(1):32–66.
- Gill, J. T. 1977. Computational complexity of probabilistic Turing machines. *SIAM J. on Computing* 6(4):675–695.
- Gottlob, G.; Lukasiewicz, T.; Martínez, M. V.; and Simari, G. I. 2013. Query answering under probabilistic uncertainty in Datalog[±] ontologies. *Ann. Math. Artif. Intell.* 69(1):37–72.
- Gottlob, G.; Orsi, G.; and Pieris, A. 2011. Ontological queries: Rewriting and optimization. In *Proc. of ICDE*, 2–13. IEEE Press.
- Hoffart, J.; Suchanek, F. M.; Berberich, K.; and Weikum, G. 2013. YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. In *Proc. of IJCAI*, 3161–3165.
- Imieliski, T., and Lipski, W. 1984. Incomplete information in relational databases. *J. ACM* 31(4):761–791.
- Jung, J. C., and Lutz, C. 2012. Ontology-based access to probabilistic data with OWL QL. In *Proc. of ISWC*, 182–197. Springer.
- Kwiatkowska, M.; Norman, G.; and Parker, D. 2002. PRISM: Probabilistic symbolic model checker. In *Proc. TOOLS*, 200–204.
- Lembo, D.; Lenzerini, M.; Rosati, R.; Ruzzi, M.; and Savo, D. F. 2010. Inconsistency-tolerant semantics for description logics. In *Proc. of RR*, 103–117. Springer.
- Levi, I. 1980. *The Enterprise of Knowledge*. MIT Press.
- Lukasiewicz, T.; Martínez, M. V.; Pieris, A.; and Simari, G. I. 2015. From classical to consistent query answering under existential rules. In *Proc. of AAI*, 40–45. AAAI Press.
- Milch, B.; Marthi, B.; Russell, S.; Sontag, D.; Ong, D. L.; and Kolobov, A. 2005. BLOG: Probabilistic models with unknown objects. In *Proc. of IJCAI*, 1352–1359. Morgan Kaufmann.
- Mitchell, T.; Cohen, W.; Hruschka, E.; Talukdar, P.; Betteridge, J.; Carlson, A.; Dalvi, B.; and Gardner, M. 2015. Never-ending learning. In *Proc. of AAI*, 2302–2310. AAAI Press.
- Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann.
- Poggi, A.; Lembo, D.; Calvanese, D.; Giacomo, G. D.; Lenzerini, M.; and Rosati, R. 2008. Linking data to ontologies. *J. Data Sem.* 10:133–173.
- Poole, D. 1997. The independent choice logic for modelling multiple agents under uncertainty. *AIJ* 94(1-2):7–56.
- Renkens, J.; Shterionov, D.; Van den Broeck, G.; Vlasselaer, J.; Fierens, D.; Meert, W.; Janssens, G.; and De Raedt, L. 2012. ProbLog2: From probabilistic programming to statistical relational learning. In *Proc. of NIPS*, 1–5.
- Richardson, M., and Domingos, P. 2006. Markov logic networks. *Mach. Learn.* 62(1-2):107–136.
- Shin, J.; Wang, F.; Sa, C. D.; Zhang, C.; and Wu, S. 2015. Incremental knowledge base construction using DeepDive. In *Proc. of VLDB*.
- Suciu, D.; Olteanu, D.; Ré, C.; and Koch, C. 2011. *Probabilistic Databases*. Morgan & Claypool.
- Torán, J. 1991. Complexity classes defined by counting quantifiers. *J. ACM* 38(3):753–774.
- Valiant, L. G. 1979. The complexity of computing the permanent. *TCS* 8(2):189–201.
- Vardi, M. Y. 1982. The complexity of relational query languages. In *Proc. of STOC*, 137–146.
- Wagner, K. W. 1986. The complexity of combinatorial problems with succinct input representation. *Acta Inf.* 23(3):325–356.
- Wu, W.; Li, H.; Wang, H.; and Zhu, K. Q. 2012. Probase: A probabilistic taxonomy for text understanding. In *Proc. of SIGMOD*, 481–492. ACM.