# Tractable Reasoning with Bayesian Description Logics

Claudia d'Amato[1], Nicola Fanizzi[1], and Thomas Lukasiewicz[2,3]

[1] Dipartimento di Informatica, Università degli Studi di Bari
Campus Universitario, Via Orabona 4, 70125 Bari, Italy
{claudia.damato,fanizzi}@di.uniba.it
[2] Computing Laboratory, University of Oxford
Wolfson Building, Parks Road, Oxford OX1 3QD, UK
thomas.lukasiewicz@comlab.ox.ac.uk
[3] Institut für Informationssysteme, Technische Universität Wien,
Favoritenstraße 9-11, 1040 Wien, Austria
lukasiewicz@kr.tuwien.ac.at

**Abstract.** The *DL-Lite* family of tractable description logics lies between the semantic web languages RDFS and OWL Lite. In this paper, we present a probabilistic generalization of the *DL-Lite* description logics, which is based on Bayesian networks. As an important feature, the new probabilistic description logics allow for flexibly combining terminological and assertional pieces of probabilistic knowledge. We show that the new probabilistic description logics are rich enough to properly extend both the *DL-Lite* description logics as well as Bayesian networks. We also show that satisfiability checking and query processing in the new probabilistic description logics is reducible to satisfiability checking and query processing in the *DL-Lite* family. Furthermore, we show that satisfiability checking and answering unions of conjunctive queries in the new logics can be done in LogSpace in the data complexity. For this reason, the new probabilistic description logics are very promising formalisms for data-intensive applications in the Semantic Web involving probabilistic uncertainty.

**Keywords:** Bayesian description logics, tractable reasoning, description logics, ontologies, *DL-Lite*, Bayesian networks, Semantic Web.

## 1 Introduction

The *Semantic Web (SW)* is a web of data to be shared by machines in order to help them to understand the information in the Web and to perform various complex tasks autonomously, such as data integration, discovery, etc. Ontology languages such as OWL have been proposed to express concepts and relations in this context. These are ultimately based on *description logics (DLs)* [1].

Intuitively, description logics model a domain of interest in terms of concepts and roles, which represent classes of individuals resp. binary relations on classes of individuals. A description logic knowledge base (or ontology) encodes in particular (i) subsumption relationships between concepts, (ii) subsumption relationships between roles, (iii) instance relationships between individuals and concepts, and (iv) instance relationships between pairs of individuals and roles.

The heterogeneity of the data sources clearly introduces degrees of uncertainty in the data manipulation process, which causes purely logical methods to fall short. Hence, extended forms of ontology languages have been proposed in order to deal with uncertainty through probabilistic reasoning [20].

There is a plethora of applications with an urgent need for handling uncertain knowledge in ontologies, especially in areas like medicine, biology, defense, and astronomy. Furthermore, there are strong arguments for the critical need of dealing with probabilistic uncertainty in ontologies in the Semantic Web (in order to encode ambiguous information, such as "John is a student with the probability 0.7 and a teacher with the probability 0.3", which is very different from vague/fuzzy information, such as "John is tall with the degree of truth 0.7"):

- Concepts of a probabilistic ontology are probabilistically related. For example, two concepts either may be logically related via a subsumption or disjointness relationship, or they may show a certain degree of overlap. Probabilistic ontologies allow for quantifying these degrees of overlap, for reasoning about them, and for using them in semantic web applications, e.g., information retrieval. The degrees of concept overlap may also be exploited in personalization and recommender systems.
- The Semantic Web will consist of a huge collection of different ontologies. So, in semantic web applications of reasoning and retrieval, one may have to align the concepts of different ontologies, which is called ontology matching/mapping. In general, the concepts of different ontologies do not match exactly, and we have to deal with degrees of concept overlap as above, which are determined by automatic or semi-automatic tools or experts. These degrees of concept overlap are then represented in probabilistic ontologies, which thus allows for inference about the degrees of overlap between other concepts and about uncertain instance relationships.
- The Semantic Web will likely contain controversial information in different web sources. This can be handled via probabilistic data integration by associating with every web source a probability describing its degree of reliability. As resulting pieces of data, such a probabilistic data integration process necessarily produces probabilistic facts, i.e., probabilistic knowledge at the instance level. Such probabilistic instance relationships can be encoded in probabilistic ontologies and there be enhanced by further classical and/or terminological probabilistic knowledge, which then allows for inference about other probabilistic instance relationships.

Although there are many previous approaches to probabilistic description logics and probabilistic ontology languages in the literature, including some that are specifically designed for the Semantic Web, there is only little work on tractable probabilistic description logics (see Section 7), and to date no work on tractable probabilistic description logics for the Semantic Web. In this paper, we try to fill this gap. We present a novel combination of description logics with probabilistic uncertainty, which is especially directed towards tractable formalisms for reasoning under probabilistic uncertainty with ontologies in the Semantic Web. More concretely, we present an extension of the *DL-Lite* family of description logics [2] by probabilistic uncertainty as in Bayesian networks. The main contributions of this paper can be summarized as follows:

- We present a probabilistic generalization of the *DL-Lite* family of description logics, which is based on Bayesian networks.

- We show that the new probabilistic description logics are rich enough to properly extend both the *DL-Lite* description logics as well as Bayesian networks.
- We also show that satisfiability checking and query processing in the new logics can be reduced to satisfiability checking and query processing in the *DL-Lite* family.
- Finally, we show that satisfiability checking and answering unions of conjunctive queries in the new logics can be done in LogSpace in the data complexity.

Compared to previous tractable probabilistic description logics, our new approach to tractable probabilistic description logics is especially well-suited for data-intensive applications in the Semantic Web, such as the ones listed above (see also Section 7).

The rest of this paper is organized as follows. In the next section, the preliminaries of the *DL-Lite* family of description logics and of Bayesian networks are presented. In Section 3, we introduce our new probabilistic description logics. Sections 4 to 6 provide semantic, computational, and data tractability results, respectively, around the new logics. In Section 7, we survey related work in neighboring research areas. Finally, Section 8 concludes the paper and outlines possible future directions of research. Note that detailed proofs of all results in this paper are given in the extended report.

## 2   Preliminaries

In this section, we first recall the main concepts of the *DL-Lite* family of tractable description logics, and we then recall the basics of Bayesian networks.

### 2.1   The *DL-Lite* Family

We now recall the *DL-Lite* family of tractable description logics [2], which include the core language *DL-Lite*$_{\text{core}}$ and its extensions *DL-Lite* (also called *DL-Lite*$_\mathcal{F}$) and *DL-Lite*$_\mathcal{R}$. They are a restricted class of classical description logics for which the main reasoning tasks in description logics can be done in deterministic polynomial time in the size of the knowledge base and some of these tasks even in LogSpace in the size of the ABox in the data complexity. The *DL-Lite* description logics are the most common tractable description logics in the semantic web context. They are especially directed towards data-intensive applications. We now first preliminarily recall the language and its semantics, and we then recall tractability results.

*Syntax.*  We first define *DL-Lite* (also called *DL-Lite*$_\mathcal{F}$). Let **A**, **R**$_A$, and **I** be pairwise disjoint sets of atomic concepts, abstract roles, and individuals, respectively.

A *basic role (in DL-Lite)* is either an atomic role $P \in \mathbf{R}_A$ or its inverse $P^-$. *Roles (in DL-Lite)* are defined as follows. Every basic role $P$ and negation of a basic role $\neg P$ is a role. A *basic concept (in DL-Lite)* is either an atomic concept from **A** or an existential restriction on a basic role $R$, denoted $\exists R.\top$ (abbreviated as $\exists R$). *Concepts (in DL-Lite)* are defined as follows. Every basic concept $B$ and negation of a basic concept $\neg B$ is a concept.

An *axiom (in DL-Lite)* is either (1) a *concept inclusion axiom* $B \sqsubseteq \phi$, where $B$ is a basic concept, and $\phi$ is a concept, or (2) a *functionality axiom* (funct $R$), where $R$ is a basic role, or (3) a *concept membership axiom* $A(a)$, where $A$ is an atomic concept

and $a \in \mathbf{I}$, or (4) a *role membership axiom* $R(a, c)$, where $R \in \mathbf{R}_A$ and $a, c \in \mathbf{I}$. A *TBox (in DL-Lite)* is a finite set of concept inclusion and functionality axioms. An *ABox (in DL-Lite)* is a finite set of concept and role membership axioms. A *knowledge base (in DL-Lite) $KB = (\mathcal{T}, \mathcal{A})$* consists of a TBox $\mathcal{T}$ and an ABox $\mathcal{A}$. A *query* $\phi$ is an open formula of first-order logic with equalities. A *conjunctive query* is of the form $\exists \mathbf{y} \, \phi(\mathbf{x}, \mathbf{y})$, where $\phi$ is a conjunction of atoms and equalities with free variables $\mathbf{x}$ and $\mathbf{y}$. A *union of conjunctive queries* is of the form $\bigvee_{i=1}^{n} \exists \mathbf{y}_i \, \phi_i(\mathbf{x}, \mathbf{y}_i)$, where each $\phi_i$ is a conjunction of atoms and equalities with free variables $\mathbf{x}$ and $\mathbf{y}_i$.

The description logic *DL-Lite$_{\mathrm{core}}$* does not allow for functionality axioms in knowledge bases, while *DL-Lite$_{\mathcal{R}}$* allows for (5) *role inclusion axioms* $R \sqsubseteq E$, rather than functionality axioms, where $R$ is a basic role, and $E$ is a role.

The following example from semantic web services illustrates the above notions.

*Example 1 (Flight Services).* Given an ontology as a shared knowledge base, we use description logic concepts to describe semantic web services, and their instances to represent the real procedures implementing the services (see [11] for more details).

More specifically, we consider flight services. The following knowledge base $KB = (\mathcal{T}, \mathcal{A})$ in *DL-Lite$_{\mathcal{R}}$* encodes an ontology with airports and air connections between them (where conjunctions are used to compactly represent several concept inclusion axioms with the same body by one concept inclusion axiom):

$\mathcal{T} = \{$*Service* $\sqsubseteq$ *Top*; *Airport* $\sqsubseteq$ *Top*; *Country* $\sqsubseteq$ *Top*;
　　*Service* $\sqsubseteq \neg$*Airport* $\sqcap \neg$*Country*; *Airport* $\sqsubseteq \neg$*Country*;
　　*Italy* $\sqsubseteq$ *Country*; *Germany* $\sqsubseteq$ *Country*; *UK* $\sqsubseteq$ *Country*;
　　*Italy* $\sqsubseteq \neg$*Germany* $\sqcap \neg$*UK*; *Germany* $\sqsubseteq \neg$*UK*;
　　*Rome* $\sqsubseteq$ *Airport*; *Cologne* $\sqsubseteq$ *Airport*; *Frankfurt* $\sqsubseteq$ *Airport*; *London* $\sqsubseteq$ *Airport*;
　　*Rome* $\sqsubseteq \neg$*Cologne* $\sqcap \neg$*Frankfurt* $\sqcap \neg$*London*; $\ldots$; *Frankfurt* $\sqsubseteq \neg$*London*;
　　*RomLon* $\sqsubseteq$ *Service*; *CgnLon* $\sqsubseteq$ *Service*; *FraLon* $\sqsubseteq$ *Service*;
　　*RomLon* $\sqsubseteq \neg$*CgnLon* $\sqcap \neg$*FraLon*; *CgnLon* $\sqsubseteq \neg$*FraLon*;
　　*FraLgw* $\sqsubseteq$ *FraLon*; *FraLhr* $\sqsubseteq$ *FraLon*; *FraLgw* $\sqsubseteq \neg$*FraLhr*;
　　*Service* $\sqsubseteq \exists$*From*; *Airport* $\sqsubseteq \exists$*From$^-$*;
　　*Service* $\sqsubseteq \exists$*To*; *Airport* $\sqsubseteq \exists$*To$^-$* $\}$,

$\mathcal{A} = \{$*Rome(FCO)*; *Rome(CIA)*; *Cologne(CGN)*; *Frankfurt(FRA)*; *London(LHR)*;
　　*FraLon(LH456)*; *CgnLon(GermanWings123)*; *RomLon(BA789)*;
　　*From(LH456,FRA)*; *From(GermanWings123,CGN)*; *From(BA789,FCO)*;
　　*To(LH456,LHR)*; *To(GermanWings123,LHR)*; *To(BA789,LHR)* $\}$.

In particular, the concepts *FraLon*, *CgnLon*, and *RomeLon* describe flight services from Frankfurt, Cologne, and Rome, respectively, to London. For each such concept, also an instance is specified. The above TBox $\mathcal{T}$ is partially illustrated in Fig. 1.

The union of conjunctive queries $Q(x) = \exists y(To(x, y) \wedge Rome(y)) \vee \exists y(From(x, y) \wedge Rome(y))$ then asks for all flight services that are ending or starting in Rome.

*Semantics.* An *interpretation* $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ consists of a nonempty (*abstract*) *domain* $\Delta^{\mathcal{I}}$ and a mapping $\cdot^{\mathcal{I}}$ that assigns to each atomic concept $C \in \mathbf{A}$ a subset of $\Delta^{\mathcal{I}}$, to each abstract role $R \in \mathbf{R}_A$ a subset of $\Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$, and to each individual $a \in \mathbf{I}$ an element
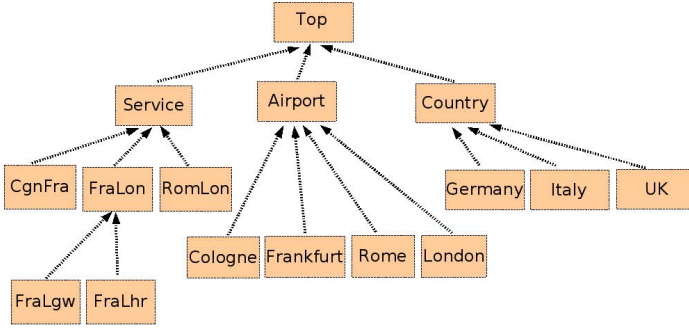
**Fig. 1.** TBox (partially) for the Flight Services Example

of $\Delta^{\mathcal{I}}$. Here, different individuals are associated with different elements of $\Delta^{\mathcal{I}}$ (*unique name assumption*). The mapping $\cdot^{\mathcal{I}}$ is extended to all concepts and roles as follows:

- $(R^{-})^{\mathcal{I}} = \{(a, b) \mid (b, a) \in R^{\mathcal{I}}\}$;
- $(\neg R)^{\mathcal{I}} = \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}} \setminus R^{\mathcal{I}}$;
- $(\exists R)^{\mathcal{I}} = \{x \in \Delta^{\mathcal{I}} \mid \exists y \colon (x, y) \in R^{\mathcal{I}}\}$;
- $(\neg B)^{\mathcal{I}} = \Delta^{\mathcal{I}} \setminus B^{\mathcal{I}}$.

The *satisfaction* of an axiom $F$ in the interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$, denoted $\mathcal{I} \models F$, is defined as follows: (1) $\mathcal{I} \models B \sqsubseteq \phi$ iff $B^{\mathcal{I}} \subseteq \phi^{\mathcal{I}}$; (2) $\mathcal{I} \models (\text{funct } R)$ iff $(o, o') \in R^{\mathcal{I}}$ and $(o, o'') \in R^{\mathcal{I}}$ implies $o' = o''$; (3) $\mathcal{I} \models A(a)$ iff $a^{\mathcal{I}} \in A^{\mathcal{I}}$; (4) $\mathcal{I} \models R(a, b)$ iff $(a^{\mathcal{I}}, b^{\mathcal{I}}) \in R^{\mathcal{I}}$; and (5) $\mathcal{I} \models R \sqsubseteq E$ iff $R^{\mathcal{I}} \subseteq E^{\mathcal{I}}$. The interpretation $\mathcal{I}$ *satisfies* the axiom $F$, or $\mathcal{I}$ is a *model* of $F$, iff $\mathcal{I} \models F$. The interpretation $\mathcal{I}$ *satisfies* a knowledge base $KB = (\mathcal{T}, \mathcal{A})$, or $\mathcal{I}$ is a *model* of $KB$, denoted $\mathcal{I} \models KB$, iff $\mathcal{I} \models F$ for all $F \in \mathcal{T} \cup \mathcal{A}$. We say that $KB$ is *satisfiable* (resp., *unsatisfiable*) iff $KB$ has a (resp., no) model. An axiom $F$ is a *logical consequence* of $KB$, denoted $KB \models F$, iff every model of $KB$ satisfies $F$. An *answer* for a query $\phi$ to $KB$ is a ground substitution $\theta$ for all free variables in $\phi$ such that $\phi\theta$ is a logical consequence of $KB$.

*Example 2 (Flight Services (cont'd)).* Consider again the knowledge base $KB = (\mathcal{T}, \mathcal{A})$ in *DL-Lite$_{\mathcal{R}}$* from Example 1. It is not difficult to verify that $KB$ is satisfiable, and that some logical consequences of $KB$ are given by *FraLhr* $\sqsubseteq$ *Service* and *Service*(*LH456*). The only answer for the query $Q(x)$ of Example 1 to $KB$ is $\theta = \{x/BA789\}$.

*Tractability.* We briefly recall the tractability results for reasoning with *DL-Lite* (resp., *DL-Lite$_{\mathcal{R}}$*) that we will use in the probabilistic generalization. The following result from [2] shows that deciding the satisfiability of *DL-Lite* (resp., *DL-Lite$_{\mathcal{R}}$*) knowledge bases can be done in LogSpace in the size of the ABox in the data complexity.

**Theorem 1 (see [2]).** *Given a DL-Lite (resp., DL-Lite$_{\mathcal{R}}$) knowledge base $KB = (\mathcal{T}, \mathcal{A})$, deciding whether $KB$ is satisfiable can be done in LogSpace in the size of the ABox $\mathcal{A}$ in the data complexity.*

The next result from [2] shows that computing the answers for unions of conjunctive queries to *DL-Lite* (resp., *DL-Lite$_\mathcal{R}$*) knowledge bases can also be done in LogSpace in the size of the ABox in the data complexity.

**Theorem 2 (see [2]).** *Given a DL-Lite (resp., DL-Lite$_\mathcal{R}$) knowledge base $KB = (\mathcal{T}, \mathcal{A})$ and a union of conjunctive queries $Q = \bigvee_{i=1}^{n} \exists \mathbf{y}_i\, \phi_i(\mathbf{x}, \mathbf{y}_i)$, computing all answers for Q to KB can be done in LogSpace in the size of the ABox $\mathcal{A}$ in the data complexity.*

## 2.2 Bayesian Networks

We now briefly recall *Bayesian networks* (see especially [25, 15]). Let $V$ be a finite set of *random variables*. Each variable $X \in V$ may take on *values* from a finite *domain* $D(X)$. A *value* for a set of variables $X = \{X_1, \ldots, X_n\} \subseteq V$ is a mapping $x \colon X \to \bigcup_{i=1}^{n} D(X_i)$ such that $x(X_i) \in D(X_i)$ (where the empty mapping $\emptyset$ is the unique value for $X = \emptyset$). The *domain* of $X$, denoted $D(X)$, is the set of all values for $X$. For $Y \subseteq X$ and $x \in D(X)$, we use $x|Y$ to denote the restriction of $x$ to $Y$. We often identify singletons $\{X_i\} \subseteq V$ with $X_i$, and their values $x$ with $x(X_i)$.

A *Bayesian network* $BN = (G, Pr)$ over $V$ is defined by a directed acyclic graph $G = (V, E)$ over the random variables in $V$ as nodes and by a conditional probability distribution $Pr(X = \cdot \mid Y = y) \colon D(X) \to [0, 1]$ for each variable $X \in V$ and each value $y \in D(Y)$ of the parents $Y \subseteq V$ of $X$ in $G$, denoted $\mathrm{Pa}(X)$. It specifies a unique joint probability distribution $Pr_{BN}$ over all values for $V$ by:

$$Pr_{BN}(V = v) = \prod_{X \in V} Pr(X = v|X \mid \mathrm{Pa}(X) = v|\mathrm{Pa}(X)) \quad \text{(for every } v \in D(V)).$$

That is, the joint probability distribution $Pr_{BN}$ is uniquely determined by the conditional probability distributions $Pr(X = \cdot \mid Y = y)$. This implicitly assumes conditional probabilistic independencies encoded in the directed acyclic graph $G$. One then specifies a probability $Pr_{BN}$ for every $X \subseteq V$ and $x \in D(X)$ as follows:

$$Pr_{BN}(X = x) = \sum_{v \in D(V),\, v|X=x} Pr_{BN}(V = v).$$

# 3 Bayesian *DL-Lite$_\mathcal{R}$* (*BDL-Lite$_\mathcal{R}$*)

In this section, we introduce the novel probabilistic description logic *Bayesian DL-Lite$_\mathcal{R}$* (or *BDL-Lite$_\mathcal{R}$*), which combines classical knowledge bases in *DL-Lite$_\mathcal{R}$* with Bayesian networks. Informally, every description logic axiom is annotated with an event, which is in turn associated with a probability value via a Bayesian network. Like *DL-Lite$_\mathcal{R}$*, *BDL-Lite$_\mathcal{R}$* is especially directed towards data-intensive applications. Note that a very similar probabilistic generalization can be defined for *DL-Lite$_\mathcal{F}$*.

## 3.1 Syntax

We first define the syntax of *BDL-Lite$_\mathcal{R}$*. As for the elementary ingredients, as in Section 2.1, let $\mathbf{A}$, $\mathbf{R}_A$, and $\mathbf{I}$ be pairwise disjoint sets of atomic concepts, abstract roles, and individuals, respectively. As in Section 2.2, we assume a finite set of random variables $V$, where each $X \in V$ may take on values from a finite domain $D(X)$.

We next define the concept of a probabilistic knowledge base, which consists of a set of probabilistic axioms and a Bayesian network. Every probabilistic axiom in turn consists of a classical axiom in *DL-Lite*$_{\mathcal{R}}$ and a probabilistic annotation, which connects it to a set of value assignments $V = v$ with $v \in D(V)$ of a Bayesian network over $V$ along with their probability values. Formally, a *probabilistic annotation* is an expression of the form $X = x$, where $X \subseteq V$ and $x \in D(X)$. We also use $\top$ to denote the probabilistic annotation for $X = \emptyset$. Informally, every probabilistic annotation represents a scenario (or an event) which is associated with the set of all value assignments $V = v$ with $v \in D(V)$ that are compatible with $X = x$ (that is, $v | X = x$) and their probability value $Pr_{BN}(V = v)$ in a Bayesian network $BN$ over $V$. We next formally define *probabilistic axioms* as follows. A *probabilistic concept membership* (resp., *role membership*, *concept inclusion*, *functionality*, *role inclusion*) *axiom in BDL-Lite*$_{\mathcal{R}}$ is an expression of the form $\phi\colon X = x$, where $\phi$ is a concept membership (resp., role membership, concept inclusion, functionality, role inclusion) axiom in *DL-Lite*$_{\mathcal{R}}$, and $X = x$ is a probabilistic annotation. Informally, such a probabilistic axiom $\phi\colon X = x$ encodes that in the scenario $X = x$, the description logic axiom $\phi$ holds. We often abbreviate probabilistic axioms of the form $\top\colon X = x$ (resp., $\phi\colon \top$) by $X = x$ (resp., $\phi$). A *probabilistic TBox in BDL-Lite*$_{\mathcal{R}}$ is a finite set of probabilistic concept inclusion and probabilistic role inclusion axioms in *BDL-Lite*$_{\mathcal{R}}$. A *probabilistic ABox in BDL-Lite*$_{\mathcal{R}}$ is a finite set of probabilistic concept and probabilistic role membership axioms in *BDL-Lite*$_{\mathcal{R}}$. A *probabilistic knowledge base $KB = (\mathcal{T}, \mathcal{A}, BN)$ in BDL-Lite*$_{\mathcal{R}}$ consists of (i) a probabilistic TBox $\mathcal{T}$ in *BDL-Lite*$_{\mathcal{R}}$, (ii) a probabilistic ABox $\mathcal{A}$ in *BDL-Lite*$_{\mathcal{R}}$, and (iii) a Bayesian network $BN = ((V, E), Pr)$.

We finally define probabilistic queries to probabilistic knowledge bases in *BDL-Lite*$_{\mathcal{R}}$. A *probabilistic query* is of the form $\psi\colon X = x$, where $\psi$ is a first-order formula, and $X = x$ is a probabilistic annotation. We often abbreviate probabilistic queries of the form $\top\colon X = x$ (resp., $\psi\colon \top$) by $X = x$ (resp., $\psi$). A *probabilistic union of conjunctive queries* is a probabilistic query $\psi\colon X = x$ such that $\psi$ is a union of conjunctive queries.

*Example 3 (Flight Services (cont'd)).* Consider again the knowledge base $KB = (\mathcal{T}, \mathcal{A})$ in *DL-Lite*$_{\mathcal{R}}$ given in Example 1. We may now know that, given that a service belongs to a concept, then it also belongs to another concept with a certain probability. For example, we may know that a service in *FraLon* is also a service in *FraLhr* with the probability 0.8. This probabilistic information may be useful, for example, when searching for services in *FraLon*, to speed up the service discovery process [6]. It can be encoded by the following probabilistic concept inclusion axiom:

$$FraLon \sqsubseteq FraLhr\colon LonLhr = \mathbf{true}\,,$$

where *LonLhr* is a random variable, which is true with the probability 0.8. Similarly, functionality axioms and role inclusion axioms can be annotated with probabilities.

In the same way, we may know that the individual *GermanWings456* is an instance of the service description *FraLon* with the probability 0.9, which can be expressed by the following probabilistic concept membership axiom:

$$FraLon(GermanWings456)\colon FraLonGW = \mathbf{true}\,,$$

where *FraLonGW* is a random variable, which is true with the probability 0.9.
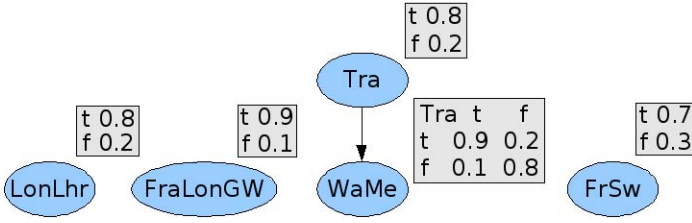
**Fig. 2.** Bayesian network (with CPTs) for the Flight Services Example

Similarly, we can express that (1) the individual *Swiss123* belongs to the concept *Transatlantic* with the probability 0.8, (2) the individual *Swiss123* belongs to the concept *WarmMeal*, given that it belongs (resp., does not belong) to the concept *Transatlantic* with the probability 0.9 (resp., 0.2), and (3) the pair of individuals (*Swiss123*, *CIA*) belongs to the role *From* with the probability 0.7 by the following probabilistic concept membership axioms and probabilistic role membership axiom:

$$Transatlantic(Swiss123)\colon Tra = \textbf{true},$$
$$WarmMeal(Swiss123)\colon WaMe = \textbf{true},$$
$$From(Swiss123,CIA)\colon FrSw = \textbf{true},$$

where (1) *Tra* is a random variable, which is true with the probability 0.8, (2) *WaMe* is a random variable, which is true with the probability 0.9 (resp., 0.2), given *Tra* is true (resp., false), and (3) *FrSw* is a random variable, which is true with the probability 0.7.

The above random variables along with their probabilities and conditional probabilities form a Bayesian network, which is shown in Fig. 2, where the probabilities and conditional probabilities are represented in conditional probability tables (CPTs).

The probabilistic union of conjunctive queries $Q(x) = \exists y (To(x, y) \land Rome(y)) \lor \exists y (From(x, y) \land Rome(y))$ then asks for all flight services that are ending or starting in Rome, along with their probabilities. Whereas the probabilistic conjunctive query $Q'(x) = \exists y (From(x, y) \land Rome(y) \land WarmMeal(x))$ asks for all flight services that are starting in Rome and are offering a warm meal, along with their probabilities.

From the engineering viewpoint, there are two different ways of designing probabilistic knowledge bases $KB = (\mathcal{T}, \mathcal{A}, BN)$ in *BDL-Lite*$_\mathcal{R}$. One is to start to model the Bayesian network $BN = ((V, E), Pr)$, and to collect for every probabilistic annotation $V = v$ with $v \in D(V)$ a set of probabilistic axioms $\phi\colon V = v$, which is then simplified to a set of probabilistic axioms of the type $\phi\colon X = x$ with $X \subseteq V$ and $x \in D(X)$. Another way is to start to model a set of probabilistic axioms $\phi\colon X = x$ with single binary random variables $X$, which are then used to form the nodes of the Bayesian network $BN = ((V, E), Pr)$. In this paper, we adopt especially the first viewpoint.

## 3.2   Semantics

We now define a formal semantics of probabilistic knowledge bases in *BDL-Lite*$_\mathcal{R}$, in terms of probability distributions over classical interpretations.

We first define annotated interpretations, which extend standard first-order interpretations (under the unique name assumption) by value assignments $V = v$ in a Bayesian network over $V$. Formally, an *annotated interpretation* $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ is defined in the same way as a classical first-order interpretation under the unique name assumption (see Section 2.1) except that $\cdot^{\mathcal{I}}$ maps additionally the set of random variables $V$ to a value $v \in D(V)$. The annotated interpretation $\mathcal{I}$ *satisfies* (or is a *model* of) a probabilistic axiom $\phi\colon X = x$, denoted $\mathcal{I} \models \phi\colon X = x$, iff $V^{\mathcal{I}}|X = x$ is equivalent to $\mathcal{I} \models \phi$.

We next define probabilistic interpretations, which are finite probability distributions over annotated interpretations. Formally, a *probabilistic interpretation* $Pr$ is a probability function over the set of all annotated interpretations that associates only a finite number of annotated interpretations with a positive probability. The *probability* of a probabilistic axiom $\phi\colon X = x$ in $Pr$, denoted $Pr(\phi\colon X = x)$, is the sum of all $Pr(\mathcal{I})$ such that $\mathcal{I}$ is an annotated interpretation that satisfies $\phi\colon X = x$. A probabilistic interpretation $Pr$ *satisfies* (or is a *model* of) a probabilistic axiom $\phi\colon X = x$ iff $Pr(\phi\colon X = x) = 1$. We say $Pr$ *satisfies* (or is a *model* of) a set of probabilistic axioms $\mathcal{F}$ iff $Pr$ satisfies all $F \in \mathcal{F}$. The probabilistic interpretation $Pr$ *satisfies* (or is a *model* of) a probabilistic knowledge base $KB = (\mathcal{T}, \mathcal{A}, BN)$ in *BDL-Lite$_{\mathcal{R}}$* iff (i) $Pr$ is a model of $\mathcal{T} \cup \mathcal{A}$ and (ii) $Pr(V = v) = Pr_{BN}(V = v)$ for all $v \in D(V)$. We say $KB$ is *satisfiable* iff it has a model $Pr$.

We finally define answers for probabilistic queries as follows. An annotated interpretation $\mathcal{I}$ *satisfies* (or is a *model* of) a ground query $\psi\colon X = x$, denoted $\mathcal{I} \models \psi\colon X = x$, iff $V^{\mathcal{I}}|X = x$ and $\mathcal{I} \models \psi$. The *probability* of a ground query $\psi\colon X = x$ in $Pr$, denoted $Pr(\psi\colon X = x)$, is the sum of all $Pr(\mathcal{I})$ such that $\mathcal{I}$ is an annotated interpretation that satisfies $\psi\colon X = x$. An *answer* for a probabilistic query $Q = \psi\colon X = x$ to a probabilistic knowledge base $KB = (\mathcal{T}, \mathcal{A}, BN)$ is a pair $(\theta, pr)$ consisting of a ground substitution $\theta$ for the variables in $Q$ and some $pr \in [0, 1]$ such that $Pr(\psi\theta\colon X = x) = pr$ for all models $Pr$ of $KB$. An answer $(\theta, pr)$ for $Q$ to $KB$ is *positive* iff $pr > 0$.

*Example 4 (Flight Services (cont'd)).* Consider again the probabilistic knowledge base $KB$ in *BDL-Lite$_{\mathcal{R}}$* and the two probabilistic queries $Q(x)$ and $Q'(x)$ described in Example 3. It is not difficult to verify that $KB$ is satisfiable. The only positive answers $(\theta, pr)$ for $Q(x)$ to $KB$ are $(\{x/BA789\}, 1)$ and $(\{x/Swiss123\}, 0.7)$, while the only positive answer $(\theta, pr)$ for $Q'(x)$ to $KB$ is $(\{x/Swiss123\}, 0.76)$. If $KB$ would additionally contain the probabilistic concept membership axiom *Transatlantic(Swiss123)*, then the only positive answer $(\theta, pr)$ for $Q'(x)$ to $KB$ would be $(\{x/Swiss123\}, 0.9)$.

## 4   Semantic Properties

An important property of hybrid knowledge representation and reasoning formalisms is that they faithfully extend their integrated formalisms. In this section, we show that *BDL-Lite$_{\mathcal{R}}$* faithfully extends both *DL-Lite$_{\mathcal{R}}$* and Bayesian networks.

The following theorem shows that probabilistic knowledge bases in *BDL-Lite$_{\mathcal{R}}$* faithfully extend Bayesian networks. That is, querying any Bayesian network is equivalent to querying any of its extensions to a satisfiable probabilistic knowledge base in *BDL-Lite$_{\mathcal{R}}$*.

**Theorem 3.** *Let $BN = ((V, E), Pr)$ be a Bayesian network, and let $KB = (\mathcal{T}, \mathcal{A}, BN')$ be any probabilistic knowledge base in BDL-Lite$_\mathcal{R}$ such that $BN' = BN$. Let $X \subseteq V$ and $x \in D(X)$. Then, the probabilistic query $Q = X = x$ to KB has the pair $(\theta, pr) = (\emptyset, Pr_{BN}(X = x))$ as an answer. If KB is satisfiable, then this pair $(\theta, pr)$ is also the only answer for Q to KB.*

We next show that probabilistic knowledge bases in *BDL-Lite$_\mathcal{R}$* also faithfully extend classical knowledge bases in *DL-Lite$_\mathcal{R}$*. In detail, querying any satisfiable knowledge base in *DL-Lite$_\mathcal{R}$* is equivalent to querying any of its extensions in *BDL-Lite$_\mathcal{R}$*.

**Theorem 4.** *Let $KB = (\mathcal{T}, \mathcal{A})$ be a satisfiable knowledge base in DL-Lite$_\mathcal{R}$, let $\psi$ be a query to KB, and let $BN = ((V, E), Pr)$ be any Bayesian network. Then, the probabilistic query $Q = \psi$ to $KB' = (\mathcal{T}, \mathcal{A}, BN)$ has as positive answers $(\theta, pr)$ exactly all pairs $(\theta, 1)$ such that $\theta$ is an answer for $\psi$ to KB.*

## 5    Computation

In this section, we show that satisfiability checking and query processing in *BDL-Lite$_\mathcal{R}$* can be reduced to satisfiability checking and query processing in *DL-Lite$_\mathcal{R}$*.

The following theorem shows that the satisfiability of probabilistic knowledge bases in *BDL-Lite$_\mathcal{R}$* can be reduced to the satisfiability of knowledge bases in *DL-Lite$_\mathcal{R}$*. Note that all negated axioms in the theorem can be simulated by positive ones.

**Theorem 5.** *Let $KB = (\mathcal{T}, \mathcal{A}, BN)$ be a probabilistic knowledge base in BDL-Lite$_\mathcal{R}$. For every $v \in D(V)$, let $\mathcal{T}_v$ (resp., $\mathcal{A}_v$) be the set of all axioms $\phi$ and $\neg\phi$ for which there exists a probabilistic axiom $\phi\colon X = x$ in $\mathcal{T}$ (resp., $\mathcal{A}$), such that $v|X = x$ and $v|X \neq x$, respectively. Then, KB is satisfiable iff the knowledge base $KB_v = (\mathcal{T}_v, \mathcal{A}_v)$ in DL-Lite$_\mathcal{R}$ is satisfiable for every $v \in D(V)$ with $Pr_{BN}(V = v) > 0$.*

The next theorem shows that query processing in probabilistic knowledge bases in *BDL-Lite$_\mathcal{R}$* can be reduced to query processing in knowledge bases in *DL-Lite$_\mathcal{R}$*. Note that all negated axioms in the theorem can be simulated by positive ones.

**Theorem 6.** *Let $KB = (\mathcal{T}, \mathcal{A}, BN)$ be a satisfiable probabilistic knowledge base in BDL-Lite$_\mathcal{R}$, and let $Q = \psi\colon X = x$ be a probabilistic query to KB. For every $v \in D(V)$, let $\mathcal{T}_v$ (resp., $\mathcal{A}_v$) be the set of all $\phi$ and $\neg\phi$ for which there exists some $\phi\colon X = x$ in $\mathcal{T}$ (resp., $\mathcal{A}$) such that $v|X = x$ and $v|X \neq x$, respectively. Let $\theta$ be a ground substitution for the variables in Q and let $pr \in (0, 1]$. Then, $(\theta, pr)$ is an answer for Q to KB iff pr is the sum of all $Pr_{BN}(V = v)$ such that (i) $v \in D(V)$ with $Pr_{BN}(V = v) > 0$, (ii) $\theta$ is an answer for $\psi$ to $KB_v = (\mathcal{T}_v, \mathcal{A}_v)$, and (iii) $v|X = x$.*

## 6    Tractability Results

As an important result of this paper, we now show that both satisfiability checking and query processing in *BDL-Lite$_\mathcal{R}$* can be done in LogSpace in the data complexity. Note that we adopt the notion of data complexity from logic programming [5].

The following theorem shows that deciding whether a probabilistic knowledge base in *BDL-Lite$_\mathcal{R}$* is satisfiable can be done in LogSpace in the data complexity.

**Theorem 7.** *Given a probabilistic knowledge base $KB = (\mathcal{T}, \mathcal{A}, BN)$ in BDL-Lite$_\mathcal{R}$, deciding whether KB is satisfiable can be done in LogSpace in the size of $\mathcal{A}$ in the data complexity.*

The next theorem shows that computing all positive answers for probabilistic unions of conjunctive queries can also be done in LogSpace in the data complexity.

**Theorem 8.** *Given a satisfiable probabilistic knowledge base $KB = (\mathcal{T}, \mathcal{A}, BN)$ in BDL-Lite$_\mathcal{R}$ and a probabilistic union of conjunctive queries $Q = \psi \colon X = x$, computing the set of all positive answers $(\theta, pr)$ for Q to KB can be done in LogSpace in the size of the ABox $\mathcal{A}$ in the data complexity.*

## 7    Related Work

There are several related approaches to probabilistic description logics in the literature, which can be classified according to the generalized description logics, the supported forms of probabilistic knowledge, and the underlying probabilistic reasoning.

Closest in spirit to this paper is perhaps the work by Koller et al. [17], which presents P-CLASSIC, which is a probabilistic generalization of the CLASSIC description logic, rather than the *DL-Lite* family. Like our approach, theirs allows for terminological probabilistic knowledge about concepts and roles, but unlike ours, theirs does not support assertional knowledge about instances of concepts and roles. Like ours, their approach is based on inference in Bayesian networks as underlying probabilistic reasoning formalism. Closely related work by Yelland [29] combines a restricted description logic close to $\mathcal{FL}$ with Bayesian networks, rather than the *DL-Lite* family. It also allows for terminological probabilistic knowledge about concepts and roles, but does not support assertional knowledge about instances of concepts and roles. The main differences to our work are summarized as follows. First, we allow for both terminological probabilistic knowledge about concepts and roles, and assertional knowledge about instances of concepts and roles. Second, as a closely related aspect, unlike the above two works, we provide LogSpace data complexity results, and we consider the problem of answering probabilistic unions of conjunctive queries. Third, the above two probabilistic description logics essentially lie in the intersection of tractable description logics and Bayesian networks, and are thus limited in their expressive power, while ours orthogonally and faithfully combine the two components, and thus keep their expressive power. For this reason, our approach allows for much richer terminological knowledge. Hence, compared to previous tractable probabilistic description logics, our new approach to tractable probabilistic description logics is especially well-suited for data-intensive applications in the Semantic Web, such as the ones listed in the introduction.

Also closely related are the probabilistic description logics in [20], which are probabilistic extensions of the expressive description logics $\mathcal{SHIF}(\mathbf{D})$ and $\mathcal{SHOIN}(\mathbf{D})$ behind OWL Lite and OWL DL, respectively, towards sophisticated formalisms for reasoning under probabilistic uncertainty in the Semantic Web.[1] They allow for expressing

---

[1] See [16] for an implementation of the probabilistic description logics in [20].

both terminological probabilistic knowledge about concepts and roles, and also assertional probabilistic knowledge about instances of concepts and roles. Our present work is more flexible in the sense that terminological and assertional pieces of probabilistic knowledge can be freely combined, while [20] partitions the probabilistic knowledge into terminological pieces of probabilistic knowledge and object-centered assertional pieces of probabilistic knowledge. Rather than on Bayesian networks, they are based on probabilistic lexicographic entailment from probabilistic default reasoning [19] as underlying probabilistic reasoning formalism, which treats terminological and assertional probabilistic knowledge in a semantically way as probabilistic knowledge about random resp. concrete instances. Differently from [20], we here provide LogSpace data complexity results, and we consider probabilistic unions of conjunctive queries.

Heinsohn [12] presents a probabilistic extension of $\mathcal{ALC}$, which allows to represent terminological probabilistic knowledge about concepts and roles, and which is essentially based on probabilistic reasoning in probabilistic logics, similar to [23, 18]. Heinsohn, however, does not allow for assertional knowledge about concept and role instances. Jaeger's work [13] proposes another probabilistic extension of $\mathcal{ALC}$, which allows for terminological and assertional probabilistic knowledge about concepts / roles and about concept instances, respectively, but does not support assertional probabilistic knowledge about role instances (although a possible extension in this direction is mentioned). The uncertain reasoning formalism in [13] is essentially based on probabilistic reasoning in probabilistic logics, as the one in [12], but coupled with cross-entropy minimization to combine terminological probabilistic knowledge with assertional probabilistic knowledge. Jaeger's recent work [14] is less closely related, as it focuses on interpreting probabilistic concept subsumption and probabilistic role quantification through statistical sampling distributions, and develops a probabilistic version of the guarded fragment of first-order logic.

Related works on probabilistic web ontology languages focus especially on combining the web ontology language OWL with probabilistic formalisms based on Bayesian networks. In particular, da Costa et al. [4, 3] propose a probabilistic generalization of OWL, called PR-OWL, which is based on multi-entity Bayesian networks.

Ding et al. [8] propose a probabilistic generalization of OWL, called BayesOWL, which is based on standard Bayesian networks. BayesOWL provides a set of rules and procedures for the direct translation of an OWL ontology into a Bayesian network that supports ontology reasoning, both within and across ontologies, as Bayesian inferences. The authors also describe an application of this approach in ontology mapping. In closely related work, Mitra et al. [22] introduce a technique to enhancing existing ontology mappings by using a Bayesian network to represent the influences between potential concept mappings across ontologies.

Yang and Calmet [28] present an integration of the web ontology language OWL with Bayesian networks. The approach makes use of probability and dependency-annotated OWL to represent uncertain information in Bayesian networks. Pool and Aikin [26] also provide a method for representing uncertainty in OWL ontologies, while Fukushige [9] proposes a basic framework for representing probabilistic relationships in RDF. Finally, Nottelmann and Fuhr [24] present two probabilistic extensions of variants of OWL Lite, along with a mapping to locally stratified probabilistic Datalog.

## 8   Summary and Outlook

We have presented probabilistic generalizations of the *DL-Lite* description logics, which are based on Bayesian networks. We have shown that the new probabilistic description logics properly extend both the *DL-Lite* description logics as well as Bayesian networks. We have also shown that satisfiability checking and query processing in the new probabilistic description logics can be reduced to satisfiability checking and query processing in the *DL-Lite* family. Furthermore, satisfiability checking and answering probabilistic unions of conjunctive queries can be done in LogSpace in the data complexity.

Other classical description logics can be extended similarly by probabilistic uncertainty in Bayesian networks. All results of this paper carry over to such extensions, except for the tractability results, which generally will not hold for extensions of classical description logics that are more expressive than those of the *DL-Lite* family.

We leave for future work the implementation of the new probabilistic description logics and the investigation of efficient algorithms for the general case beyond the data complexity (where tractable cases and efficient techniques from Bayesian networks may come into play). Another interesting topic for future research is to investigate the use of the new tractable probabilistic description logics in important tasks such as web search and database querying. Furthermore, it would be very interesting to develop techniques for learning the new tractable probabilistic description logics (e.g., from web data).

## References

[1] Baader, F., Calvanese, D., McGuinness, D., Nardi, D., Patel-Schneider, P.F. (eds.): The Description Logic Handbook: Theory, Implementation, and Applications. Cambridge University Press, Cambridge (2003)

[2] Calvanese, D., De Giacomo, G., Lembo, D., Lenzerini, M., Rosati, R.: DL-Lite: Tractable description logics for ontologies. In: Proceedings AAAI 2005, pp. 602–607. AAAI Press/MIT Press (2005)

[3] da Costa, P.C.G., Laskey, K.B.: PR-OWL: A framework for probabilistic ontologies. In: Proceedings FOIS 2006, pp. 237–249. IOS Press, Amsterdam (2006)

[4] da Costa, P.C.G., Laskey, K.B., Laskey, K.J.: PR-OWL: A Bayesian ontology language for the Semantic Web. In: Proceedings URSW 2005, pp. 23–33 (2005)

[5] Dantsin, E., Eiter, T., Gottlob, G., Voronkov, A.: Complexity and expressive power of logic programming. ACM Comput. Surv. 33(3), 374–425 (2001)

[6] d'Amato, C., Staab, S., Fanizzi, N., Esposito, F.: Efficient discovery of services specified in description logics languages. In: Proceedings of the ISWC-2007 Workshop on Service Matchmaking and Resource Retrieval in the Semantic Web (SMR$^2$ 2007) (2007)

[7] Ding, Z., Peng, Y.: A probabilistic extension to ontology language OWL. In: Proceedings HICSS 2004 (2004)

[8] Ding, Z., Peng, Y., Pan, R.: BayesOWL: Uncertainty modeling in Semantic Web ontologies. In: Ma, Z. (ed.) Soft Computing in Ontologies and Semantic Web. Studies in Fuzziness and Soft Computing, vol. 204, Springer, Heidelberg (2006)

[9] Fukushige, Y.: Representing probabilistic knowledge in the Semantic Web. In: Proceedings of the W3C Workshop on Semantic Web for Life Sciences, Cambridge, MA, USA (2004)

[10] Lukasiewicz, T., Giugno, R.: P-$\mathcal{SHOQ}(D)$: A Probabilistic Extension of $\mathcal{SHOQ}(D)$ for probabilistic ontologies in the Semantic Web. In: Flesca, S., Greco, S., Leone, N., Ianni, G. (eds.) JELIA 2002. LNCS (LNAI), vol. 2424, pp. 86–97. Springer, Heidelberg (2002)

[11] Grimm, S., Motik, B., Preist, C.: Variance in e-business service discovery. In: Proceedings of the ISWC 2004 Workshop on Semantic Web Services (2004)

[12] Heinsohn, J.: Probabilistic description logics. In: Proceedings UAI 1994, pp. 311–318. Morgan Kaufmann, San Francisco (1994)

[13] Jaeger, M.: Probabilistic reasoning in terminological logics. In: Proceedings KR 1994, pp. 305–316. Morgan Kaufmann, San Francisco (1994)

[14] Jaeger, M.: Probabilistic role models and the guarded fragment. In: Proc. IPMU 2004, pp. 235–242 (2004); Extended version in Int. J. Uncertain. Fuzz., 14(1), 43–60 (2006)

[15] Jensen, F.V.: Bayesian Networks and Decision Graphs. Springer, Heidelberg (2001)

[16] Klinov, P.: Pronto: A non-monotonic probabilistic description logic reasoner. In: System demo at ESWC 2008 (2008)

[17] Koller, D., Levy, A., Pfeffer, A.: P-Classic: A tractable probabilistic description logic. In: Proceedings AAAI 1997, pp. 390–397. AAAI Press/MIT Press (1997)

[18] Lukasiewicz, T.: Probabilistic deduction with conditional constraints over basic events. J. Artif. Intell. Res. 10, 199–241 (1999)

[19] Lukasiewicz, T.: Probabilistic default reasoning with conditional constraints. Ann. Math. Artif. Intell. 34(1–3), 35–88 (2002)

[20] Lukasiewicz, T.: Expressive probabilistic description logics. Artif. Intell. 172(6/7), 852–883 (2008)

[21] Lukasiewicz, T., Straccia, U.: Managing uncertainty and vagueness in description logics for the Semantic Web. J. Web Sem. (in press)

[22] Mitra, P., Noy, N.F., Jaiswal, A.: OMEN: A probabilistic ontology mapping tool. In: Gil, Y., Motta, E., Benjamins, V.R., Musen, M.A. (eds.) ISWC 2005. LNCS, vol. 3729, pp. 537–547. Springer, Heidelberg (2005)

[23] Nilsson, N.J.: Probabilistic logic. Artif. Intell. 28(1), 71–88 (1986)

[24] Nottelmann, H., Fuhr, N.: Adding probabilities and rules to OWL Lite subsets based on probabilistic Datalog. Int. J. Uncertain. Fuzz. 14(1), 17–42 (2006)

[25] Pearl, J.: Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Morgan Kaufmann, San Francisco (1988)

[26] Pool, M., Aikin, J.: KEEPER and Protégé: An elicitation environment for Bayesian inference tools. In: Proceedings of the Workshop on Protégé and Reasoning held at the 7th International Protégé Conference (2004)

[27] Smyth, C., Poole, D.: Qualitative probabilistic matching with hierarchical descriptions. In: Proceedings KR 2004, pp. 479–487. AAAI Press, Menlo Park (2004)

[28] Yang, Y., Calmet, J.: OntoBayes: An ontology-driven uncertainty model. In: Proceedings IAWTIC 2005, pp. 457–463. IEEE Computer Society Press, Los Alamitos (2005)

[29] Yelland, P.M.: An alternative combination of Bayesian networks and description logics. In: Proceedings KR 2000, pp. 225–234. Morgan Kaufmann, San Francisco (2000)