# Probabilistic Modeling and OWL: A User Oriented Introduction to P-$\mathcal{SHIQ}$(D)

Pavel Klinov and Bijan Parsia

School of Computer Science
The University of Manchester
Email: {pklinov,bparsia}@cs.man.ac.uk

**Abstract.** This paper presents a non-technical, user oriented introduction into P-$\mathcal{SHIQ}$(D)— an expressive formalism that allows modelers to incorporate probabilistic background knowledge in OWL ontologies. Instead of providing formal description of the language, its syntax, semantics, and reasoning procedures, we explain the basic principles, potentially useful practices and pitfalls of probabilistic modeling. The goal of the paper is to present P-$\mathcal{SHIQ}$(D) and the reasoning tools, such as Pronto, in an accessible form to encourage their usage in practical applications. We describe features and benefits of P-$\mathcal{SHIQ}$(D) using examples from the HealthCare and Life Sciences domain including the prototype of a probabilistic ontology for breast cancer risk assessment.

## 1 Introduction

Extensions of Description Logics (DLs) aimed at handling imprecision in ontologies have received significant research attention over the last several years, strongly driven by Bioinformatics and Semantic Web applications. There are now a variety of formalisms based on different mathematical approaches. There is a daunting number of technical papers containing complete discussions of the syntax, semantics, and computational properties of those formalisms.

This paper presents one of the approaches to probabilistic description logic — namely, P-$\mathcal{SHIQ}$(D) [GL02][Luk08] — but from a slightly different perspective. Instead of a presentation of the formal properties of P-$\mathcal{SHIQ}$(D), we examine it from a modeler's perspective. We attempt to show that P-$\mathcal{SHIQ}$(D) is a rather natural extension to OWL that allows ontology designers to incorporate statistical knowledge into OWL ontologies so that it can be used in applications.

P-$\mathcal{SHIQ}$(D) suffers from the standard problem of novel formalisms: its modeling principles are unclear, which makes people reluctant (or unable) to use it, which, in turn, prevents the development of a significant number of real P-$\mathcal{SHIQ}$(D) ontologies which is essential for the emergence of design principles. This paper tries to address this problem by providing the first realistic prototype of a P-$\mathcal{SHIQ}$(D) ontology and discussing its design and use. The ontology will represent statistical background knowledge about breast cancer and support breast cancer risk assessment (BRCA) of individual women.

The paper consists of two parts. We first provide a brief and informal overview of the representation and reasoning features that P-$\mathcal{SHIQ}$(D) adds to OWL. Then we proceed to the modeling issues illustrated through the process of constructing the BRCA ontology. We do so step by step highlighting principles that seem to be useful for building P-$\mathcal{SHIQ}$(D) ontologies.

The principles we have developed seem quite sensible and general, but the strength of our conclusions is necessarily limited. The problem we tackled is real with a real application, but we are not, ourself, domain experts. Thus, our recommendations are meant more to help people bootstrap themselves into using P-$\mathcal{SHIQ}$(D) than to provide the last word in good modeling style.

## 2 Probabilistic Language

First of all, P-$\mathcal{SHIQ}$(D) is a new ontological language built on top of OWL in order to represent probabilistic statements. The syntax is a simple extension of OWL: we add a new kind of axiom called *conditional constraints* which are responsible for the representation of all types of uncertainty in P-$\mathcal{SHIQ}$(D).

**Conditional Constraints** Conditional constraints are expressions of the form $(D|C)[l, u]$ where $C, D$ are OWL class expressions[1] called the *evidence* and the *conclusion* respectively, and $[l, u] \subseteq [0, 1]$ is a probability interval.

Informally, conditional constraints express probabilistic relationships in a domain. In spite of seeming similar to arcs in Bayesian, networks they can capture relationships between first order class expressions rather than propositional random variables. For example, we can represent that flying objects that have wings are birds in at least 80% of the cases: $(Bird|\exists hasWing.\top \sqcap FlyingObject)[0.8, 1]$

There are two types of conditional constraints: generic constraints that represent relationships between concepts (OWL classes) and individual constraints that represent relationships between an OWL class and an individual. Similarly to classical DL we will call the first type probabilistic TBox (PTBox) constraints and the second — probabilistic ABox (PABox) constraints. A probabilistic knowledge base (ontology) in P-$\mathcal{SHIQ}$(D) consists of three parts: an normal OWL ontology, a set of PTBox constraints, and a set of PABoxes with each one representing information about a distinct (probabilistic) individual, e.g.,:

*Example 1.*

$T = \{Penguin \sqsubseteq Bird\}$
$P = \{(FlyingObject|Bird)[0.8, 0.9], (FlyingObject|Penguin)[0, 0.1],$
$(\exists livesIn.Antarctida|Bird)[0.1, 0.2],$
$(\exists livesIn.Antarctida|Penguin)[0.7, 0.8],$
$P_{tweety} = \{(Penguin|\top)[0.9, 0.9]\}$

Next we explain the meaning of such statements in more detail.

---

[1] No nominals can be used for the time being

**Objective Probabilities** PTBox axioms express *generic* uncertainty and can be seen as probabilistic generalizations of normal subclass axioms in standard DLs. Similarly to TBox axioms, generic constraints (or PTBox axioms) express background knowledge, i.e. general relationships that are known for a particular domain. The informal meaning of a constraint $(D|C)[l, u]$ is [GL02][Luk08] *Given that a fresh, randomly chosen object is an instance of class $C$ its probability of being an instance of class $D$ is within $[l, u]$.*

The words "randomly chosen" are the key here. They mean that the statement does not directly apply any of the named individuals in the knowledge base. Neither does it have to agree with the probabilistic knowledge about such individuals. For example, consider the constraint: $(FlyingObject \sqcap Penguin | \top)[0, 0]$ which states that the absolute probability that a randomly chosen object is a flying penguin is zero. Somewhat surprisingly it does not contradict statements about some *specific* flying penguins, for example, an ABox axiom ($sam$ : $FlyingObject \sqcap Penguin$) or individual conditional constraints (see the next subsection) $(FlyingObject \sqcap Penguin | \top)[1, 1]$ for Sam. The key is that *Sam is not* a random object. Even though one flying penguin is known, it can still be the case that the probability of a *fresh* object to be a flying penguin is zero.

This issue is both, a tricky point and a powerful feature of P-$\mathcal{SHIQ}$(D) that allows for the treatment of exceptions (see Section 3.2). Knowledge expressed by the PTBox constraints is *default* in the sense that it is expected to hold in general, but at the same time might fail in specific cases, e.g., for *Sam*. Note, however, that even though it is admissible to believe that *Sam* is an exception, it would be logically inconsistent to believe that *all* or some fraction of objects are exceptions, i.e., believe that there is a class of flying penguins.

Finally, note that PTBox constraints are aimed to represent *objective* probabilities. Their values can come from all sorts of experiments, for example, clinical trials that yield statistics about relationships in some medical domain. Such probabilities may also be produced by various learning algorithms. Being objective they do not depend on one's belief about individual objects but may (and typically do) affect such beliefs.

**Subjective Probabilities** Continuing the analogy with classical DLs (and OWL) we now proceed to PABox constraints which are probabilistic counterparts of ABox class assertions. PABox constraints are expressed in a special form: $a : (D|\top)[l, u]$. Each of them pertains to a single individual $a$ (such individual is called *probabilistic*). Their informal meaning is as follows: *The degree of belief that the individual $a$ is an instance of the class $D$ is within $[l, u]$.*

The key difference between PTBox and PABox constraints is the kind of probabilities they express. PABox constraints represent degrees of belief rather than statistics, i.e., the interpretation is subjective in this case. In particular, experts such as physicians can create such statements basing on their certainty about symptoms and diagnoses (see the next section for more details on how PABox constraints are used to represent individual risk factors).

PABox constraints characterize only a single individual. They do not have any influence on entailments of PTBox constraints. Also each PABox (a collection of individual constraints) is isolated from other PABoxes. In other words probabilistic knowledge about *specific* individuals is always irrelevant to other individuals. This may be seen as a shortcoming of the formalism because it prohibits probabilistic property assertions between two probabilistic individuals.[2] At the same time it simplifies the formalism and substantially improves scalability with respect to the number of probabilistic individuals.

Another important difference is that PABox constraints are *strict* as opposed to default. It means that there cannot be exceptions for PABox knowledge. Also only PABox constraints can override default PTBox constraints but not vice versa (see Section 3.2).

**Reasoning Services** P-$\mathcal{SHIQ}$(D) offers the following reasoning services: consistency checking and both PTBox and PABox entailment checking. The entailment services are directly exposed to users whereas consistency is an internal procedure carried out by a reasoner to determine whether the ontology can be reasoned with.

The formal description of P-$\mathcal{SHIQ}$(D) entailment relation including its semantic properties is beyond the scope of this paper. This relation specifies how both generic and individual constrains can be inferred from an ontology. The services are query oriented, i.e., reasoners accept queries of the form: $(D|C)[?, ?]$ [GL02] [KP08] and compute the tightest (i.e. the most informative) probability interval. PABox entailments are done by combining PTBox and PABox constraints whereas PTBox constraints are entailed from the PTBox only.

In the next section we will show how the reasoning in P-$\mathcal{SHIQ}$(D) can be used for the probabilistic assessment of breast cancer risk.

## 3 Probabilistic Modeling and the BRCA Ontology

When any new knowledge representation formalism is introduced it is critically important to develop the first realistic model because it greatly helps to reveal modeling principles, discover both good and bad practices and understand the complexity of modeling. We have chosen the domain of breast cancer and the corresponding problem of risk assessment as a target for such prototype.

### 3.1 The BRCA Problem

The problem can be informally stated as follows:

> *Given the statistical knowledge about breast cancer, e.g., causes, associations, etc., and facts about a specific woman, e.g., her age, medical*

---

[2] Probabilistic relationships between probabilistic and classical individuals are allowed but require support of nominals [GL02]

*history, etc., approximately estimate the chances that she will develop breast cancer within some future period of time* [Kom07].

We chose this problem for several reasons: First, the domain is uncertain because not all the risk factors are known and some relationships between them are still to be investigated. Second, the medical domain has a long and successful history of using ontologies, in particular, OWL, so the ability to augment them with uncertain background knowledge is interesting. Third, the problem has been well studied, and there are clear statistical relationships to be formalized (see [Kom07]) as well as developed probabilistic models, e.g., the Gail model [GBB+89]. Finally, there are known entailments of interest to applications (e.g., individual patient risk calculation). All of this makes it an appropriate ground for introducing and evaluating P-$\mathcal{SHIQ}$(D).

Current approaches to the BRCA problem are mostly based on statistical models. They use statistics about risk factors to assess the risk of specific women developing breast cancer. There are some online tools built on top of the Gail model that are available for personal use[3]. The scenario is simple: a woman supplies her relevant risk factors and gets back her approximate risk.

This approach has strong advantages but is not ideal. The underlying statistical models are not transparent or easy to extend. They support only a limited number of risk factors and have been criticized for underestimating risk for certain groups of women, e.g., African Americans. The reasoning results are not machine explainable. Many of such issues are addressable through the use of ontologies and formal reasoning.

## 3.2  P-$\mathcal{SHIQ}$(D) Formalization of the BRCA Problem

The purpose of building a P-$\mathcal{SHIQ}$(D) ontology to serve as probabilistic model of the BRCA problem is to represent the assessment of risk as a reasoning problem in P-$\mathcal{SHIQ}$(D). The model needs to contain knowledge of few kinds: knowledge about important concepts such as risk factors, womens, types of risk; knowledge about statistical relationships between risk factors including their combination; and finally, knowledge about individual women, e.g, their age or ethnicity. Clearly there is a correspondence between such needs and the structure of P-$\mathcal{SHIQ}$(D) knowledge bases: the classical part of the ontology models conceptual knowledge, the PTBox models our knowledge of statistical relationships, and PABoxes model our beliefs about particular women.

It may not be immediately obvious how to start creating a P-$\mathcal{SHIQ}$(D) ontology. For example, it is unclear what OWL vocabulary is needed to support the probabilistic part and how the classes should be organized in a taxonomy. Thus we will begin with the representation of the facts about particular women and the judgment of their risks, and add the remaining knowledge along the way.

---

[3] http://http://www.cancer.gov/bcrisktool/

**Representing Individuals** The BRCA ontology requires some representation of the "input", i.e., a way to describe women and their relevant risk factors. For example, we might say: "Ann is a woman of 57 years old, who is post menopausal and regularly taking hormones including estrogen". Such knowledge can be put into a P-$\mathcal{SHIQ}$(D) ontology by using PABox constraints for Ann. The constraints will look like: $ann : (WomenWithRiskFactorX|\top)[l, u]$; they mean that Ann belongs to a certain category of objects denoted as $WomenWithRiskFactorX$, in particular, women that have risk factor $X$. Such classes are the first ones to be provided by the OWL part of any P-$\mathcal{SHIQ}$(D) ontology. In our case the following will be used: $WomenBetween50And60$, $PostmenopausalWomen$, $WomenWithHighLevelOfEstrogen$.

We will call such classes *primary evidence* classes. They are used to capture the evidence, i.e., what is known about a woman. Note that some discretization might be necessary in the case risk factors represented by continuously distributed numerical attributes like age. So the OWL part will supply classes like $WomenUnder20$, $WomenBetween30And40$, ..., $WomenOver70$ to represent such risk factors.

Also observe that P-$\mathcal{SHIQ}$(D) allows for the representation of uncertain evidence. For example, it could be hard to state with 100% confidence that Ann's level of estrogen is high enough to be a risk factor. Even if she is not taking estrogen pills it can be contained in low quantities in her typical food or drinks, e.g., soya, beer, etc. In that case it might be a better option to specify a degree of belief that she belongs to the class $WomenWithHighLevelOfEstrogen$.

All the primary evidence classes of the BRCA ontology form an OWL taxonomy with the root class $WomenWithRiskFactors$. The taxonomy can be easily extended as soon as the statistics about some new risk factors becomes available.


**Representing Results** The next step will be to define the representation of results produced by the model. It is necessary to ensure that the entailed conditional constraints corresponds to risk assessment statements. Such statements can be of the following kinds [Kom07]:

– Absolute risk assessments, i.e., estimations made without reference to other categories of women. For example, *"an average woman has up to 12.3% chance of developing breast cancer in her lifetime"* [Kom07].
– Relative risk assessments, i.e., estimations describing the increase in risk relatively to the women who do not have some particular risk factors. Statements like *"having BRCA1 gene mutation increases the risk of developing breast cancer by a factor of four"* is an example of relative risk [Kom07].

PABox constraints entailed from the ontology should capture the meaning of one of the above kinds of statements. For example, given some individual, say $Ann$, we want to be able to imply that "Ann's chance of developing breast cancer within some future time is between $l$ and $u$" (note, that this is exactly the kind of output produced by the NCI risk calculator). Recall that PABox constraints represent the probability of membership. So if we define OWL classes

to represent all women that *will* have developed BRC within some future period then PABox constraints having such classes as conclusions will represent the desired probabilities. We distinguish between long-term (lifetime) and short-term (10 years) risk, thus provide two classes: $WomenUnderLifetimeBRCRisk$ and $WomenShorttermBRCRisk$.

The situation with relative risk is less obvious. We need to define categories of women that will be under the risk increased by a certain factor. Factors continuously range from 1 (normal risk) to over 10 (very strong increase in risk, e.g., in the case of BRCA1(2) gene mutations). We handle it similarly to continuously distributed values of risk factors, such as age, i.e., by splitting the full spectrum of values on a finite set of categories (OWL classes). Table 1 lists the OWL classes used to represent categories of women at a certain relative risk:

**Table 1.** Conclusion classes that represent categories of women with respect to their relative risk of developing breast cancer. Each class is associated with an interval which represents the factor of increase in risk.

| Conclusion class | Factor of risk increase |
|---|---|
| $WomanUnderWeaklyIncreasedBRCRisk$ | $[1.2 - 1.5]$ |
| $WomanUnderModeratelyIncreasedBRCRisk$ | $[1.5 - 3.0]$ |
| $WomanUnderStronglyIncreasedBRCRisk$ | $[3.0 - +\infty)$ |

The categories are made disjoint in the BRCA ontology. This enables potentially useful inferences, for example, if a woman has 90% chance of being in the top risk category then it will be inferred that her chances of being under moderately increased risk are up to 10%. Another option might be to organize the categories in a hierarchy. In that case the semantics of, for example, $WomanUnderModeratelyIncreasedBRCRisk$ would change to: class of women whose risk of breast cancer is *at least* moderately increased.

We will call classes used to represent categories of women with respect to risk *ultimate conclusions*. Similarly to the evidence classes they form a taxonomy in the OWL part of the ontology. The root class is $WomenUnderRisk$ which is a parent of classes $WomenUnderAbsoluteRisk$ and $WomenUnderRelativeRisk$.

It is important to notice that such separation on the evidence and conclusion classes is not a domain-specific thing. It seems to be a rather generic principle of P-$\mathcal{SHIQ}$(D) modeling. It can be applied to a range of problems that can be reduced to an uncertain classification of objects (e.g., women in the BRCA case). Then evidence classes can be used to represent known facts about the objects and conclusions — the classification categories. Thus it may be useful to start model design by figuring out possible evidence and conclusion classes similarly to how it is done for the BRCA ontology.

**Representing BRCA Statistics** So far we have described the representation of the input and the output. The missing part is the probabilistic model itself, i.e., the statistical knowledge about the relationships between risk factors or their

combinations and absolute or relative risk of breast cancer. Such relationships are usually inferred by means of experiments such as clinical trials, learning or mining medical data.

At the moment there are dozens of such relationships known. The simplest of them specify how a single risk factor affects the overall risk. Such associations can be straightforwardly represented by PTBox constraints of the form: $(WomenUnderRiskY|WomenWithRiskX)[l, u]$ which can be interpreted as: "the probability that a woman having risk factor $X$ is under risk category $Y$ is between $[l, u]$". Here $X$ can be one of the primary evidence classes and Y one of the ultimate conclusion classes. For example [Kom07]:

$(WomanUnderLifetimeBRCRisk|WomanWithBRCA1Mutation)[0.6, 0.8]$

The BRCA ontology currently contains PTBox statements covering over 20 different risk factors.[4] However such collection of constraints is not yet a model because all the factors are treated separately. Two important things are missing: co-occurrence of risk factors and their combined influence on the risk.

Representing co-occurrence is important for dealing with risk factors which a woman may not be aware of. For example, a woman cannot always be expected to know such her factors as bone density, level of estrogen in blood, etc. However, by capturing statistics about co-occurrence of factors it might be possible to guess on the presence of some factors given probabilistic facts about others. One such example is the statistics that Ashkenazi Jews are more likely to develop BRC1(2) gene mutations which is a critically important factor [Kom07]. To represent such associations we add PTBox constrains that link different evidence classes, e.g.: $(WomenWithBRCAMutation|AshkenaziJewishWoman)[0.025, 0.025]$.

Even more importantly, numerous experiments have revealed that a combination of certain risk factors can be a stronger risk factor than each factor by itself. In other words, risk factors may strengthen or weaken each other's influence on the overall risk. For example, it is known that the harmful effect of estrogen can be substantially worsened by another hormone — progestin [Kom07]. To capture this we need to treat women that have both kinds of evidence (estrogen and progestin) differently from those who only have one.

A straightforward way to capture this would be to add another evidence class to the OWL part which will represent the women that are in both categories — those who have high level of estrogen in blood and those who have high level of progestin: $PWEP \equiv PWE \sqcap PWP$[5] (where classes $PWE, PWP, PWEP$ respectively define post menopausal women who are exposed to high levels of estrogen, progestin or both). Then it can be used in the constraint that represents the boosted risk:

$(WomenUnderModeratelyIncreasedBRCRisk|PWEP)[0.35, 0.35]$.

However it would be problematic for a monotonic formalism because the constraint above contradicts the statistics about estrogen alone:

---

[4] Data taken from the BRC risk factors summary table at:
http://cms.komen.org/komen/AboutBreastCancer/RiskFactorsPrevention

[5] P-$\mathcal{SHIQ}$(D) allows usage of arbitrary class expressions in constraints but for the moment we need to define class names because of implementation limitations.

$(WomenUnderModeratelyIncreasedBRCRisk|PWE)[0.25, 0.25]$.

Women having high levels estrogen and progestin are a *subclass* of those exposed to estrogen only, thus one may expect that any property that holds about the superclass should also hold for the subclass. Here it is not the case. The probability intervals for the two constraints are incompatible, i.e., do not intersect.

One of the most attractive features of P-$\mathcal{SHIQ}$(D) is that this situation does not lead to inconsistency. The entailment relation in P-$\mathcal{SHIQ}$(D) is defined to meaningfully resolve such conflicts in an expected way, i.e., by *overriding* more generic knowledge by more specific one. This is similar to object-oriented programming or reference class reasoning [Luk02]. So if we add the constraint about both hormones then any woman that has both pieces of evidence will be correctly characterized by a higher risk.

Such overriding can, and should, be heavily used in P-$\mathcal{SHIQ}$(D) models. It gives model designers a lot more freedom in representing knowledge without being too much concerned about conflicts that can be brought about by exceptional subclasses. This is not possible in monotonic reasoning systems (recall the famous "birds and penguins" example). By using overriding one can add to the evidence taxonomy as many classes representing different combinations of risk factors as needed.

Finally, overriding also works for individuals, not only for subclasses. PABox constraints about individuals will always override conflicting PTBox constraints (in other words, knowledge about a concrete objects is always more specific than knowledge about any class of objects). For example, in the BRCA ontology one can assert that some woman has low estrogen level even if this contradicts the statistics about her class (e.g., she is taking post menopausal hormones).

Summing up, the BRCA ontology consists of the following major parts:

– A normal OWL ontology that is made up of two main sub-taxonomies — evidence classes (children of $WomenWithRiskFactors$) and conclusion classes (children of $WomenUnderRisk$).
– A set of PTBox constraints that encode several kinds of relationships: associations between single risk factors and the risk, between combinations of risk factors and the risk, and associations between different risk factors.
– A set of individuals (women); each contains a number of PABox constraints.

## 4   Summary and Conclusions

The BRCA ontology is interesting from several points of view. First it offers an alternate view on probabilistic models for risk assessment. But, perhaps even more importantly, it illustrates some principles of probabilistic ontological modeling. We expect the following principles to be generalizable to other use cases:

– **Using probabilistic entailment as the main tool of computation**. A lot of computational problems can be formulated in terms of probabilistic queries, e.g. risk assessment, diagnosing under uncertainty, decision making,

etc. They are all interpretable as special cases of classification under uncertainty which would be similar to classifying women with respect to relative risk categories.

 – **Separation between evidence and conclusion class hierarchies.** It is useful to clearly define classes that will serve for representing facts about individuals (evidence classes) and those that will be used in probabilistic queries (conclusion classes). Such organization helps to reduce the original problem to probabilistic entailment. Once this is done conditional constraints can be naturally used to *connect* classes from the two hierarchies by representing statistical domain knowledge.

 – **Extensive use of overriding.** Exceptional subclasses and individuals exist in any large domain due to incompleteness of knowledge. It is important that they do not prevent modelers from capturing typical relationships in the domain which would lead to missing important entailments.

Without making too general claims we anticipate that at least other problems related to risk assessment can be effectively approached using P-$\mathcal{SHIQ}$(D) and Pronto. They include risks of developing other diseases, risks of privacy breaches, investment risks, etc. All such problems have similar pre-conditions that make probabilistic ontologies an attractive option for modeling.

It is certainly not expected that P-$\mathcal{SHIQ}$(D) will quickly become a standard for modeling under uncertainty similar to OWL for classical modeling. The goal was rather to attract interest from people who are concerned about dealing with uncertain knowledge. The message addressed to such people was that P-$\mathcal{SHIQ}$(D) might help to retain all the benefits of ontological modeling but extend it to cover uncertain domains.

# References

[GBB⁺89]  M H Gail, L A Brinton, D P Byar, D K Corle, S B Green, C Shairer, and J J Mulvihill. Projecting individualized probabilities of developing breast cancer for white females who are being examined annually. *Journal of the National Cancer Institute*, 81(25):1879–1886, 1989.

[GL02]  R Giugno and T Lukasiewicz. $P-\mathcal{SHOQ}(D)$: A probabilistic extension of $\mathcal{SHOQ}$(D) for probabilistic ontologies in the semantic web. Technical Report Nr. 1843-02-06, Institut fur Informationssysteme, Technische Universitat Wien, 2002.

[Kom07]  S G Komen. Breast cancer risk factors table, 2007. Retrieved from: http://cms.komen.org/Komen/AboutBreastCancer/.

[KP08]  P Klinov and B Parsia. Optimization and evaluation of reasoning in probabilistic description logic: Towards a systematic approach. Accepted to the International Semantic Web Conference, 2008.

[Luk02]  T Lukasiewicz. Probabilistic default reasoning with conditional constraints. *Annals of Mathematics and Artificial Intelligence*, 34(1-3):35–88, 2002.

[Luk08]  T Lukasiewicz. Expressive probabilistic description logics. *Artificial Intelligence*, 172(6-7):852–883, 2008.