

# Optimisation des approximations de probabilité des requêtes en XML probabiliste

Asma Souihli

Institut Mines–Télécom

Télécom ParisTech; CNRS LTCI  
asma.souihli@telecom-paristech.fr

Pierre Senellart

Institut Mines–Télécom

Télécom ParisTech; CNRS LTCI  
pierre.senellart@telecom-paristech.fr

XML probabiliste est un modèle probabiliste pour les bases de données incertaines semi-structurées, avec des applications telles que l’intégration incertaine de données, l’extraction d’informations ou le contrôle probabiliste de versions. Nous explorons dans ce travail une solution efficace pour l’évaluation des requêtes *tree-pattern* avec jointures sur ces documents, ou, plus précisément, pour l’approximation de la probabilité d’une requête booléenne sur un *document probabiliste*. L’approche repose sur, d’une part, la production de la provenance probabiliste de la requête posée, et, d’autre part, la recherche d’une stratégie optimale pour estimer la probabilité de cette provenance. Cette deuxième partie s’inspire des approches des optimiseurs de requêtes : l’exploration de différents plans d’évaluation pour différentes parties de la formule et l’estimation du coût de chaque plan, suivant un modèle de coût établi pour les algorithmes de calcul utilisés. Nous démontrons l’efficacité de cette approche sur des jeux de données utilisés dans des travaux précédents sur l’interrogation des bases de données XML probabilistes, ainsi que sur des données synthétiques.

## 1 Introduction

Uncertainty comes along with data generated by imprecise automatic tasks such as data extraction and integration, data mining, or natural language processing. A measure of this uncertainty may be induced from the trustworthiness of the resources, the quality of the data mapping procedure, etc. Often, information is described in a semistructured manner because representation by means of a hierarchy of nodes is natural, especially when the source (e.g., XML or HTML) is already in this form. One possible way, among the most natural, to represent this uncertainty is through probabilistic databases, and probabilistic XML [23] in particular.

*Probabilistic documents* or *p-documents* [3, 20] are a general representation system for probabilistic XML, based on probabilistic XML trees with ordinary and distributional nodes. The latter define the process of generating a random XML instance following the specified distribution at the level of each node. The model is a compact and complete representation of a probabilistic space of documents (i.e., a finite set of possible worlds, each with a particular probability).

Probabilistic documents have been used in various applications [23], such as uncertain data integration [35], XML warehousing [30], uncertain version control systems [1], or Web information extraction. We do not deal here with obtaining probabilistic documents; we assume the p-document given and investigate how to efficiently query it. To simplify the presentation, we consider Boolean queries. Over a deterministic XML document, a Boolean query returns *true* or *false*. Over a p-document, a Boolean query returns the *probability* that the query is true, i.e., the sum of the probabilities of all possible worlds where the query is true.

In contrast with existing work [20–22] that has proposed algorithms for tractable subcases and characterized the complexity of the problem, we consider in this work a very general form of p-documents (involving arbitrary correlations between nodes of the tree) together with a large class of queries (tree-pattern queries with joins, with results extensible to the even more general class of *locally monotone* [30] queries) and aim at a practical solution for querying p-documents. Computing the exact probability of a Boolean query over a p-document is #P-hard [21], but, under data complexity, there are fully polynomial-time randomized approximation schemes [21]: we thus focus on efficiently approximating the probability of a query. Indeed, for many applications, one simply needs a good estimate of the probability value, i.e., an approximation *to a multiplicative factor* of the correct probability, *with high confidence*. Approximations *to an additive factor* are of lesser interest, since they make it hard to distinguish between, say, probabilities of  $10^{-2}$  and of  $10^{-5}$ .

Following ideas from [21, 30], we reduce the problem of approximating the probability of query  $Q$  over p-document  $\mathcal{P}$  to approximating the probability of a propositional formula  $\varphi$  in disjunctive normal form (DNF). We first rewrite the initial (XPath) query  $Q$  into an XQuery query  $Q'$  that returns, for every match of  $Q$  over the deterministic document underlying  $\mathcal{P}$ , the conjunction of events conditioning this match (labeled probabilities related to the probabilistic nodes of the p-document). This rewritten query is then evaluated by a standard XQuery processor, allowing the use of all standard XML indexing techniques. This step can be done in time polynomial in the size of  $\mathcal{P}$ . The disjunction of all conjunctions of events associated to matches to the query constitutes the propositional *lineage*  $\varphi$  of the query. Note that computing the probability of a propositional formula relates to computing the number of satisfying assignments of the formula, a problem known to be #P-complete [34].

Our system, ProApproX, has the characteristic of not relying on a single algorithm to evaluate the probability of a lineage formula but of deciding on the algorithm to be used based on a *cost model* of the various potential evaluation algorithms. In addition, the formula is compiled into an evaluation plan, different evaluation algorithms can be used on different subparts of the formula, and the overall cost of the whole plan is estimated. As with regular query optimizers, the space of evaluation plans (for a user-specified approximation guarantee) is searched for one of optimal cost. ProApproX thus reveals its originality through the following major features:

1. The support of a broader range of XPath queries over a more general data model than current probabilistic XML systems (Sections 2 and 3);
2. A cost model for a variety of probability evaluation algorithms, based on which the selection of the most efficient algorithm is performed (Section 4);
3. Simplification of the computation process via a decomposition of the probabilistic lineage into independent computational cells (Section 5);
4. The possibility of setting arbitrary error bound  $\varepsilon$  and confidence  $\delta$  for the desired probabilistic approximation, with well-grounded propagation mechanisms of error and reliability between computational units (Section 6);
5. An exploration of the space of evaluation plans based on the proposed cost model (Section 7).

We start with a brief introduction to the probabilistic XML model, and to probabilistic lineage. After presenting all features of the algorithm, we present an extensive experimental evaluation of the proposed system in Section 8, and go over the related work in Section 9 before concluding.

## 2 Probabilistic XML

We recall here some basic notions about probabilistic XML as an uncertain data model. The model presented here was introduced and studied in [3, 20, 21] as a generalization of previously existing models [4, 26, 36].

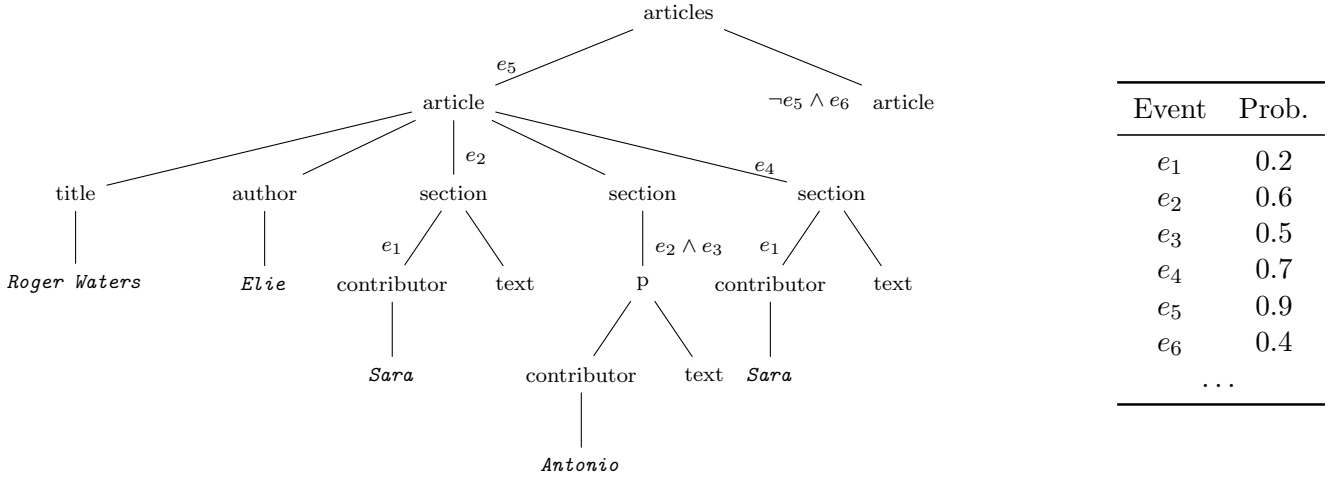


Figure 1: Tree representation of a p-document with long-distance dependencies

### Documents and p-Documents

We model XML documents as unranked, unordered, labeled trees. Not taking into account the order between sibling nodes in an XML document is a common but non-crucial assumption. The same modeling can be done for ordered trees, without much change to the theory.

A *probabilistic XML document* (or *p-document*) is similar to a document, with the difference that it has two types of nodes: ordinary and distributional. Distributional nodes are fictive nodes that specify how their children can be randomly selected. Given the criteria of dependency between elements, we distinguish two types of data representation models. In the *local dependency* model [26, 36] (named  $\text{PrXML}^{mux,ind}$  in [3]), choices made for different nodes are independent or locally dependent, i.e., a distributional node chooses one or more children independently from other choices made at other levels of the tree. In the *long-distance dependency* model [4] ( $\text{PrXML}^{cie}$  in [3]), the generation at a given distributional node might also depend on different conditions (i.e., probabilistic choices) related to other parts of the tree. This model, more general than the local dependency model, is based on one distributional node, *cie* (for *conjunction of independent events*): nodes of this type are associated with a conjunction of independent (possibly negated) random Boolean variables  $e_1 \dots e_m$  called *events*; each event has a global and independent probability  $\text{Pr}(e_i)$  that  $e_i$  is true. Note that different *cie* nodes share common events, i.e., choices can be correlated.

The *semantics* of a p-document is a probability distribution over a set of possible documents, defined by the following process for the long-distance dependency case. First, randomly draw a truth assignment for each of the  $e_i$ 's, following  $\text{Pr}(e_i)$ . Then remove all distributional nodes whose condition is falsified by this assignment. Descendant of removed nodes are removed, and the remaining ordinary nodes are connected to their lowest ordinary ancestor, yielding a regular document. The probability of this document is that of all truth assignments that generate it.

As noted in [3], local dependencies expressed by *ind* (independent) or *mux* (mutually exclusive) nodes can be tractably translated into *cie* nodes. We use these results to efficiently translate datasets in the local-dependency model into *p-documents* that rely exclusively on *cie* distributional nodes. This allows us to consider the long-distance dependency model only in the remaining of the paper.

As an example, Figure 1 presents a fragment of a probabilistic XML document describing a given Wikipedia article as a merge of all its (uncertain) revisions, reproduced (with elaboration) from [1]. For ease of presentation, *cie* nodes are shown by simply annotating edges with conjunctions of literals. For example, the first and third “section” elements in the tree appear under *cie* nodes that retain them only if  $e_1$  is true; the whole first “article” is kept if and only if  $e_5$  is true; etc. Events  $e_i$ 's appearing in the tree correspond to particular contributors, or probabilistic events that particular revisions are correct given that the authors are reliable [1], with the probability distribution of the  $\text{Pr}(e_i)$ 's typically inferred by a trust inference algorithm.

**Querying Probabilistic XML** The query language we consider in this work is *tree-pattern queries with joins* [4]. A tree-pattern query with joins is given by a tree pattern (a tree whose edges are labeled with either *child* or *descendant* and whose nodes are labeled, possibly by a wildcard) together with a set of *value joins* that impose that two nodes of the tree pattern are matched to nodes with the same value. A *match* of a query to a (deterministic) document must map all nodes of the query to the document, respecting the following constraints: (i) the root is mapped to the root; (ii) child edges are mapped to edges of the tree; (iii) descendant edges are mapped to a sequence of child edges; (iv) non-wildcard labels are preserved; (v) nodes in a join condition must have the same label. As in [4], we view each match as the minimal subtree containing all nodes of a document that are mapped to a node of the query. Queries are given using the standard XPath syntax (note that some tree-pattern queries with joins cannot be expressed in XPath 1.0 but they can all be expressed in XPath 2.0).

**Example 1.** Given the p-document of Figure 1, we can search for all contributors to a given article using the tree-pattern query  $Q_1$ :

```
//article[title='Roger Waters']//contributor
```

Similarly, the tree-pattern query with join  $Q_2$ :

```
/articles/article[author=../contributor]
```

looks for articles where the author appears among the contributors.

Though we focus on tree-pattern queries with joins for simplicity, a similar processing can be applied to a larger class of queries, namely *locally monotone* [30] queries, that includes, for example, the whole positive fragment of first-order XPath [7].

Again for simplicity, we restrict to the Boolean projection of tree-pattern queries with joins in what follows: extension to queries that return a set of nodes is completely straightforward, see, e.g., [21]. A Boolean query over a deterministic document returns either true or false. A Boolean query  $Q$  over a p-document  $\mathcal{P}$  returns the *probability* that  $Q$  is true in  $\mathcal{P}$ , that is, the sum of the probabilities of all possible documents of  $\mathcal{P}$  where  $Q$  is true.

A fundamental observation [30], due to the locally monotone nature of the query language, is that the probability of a query over a p-document is the probability that one of the match of the query over the underlying deterministic document remains in  $\mathcal{P}$ . We will use this fact to transform our probabilistic XML querying problem into a probability computation problem over a propositional formula in DNF.

### 3 Probabilistic XML Query Lineage

Recall the example query  $Q_1$  that looks for the list of contributors to the revisions of the article entitled “Roger Waters”. Looking at Figure 1 we can see three matches to  $Q_1$  on the deterministic document underlying the p-document, corresponding to the three different “contributor” nodes. For each match, we can construct the conjunction of all probabilistic literals on the path from the root to one of the node matched by the query:  $c_1 = e_5 \wedge e_1 \wedge e_2$ ,  $c_2 = e_5 \wedge e_2 \wedge e_3$ ,  $c_3 = e_5 \wedge e_1 \wedge e_4$ . These conjunctions are called the *probabilistic lineage* of each match.

As noted, the probability of a locally monotone query  $Q$  over a p-document is exactly the probability of the disjunction of probabilistic lineages of all matches to the query on the underlying deterministic document. The probability of  $Q_1$  over the p-document of Figure 1 is thus the probability of the following propositional formula, called the *probabilistic lineage*  $\varphi_1$  of the Boolean query  $Q_1$  over this document:

$$\varphi_1 = c_1 \vee c_2 \vee c_3 = (e_5 \wedge e_1 \wedge e_2) \vee (e_5 \wedge e_2 \wedge e_3) \vee (e_5 \wedge e_1 \wedge e_4).$$

Because we only have conjunctions of literals on distributional nodes in PrXML<sup>cie</sup>, probabilistic lineages are always in disjunctive normal form (DNF). In the following, we will often note this DNF  $\varphi = \bigvee_{i=1}^m c_i$  and will call each  $c_i$  a *clause* of the DNF.

In ProApproX, probabilistic events for a given node are stored in an attribute node named “prob” of a regular XML document, stored and managed by an ordinary native XML DBMS. To retrieve the lineage of a query, we need to transform our XPath query  $Q$  into an XQuery query  $Q'$  that returns the concatenation of all “prob” attributes of each match. In our example,  $Q_1$  is transformed into query  $Q'_1$ :

```

for $a in //article
for $b in $a/title/text()[.= 'Roger Waters']
for $c in $a//contributor
let $leaves:=( $b, $c)
let $atts:=(for $i in $leaves
  return $i/ancestor-or-self::*/@prob)
return <match> {
  distinct-values(
    for $att in $atts
    return string($att))} </match>

```

A generic translator from tree-pattern queries with joins encoded in XPath into XQuery lineage queries is implemented as part of ProApproX.

We have then separated the problem of querying probabilistic XML into two independent problems: (i) evaluating an XQuery query over a regular XML document; (ii) computing the probability of a formula in DNF. The former problem can be solved using any XQuery processor; we use the native XML DBMS BaseX<sup>1</sup> after comparing its performance on our datasets to other XQuery engines. The latter problem is the focus on the remaining of this paper.

Computing the exact probability of a formula in DNF is a #P-hard problem in general, even when all events have the same probability [34], and it is easy to see that any DNF formula can be generated as the lineage of even a trivial XPath query such as //A [20]. However, we are going to detect some easy cases and exploit polynomial-time randomized approximation algorithms [21] to build an efficient query processor over probabilistic XML.

## 4 Algorithms and Cost Models

In this section we present a collection of algorithms for computing or approximating the probability of a DNF formula  $\varphi$ . Based on the computational complexity of each of these algorithms, we elaborate a cost model that estimates the runtime (in ms) of each algorithm  $\text{alg}$  as a function  $\text{cost}_{\text{alg}}$  of different features of  $\varphi$ , of approximation parameters defined further, and of a cost constant  $C_{\text{alg}}$  representing the time needed for the elementary, inner, parts of each algorithm. A value for  $C_{\text{alg}}$  is measured experimentally by varying the size of synthetic formulas (see Section 8 for precisions on generation) and recording the processing time of the actual implementation of the algorithm on a given machine. We only consider sequential processing here, see the conclusion for discussions about parallelization. We give in Table 1 a summary of the notation used and in Table 2 a summary of the cost models of all algorithms.

### 4.1 Exact Computation

We start with exact computation algorithms. Since the problem is #P-hard, only exponential-time algorithms are known, but they are typically very fast on DNFs of small size. We present two such algorithms, that are folklore, and are well-suited, respectively, to the case when there are few variables, or few clauses, in  $\varphi$ .

---

<sup>1</sup><http://basex.org/>

Table 1: Notation

$\varphi$	DNF formula $\varphi = \bigvee_{i=1}^m c_i$
$m$	number of clauses in $\varphi$
$N$	number of event variables used in $\varphi$
$c_i$	a clause of $\varphi$
$k_i$	size of clause $c_i$
$L$	total size of $\varphi$ , $L = \sum_{i=1}^m k_i$
$t$	number of trials (samples)
$\varepsilon$	approximation error
$1 - \delta$	probabilistic approximation guarantee
$\ell$	lower bound on $\varphi$ 's probability
$\text{cost}_{\text{alg}}$	cost of an algorithm alg
$C_{\text{alg}}$	cost constant for algorithm alg

Table 2: Cost models and cost constants

Algorithm alg	$\text{cost}_{\text{alg}}$	$C_{\text{alg}}$ (ms)
naïve	$C_{\text{naïve}} \times 2^N \times L$	$4 \cdot 10^{-5}$
sieve	$C_{\text{sieve}} \times 2^m \times \frac{L}{m}$	$5 \cdot 10^{-5}$
AddMC	$C_{\text{AddMC}} \times \ln \frac{2}{\delta} \times \frac{L}{\varepsilon^2}$	$4 \cdot 10^{-5}$
MulMC	$C_{\text{AddMC}} \times \ln \frac{2}{\delta} \times \frac{L}{\ell^2 \varepsilon^2}$	$4 \cdot 10^{-5}$
coverage	$C_{\text{coverage}} \times \ln \frac{2}{\delta} \times \frac{(1+\varepsilon) \times L}{\varepsilon^2}$	$10^{-3}$

**Possible worlds (Naïve algorithm)** This algorithm corresponds to the naïve, exponential-time, iteration over all possible  $2^N$  truth value assignments to the  $N$  variables used in  $\varphi$ , i.e., over all possible worlds. It simply consists in summing up the probabilities of the satisfying assignments. In the example of the previous sections, the formula  $\varphi_1$  has 5 different variables, and we can enumerate all 32 possible worlds; 8 of them make  $\varphi_1$  true, and one can check that their total probability is 0.3744.

The computational complexity of the naïve algorithm is obviously  $O(2^N \times L)$ : we enumerate all possible worlds, and for each one we evaluate the truth value of the formula in linear time. We therefore set the following cost:

$$\text{cost}_{\text{naïve}} = C_{\text{naïve}} \times 2^N \times L$$

where  $C_{\text{naïve}}$  is determined experimentally as previously explained and is equal to  $4 \cdot 10^{-5}$  ms (see Table 2 for all cost constants). This means, for instance, that the expected runtime of the naïve algorithm for a formula with 10 variables and of length 500 is  $C_{\text{naïve}} \times 2^{10} \times 500 \approx 20$  ms (for the machine for which the cost constants were computed), which is reasonably low and shows that, at least for some DNF formulas, the naïve algorithm can be adapted.

**Inclusion–exclusion (sieve)** It is also possible to apply the *inclusion–exclusion* or *sieve* principle to compute the probability of  $\varphi$ . The *sieve* decomposition of the DNF  $\varphi$  is:

$$\Pr\left(\bigvee_{i=1}^n c_i\right) = \sum_{k=1}^n (-1)^{k-1} \sum_{\substack{I \subseteq \{1, \dots, n\} \\ |I|=k}} \Pr(c_I), \text{ where } c_I := \bigwedge_{i \in I} c_i.$$

On our example,  $\Pr(\varphi_1) = \Pr(c_1) + \Pr(c_2) + \Pr(c_3) - \Pr(c_1 \wedge c_2) - \Pr(c_2 \wedge c_3) - \Pr(c_1 \wedge c_3) + \Pr(c_1 \wedge c_2 \wedge c_3) = 0.108 + 0.27 + 0.126 - 0.054 - 0.0378 - 0.0756 + 0.0378 = 0.3744$ . The computational complexity is  $O(2^m \times \frac{L}{m})$ : the probability of the conjunction of each set of clauses, whose typical size is of the order

of  $\frac{L}{m}$ , can be computed in linear time. We set:

$$\text{cost} = C_{\text{sieve}} \times 2^m \times N,$$

which becomes competitive with respect to the naïve algorithm when clauses are few but long.

Note that, as described here, the sieve method is numerically instable. Implementing the sieve formula presented earlier in floating-point arithmetic results in a low accuracy because of rounding errors in sequences of additions and deletions. The sieve method can be improved towards better numerical stability using an algorithm proposed by Heidtmann [14].

## 4.2 Approximation Algorithms

We present in this section two randomized approximation algorithms, giving respectively *additive* and *multiplicative* guarantees. For fixed  $\varepsilon, \delta$ , we say  $A(\varphi)$  is an additive  $\varepsilon$ -approximation of  $\Pr(\varphi)$  with probabilistic guarantee  $1 - \delta$  if, with probability at least  $1 - \delta$ , the following holds:

$$\Pr(\varphi) - \varepsilon \leq A(\varphi) \leq \Pr(\varphi) + \varepsilon.$$

Similarly, for fixed  $\varepsilon, \delta$ ,  $A(\varphi)$  is a multiplicative  $\varepsilon$ -approximation of  $\Pr(\varphi)$  with probabilistic guarantee  $1 - \delta$  if, with probability at least  $1 - \delta$ :

$$(1 - \varepsilon) \times \Pr(\varphi) \leq A(\varphi) \leq (1 + \varepsilon) \times \Pr(\varphi).$$

In Section 6, we explain how to turn additive guarantees into multiplicative ones.

**Additive Monte Carlo** The simplest way to implement an additive approximation is by sampling, Monte-Carlo-style, the space of all possible assignments following the distribution of the literals, and test whether the picked assignment satisfies at least one clause. If we conduct  $t$  trials, an (unbiased) estimator  $A(\varphi)$  for the probability of  $\varphi$  will be the proportion of the trials that led to a satisfaction of  $\varphi$ . Following Hoeffding's bound [15], we have:

$$\Pr(|A(\varphi) - \Pr(\varphi)| > \varepsilon) \leq 2e^{-2t\varepsilon^2}.$$

Therefore, the number of trials  $t$  needed to perform an  $(\varepsilon, \delta)$ -additive approximation is:

$$t = \left\lceil \frac{\ln 2 - \ln \delta}{2\varepsilon^2} \right\rceil.$$

For each sample, the algorithm linearly scans the formula, which gives an overall complexity of  $O(t \times L)$  and a cost model:

$$\text{cost}_{\text{AddMC}} = C_{\text{AddMC}} \times \ln \frac{2}{\delta} \times \frac{L}{\varepsilon^2}.$$

**Self-Adjusting Coverage Algorithm** This algorithm was introduced by Karp, Luby, and Madras [19] and is a fully-polynomial randomized approximation scheme (FPRAS), which means it produces a multiplicative estimate of the probability. We base our cost model on the following result:

**SELF-ADJUSTING COVERAGE ALGORITHM THEOREM I** [19]. *When  $\varepsilon < 1$  and  $t = (8 \times (1 + \varepsilon) \times m \ln(2/\delta))/\varepsilon^2$ , the self-adjusting coverage algorithm yields an  $(\varepsilon, \delta)$ -approximation.*

Since each trial requires evaluating one clause of  $\varphi$ , we have:

$$\text{cost}_{\text{coverage}} = C_{\text{coverage}} \times \frac{L}{m} \times \frac{(1 + \varepsilon)m}{\varepsilon^2} \times \ln \frac{2}{\delta}.$$

As can be seen on Table 2, the cost constant computed for this self-adjusting coverage algorithm is much higher than the cost constants for other algorithms, reflecting the fact that though it has an excellent algorithmic complexity, this can still be a costly algorithm to use in practice.

As we shall see, each of these four algorithms is most efficient on some of the possible  $\varphi$ 's. However, we do not apply them on the whole formula, but start by decomposing  $\varphi$  into an evaluation tree built out of simpler formulas.

## 5 Building the Evaluation Tree

We explain in this section how ProApproX decomposes a DNF formula  $\varphi$  into an evaluation tree formed of simpler subformulas. We assume in this section that no contradiction appears in  $\varphi$  (a clause that contains both  $x$  and  $\neg x$  and is removed from  $\varphi$ ). To illustrate, we use the following running example:

**Example 2.** Consider the formula  $\varphi_2 = (e_1 \wedge e_2) \vee (e_1 \wedge e_3) \vee (\neg e_4 \wedge e_5) \vee (e_4 \wedge e_6) \vee (e_4 \wedge e_7)$ . To compute  $\Pr(\varphi_2)$ , we rewrite  $\varphi_2$  as the composition of subformulas whose probabilities are simpler to compute.

Our evaluation tree is based on three operations (independence detection, factorization, inconsistency detection) that are repeatedly applied on the original formula. All three operations yield trees whose leaves are formula in DNF and whose internal nodes are Boolean operations (independent disjunction  $\vee$ , independent conjunction  $\wedge$ , and mutually exclusive disjunction  $\oplus$ ) that support efficient probability computations.

**Independence** We say that two clauses are *independent* if they do not share any variables. Our first operation attempts to write the DNF formula  $\varphi$  as an *independent disjunction*  $\varphi' \vee \varphi''$  of two DNF formulas that form a partition of  $\varphi$  such that every clause of  $\varphi'$  is independent from every clause of  $\varphi''$ . In our example, we can write  $\varphi_2$  as  $((e_1 \wedge e_2) \vee (e_1 \wedge e_3)) \vee ((\neg e_4 \wedge e_5) \vee (e_4 \wedge e_6) \vee (e_4 \wedge e_7))$ .

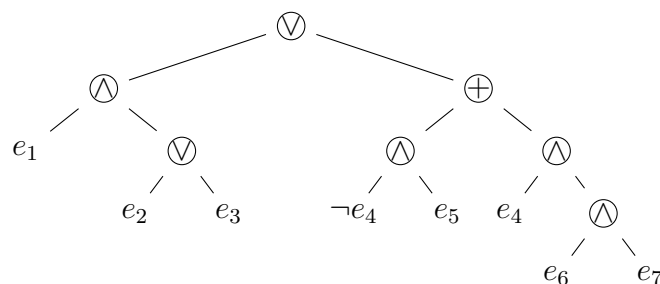
**Factorization** Our second operation consists in *factoring* the intersection  $c$  of all clauses  $c_i$  of formula  $\varphi$  out of the clauses, obtaining a rewriting of  $\varphi$  as  $\varphi' = c \wedge (c'_1 \vee \dots \vee c'_m)$  where  $\wedge$  denotes independent conjunction:  $c$  is independent of every  $c_i$ . Note that both operands of  $\wedge$  are in DNF: the left-hand side is a simple conjunctive clause, and the right-hand side is a DNF formula. If we apply this factorization to the subformulas appearing in our previous rewriting of  $\varphi_2$ , we obtain  $(e_1 \wedge (e_2 \vee e_3)) \vee ((\neg e_4 \wedge e_5) \vee (e_4 \wedge e_6) \vee (e_4 \wedge e_7))$ .

**Inconsistency** Our final operator looks for a partition of a DNF formula into two inconsistent subformulas, i.e., a rewriting of  $\varphi$  into  $\varphi' \oplus \varphi''$  where  $\varphi'$  and  $\varphi''$  are partitions of the clauses of  $\varphi$ , and there is a variable  $x$  such that all clauses of  $\varphi'$  have a literal  $\neg x$  and all clauses of  $\varphi''$  have a literal  $x$ . Our example formula becomes  $(e_1 \wedge (e_2 \vee e_3)) \vee ((\neg e_4 \wedge e_5) \oplus ((e_4 \wedge e_6) \vee (e_4 \wedge e_7)))$ .

**Evaluation tree** ProApproX repeatedly applies these three operations until no further rewritings are possible. We obtain thus for  $\varphi_2$ :

$$(e_1 \wedge (e_2 \vee e_3)) \vee ((\neg e_4 \wedge e_5) \oplus (e_4 \wedge (e_6 \vee e_7))).$$

This defines an evaluation tree for  $\varphi_2$  in a straightforward manner; inner nodes are the Boolean operations  $\vee$ ,  $\wedge$ ,  $\oplus$  (assumed binary for simplicity), and leaves are DNF formulas that cannot be simplified any further (in the simple case of  $\varphi_2$ , these are just trivial formulas formed of a single literal):





Note that the process is nondeterministic, since it depends on the order of application of operations. We choose arbitrarily one such order and do not make any claim at optimality of the simplification obtained.

These evaluation trees can be used to compute the probability of the main formula in terms of the probabilities of the formulas on the leaves, thanks to the following observations:

$$\begin{aligned}\Pr(\psi_1 \odot \psi_2) &= 1 - (1 - \Pr(\psi_1)) \times (1 - \Pr(\psi_2)) \\ \Pr(\psi_1 \otimes \psi_2) &= \Pr(\psi_1) \times \Pr(\psi_2) \\ \Pr(\psi_1 \oplus \psi_2) &= \Pr(\psi_1) + \Pr(\psi_2)\end{aligned}$$

**Towards an evaluation plan** In general, the evaluation tree of a formula still contains as leaves DNF formulas that are hard to compute. To compute the probability of the global formula, we turn these evaluation trees into evaluation plans by assigning to every leaf of the tree one algorithm (either exact or approximate) that will be used to compute the probability of this leaf. This assignment will obviously use the cost model for the different algorithms, but we need to be careful about the following aspects of the problem: (i) We need to propagate approximation parameters down the tree in a principled manner, so that we keep the approximation guarantees on the global probability (see next section). (ii) It is sometimes worth it not going down to the level of leaves of the approximation tree to evaluate the probabilities, but to do it at a higher level; in other words, we might want to assign evaluation algorithms to internal nodes of the tree rather than to leaves (see Section 7).

## 6 Propagation of approximation parameters

Our overall objective is to obtain a multiplicative approximation of the probability of a formula  $\varphi$  in DNF, given an approximation tolerance  $\varepsilon$  and probabilistic approximation guarantee  $1 - \delta$ . We have explained in the previous section that we will achieve this by decomposing  $\varphi$  into its evaluation tree and assigning different algorithms, some exact, and some approximate, to different parts of the tree. These assignments must be decided in so that the produced precisions at each node yield an overall approximation that does not exceed the tolerance set for the approximation of the DNF formula as a whole. In this section, we establish the different conditions for correct error allocation for children of  $\odot$ ,  $\otimes$ , or  $\oplus$  nodes, in the case when the probability of both nodes is approximated. We also explain how the approximation guarantee  $1 - \delta$  is propagated, and how to turn an additive approximation into a multiplicative one.

### Propagation of approximation error $\varepsilon$ .

**Proposition 1.** *Let  $\varphi = \psi_1 \odot \psi_2$ , and assume  $\tilde{p}_1$  and  $\tilde{p}_2$  are multiplicative approximations of  $\Pr(\psi_1)$  and  $\Pr(\psi_2)$ , to a factor of  $\varepsilon_1$  and  $\varepsilon_2$ , respectively. Then  $1 - (1 - \tilde{p}_1)(1 - \tilde{p}_2)$  is a multiplicative approximation of  $\Pr(\varphi)$  to a factor of  $\varepsilon$  if  $\varepsilon = \varepsilon_1 + \varepsilon_2$ .*

*Proof.* Let us denote  $p$  the probability of  $\varphi$  and similarly  $p_1$  and  $p_2$  the probabilities of  $\psi_1$  and  $\psi_2$ ;  $\tilde{p}_1$  and  $\tilde{p}_2$  are multiplicative approximations of  $p_1$  and  $p_2$ , and

$$\begin{cases} (1 - \varepsilon_1)p_1 \leq \tilde{p}_1 \leq (1 + \varepsilon_1)p_1 \\ (1 - \varepsilon_2)p_2 \leq \tilde{p}_2 \leq (1 + \varepsilon_2)p_2 \end{cases}$$

Then:

$$\begin{cases} 1 - (1 + \varepsilon_1)p_1 \leq 1 - \tilde{p}_1 \leq 1 - (1 - \varepsilon_1)p_1 \\ 1 - (1 + \varepsilon_2)p_2 \leq 1 - \tilde{p}_2 \leq 1 - (1 - \varepsilon_2)p_2 \end{cases}$$

And thus:

$$\begin{aligned}1 - [(1 - (1 - \varepsilon_1)p_1)(1 - (1 - \varepsilon_2)p_2)] &\leq 1 - (1 - \tilde{p}_1)(1 - \tilde{p}_2) \\ 1 - (1 - \tilde{p}_1)(1 - \tilde{p}_2) &\leq 1 - [(1 - (1 + \varepsilon_1)p_1)(1 - (1 + \varepsilon_2)p_2)]\end{aligned}$$

Let  $\alpha$  be  $1 - [(1 - (1 - \varepsilon_1)p_1)(1 - (1 - \varepsilon_2)p_2)]$ . We have:

$$\begin{aligned}\alpha &= 1 - [(1 - p_1 + \varepsilon_1 p_1)(1 - p_2 + \varepsilon_2 p_2)] \\ &= 1 - \underbrace{(1 - p_1)(1 - p_2)}_{=p} \\ &\quad - [(1 - p_1)\varepsilon_2 p_2 + (1 - p_2)\varepsilon_1 p_1 + \varepsilon_1 \varepsilon_2 p_1 p_2].\end{aligned}$$

Let  $\lambda = [(1 - p_1)\varepsilon_2 p_2 + (1 - p_2)\varepsilon_1 p_1 + \varepsilon_1 \varepsilon_2 p_1 p_2]$ .  $\frac{\lambda}{p} = \frac{(1-p_1)\varepsilon_2 p_2 + (1-p_2)\varepsilon_1 p_1 + \varepsilon_1 \varepsilon_2 p_1 p_2}{p_1 + p_2 - p_1 p_2}$ . Let  $\Pi = \max(p_1, p_2)$ ; we have  $p \geq \Pi$ . We write  $\lambda$  as  $\varepsilon_1 p_1 [1 + p_2(\varepsilon_2 - 1)] + \varepsilon_2 p_2 (1 - p_1)$ ; clearly,  $\lambda \leq \varepsilon_1 \Pi + \varepsilon_2 \Pi$ . Then  $\frac{\lambda}{p} \leq (\varepsilon_1 + \varepsilon_2) \frac{p}{p} \leq \varepsilon_1 + \varepsilon_2 = \varepsilon$  and  $\lambda \leq \varepsilon p$ .

We proceed similarly for the upper bound.  $\square$

In practice, when we propagate an  $\varepsilon$  factor down the tree, we look for different combinations of  $\varepsilon_1$  and  $\varepsilon_2$  that sum to 1, evaluate the cost of each combination, and choose the best found (this is one example of the general exploration of evaluation plans described in the following section). Note that if  $\varepsilon_2 = 0$  (that is, if we compute the exact value of  $\Pr(\psi_1)$ ), we can set  $\varepsilon_2 = \varepsilon$ .

**Proposition 2.** *Let  $\varphi = \psi_1 \oplus \psi_2$ , and assume  $\tilde{p}_1$  and  $\tilde{p}_2$  are multiplicative approximations of  $\Pr(\psi_1)$  and  $\Pr(\psi_2)$ , to a factor of  $\varepsilon_1$  and  $\varepsilon_2$ , respectively. Then  $\tilde{p}_1 + \tilde{p}_2$  is a multiplicative approximation of  $\Pr(\varphi)$  to a factor of  $\varepsilon$  if  $\varepsilon = \max(\varepsilon_1, \varepsilon_2)$ .*

*Proof.* We use the same notation as before.

$$\begin{cases} p_1 - \varepsilon_1 p_1 \leq \tilde{p}_1 \leq p_1 + \varepsilon_1 p_1 \\ p_2 - \varepsilon_2 p_2 \leq \tilde{p}_2 \leq p_2 + \varepsilon_2 p_2 \end{cases}$$

and then:

$$p_1 + p_2 - \varepsilon_1 p_1 - \varepsilon_2 p_2 \leq \tilde{p}_1 + \tilde{p}_2 \leq p_1 + p_2 + (\varepsilon_1 p_1 + \varepsilon_2 p_2)$$

Note that  $(\varepsilon_1 p_1 + \varepsilon_2 p_2) \leq \varepsilon(p_1 + p_2) = \varepsilon p$  which gives the required lower bound; the upper bound is similar. In particular, if  $\varepsilon_2 = 0$  we can set again  $\varepsilon_2 = \varepsilon$ .  $\square$

For the  $\hat{\wedge}$  operator, observe that it never appears between two arbitrary DNF formulas, but between a formula in DNF and a conjunction, the probability of the latter can be evaluated exactly in a tractable manner. We thus just state the propagation condition in this case:

**Proposition 3.** *Let  $\varphi = \psi_1 \hat{\wedge} \psi_2$ , and assume  $\tilde{p}_1$  is a multiplicative approximation of  $\Pr(\psi_1)$  to a factor of  $\varepsilon_1$ . Then  $\tilde{p}_1 \times \Pr(\psi_2)$  is a multiplicative approximation of  $\Pr(\varphi)$  to a factor of  $\varepsilon$  if  $\varepsilon = \varepsilon_1$ .*

**Propagation of approximation guarantee  $1 - \delta$**  The propagation of the guarantee  $1 - \delta$  is quite straightforward. We assume approximations of different subformulas are going to be carried out in an independent manner, using different samples. Then, if:

$$\begin{cases} \Pr(|p_1 - \tilde{p}_1| \leq \varepsilon_1 p_1) \geq 1 - \delta_1 \\ \Pr(|p_2 - \tilde{p}_2| \leq \varepsilon_2 p_2) \geq 1 - \delta_2 \end{cases}$$

we have:

$$\Pr((|p_1 - \tilde{p}_1| \leq \varepsilon_1 p_1) \wedge (|p_2 - \tilde{p}_2| \leq \varepsilon_2 p_2)) \geq (1 - \delta_1)(1 - \delta_2).$$

This gives the propagation rule:  $1 - \delta = (1 - \delta_1)(1 - \delta_2)$ .

In ProApprox, in any case where we use approximations for both operands of one of the internal nodes of the evaluation tree, we simply set  $\delta_1 = \delta_2 = 1 - \sqrt{1 - \delta}$ .

**Multiplicative guarantee from additive approximations** Resorting to an execution plan using additive algorithms with the ultimate goal of producing a multiplicative tolerance, might sometimes be more efficient than running a multiplicative approximation. In that case, we would like the result to be within a multiplicative error interval  $[p - p\varepsilon, p + p\varepsilon]$  of the probability  $p$ , for a given  $\varepsilon$ . Thus, we need to set an input error  $\varepsilon_{\text{add}}$  for the additive algorithm, so that:

$$\varepsilon_{\text{add}} = \varepsilon p.$$

It is not possible to exactly determine this value since  $p$  is the quantity that we are actually looking for. What we propose in this case, is to use a *lower bound*  $\ell$  for  $p$  and set  $\varepsilon_{\text{add}}$  to  $\ell\varepsilon$ . This guarantees us that  $\varepsilon_{\text{add}}$  is small enough to obtain a multiplicative error of  $\varepsilon$ . One such simple lower bound, used in ProApproX, is the maximum probability of a clause of the DNF, that can be easily determined by a linear scan of the DNF. Using this trick, we can see Monte-Carlo sampling as a multiplicative approximation algorithm, with a new cost model of

$$\text{cost}_{\text{MulMC}} = C_{\text{AddMC}} \times \ln \frac{2}{\delta} \times \frac{L}{\ell^2 \varepsilon^2}.$$

## 7 Evaluation Plans

We now have all the necessary components for defining *evaluation plans* of a DNF formula  $\varphi$ . An  $(\varepsilon, \delta)$ -evaluation plan for  $\varphi$  is a subtree of the evaluation tree of  $\varphi$  presented in Section 5, with the same root, where every leaf is associated with an algorithm among the four described in Section 4, and, when this algorithm is an approximation algorithm, with approximation parameters  $\varepsilon'$ ,  $\delta'$ . We further require that these approximation parameters guarantee a multiplicative  $(\varepsilon, \delta)$ -approximation of  $\varphi$ , following the constraints of Section 6. The cost of a leaf of an evaluation plan is the cost of applying the given algorithm on the formula at that leaf. The cost of an evaluation plan is defined as the sum of the costs of all leaves (we thus ignore the cost of combining the probabilities at each internal node, which amounts to evaluating a constant number of arithmetic operations, see Section 5). There are many (indeed, infinitely many in general) evaluation plans for a given formula  $\varphi$ , corresponding to different subtrees of the evaluation tree, and to different assignments of approximation parameters. Our goal is to find one of minimum cost.

When propagating an error bound  $e$  from one node to its children, many possible different values for  $\varepsilon_1$  and  $\varepsilon_2$  can be used (e.g., under the condition  $\varepsilon_1 + \varepsilon_2 + \varepsilon_1\varepsilon_2 \leq \varepsilon$  for two approximated nodes linked by and independent disjunction, different values for  $\varepsilon_1$  and  $\varepsilon_2$  satisfy the condition). We can find the best assignment locally for the direct descendant of a node through finding the global minimum of the  $(\varepsilon_1, \varepsilon_2)$  function. However, this does not necessarily guarantee producing the most optimal propagation for the whole tree. Indeed, if we decide also to vary the allocation of these parameters, the number of possible plans is unbounded. Inspired by query optimization techniques, we can explore the space of evaluation plans for such a case by sampling a large number of them, evaluating the cost of each, and choosing the best one found, which is thus not necessarily optimal, but can be good enough in practice.

After experimenting with such a sampling technique, we decided instead to use in ProApproX a more deterministic approach: we first decide on an adequate allocation of the approximation parameters, and then look for the best evaluation plan with respect to these. Our exploration of the search space is based on the following principles. We start with the evaluation tree, and a global  $\varepsilon$ ,  $\delta$ :

1. First visit of the tree (top-down): We recursively assign  $\varepsilon$  and  $\delta$  values down the tree, following the logical operator at the direct parent node, and propagation rules from Section 6; we always set  $\varepsilon_1 = \varepsilon_2$  and  $\delta_1 = \delta_2$ ;
2. Second visit of the tree (bottom-up):
  - a) We assign the best evaluation algorithm for a given node, following local parameters and based on the cost model;

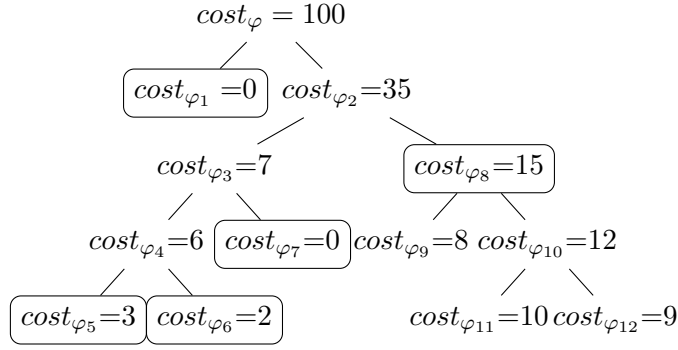


Figure 2: An example best evaluation plan in a deterministic case.

- b) If a given node is assigned with an exact algorithm, we re-assign unused  $\varepsilon$ ,  $\delta$  to nodes already assigned with an approximation algorithm, or to remaining nodes of the tree;
  - c) We restart (a,b) for the next node (including inner nodes);
3. Third visit of the tree (bottom-up): We then decide whether we want to chose to apply the algorithm at a given node (and forget about the algorithms assigned to its descendant) or at one of its ancestors.

**Example 3.** *Fig. 2 illustrates a mock-up example of a best evaluation plan, or best-tree, where computation nodes are marked with flags (framed nodes in the example compilation tree). Note that the algorithm for finding the best-tree in the deterministic case amounts to comparing the cost of a node with that of its children (a child node is never more expensive than its parent; however, the sum of children costs can exceed a parent-node cost).*

As we shall see, proceeding this way leads to better results than always choosing, say, the full evaluation tree.

## 8 Experiments

In this section we evaluate our approach experimentally. We address the following questions:

- (1) How does our new query evaluation method compare to the current state-of-the-art probabilistic XML database query processing?
- (2) How accurate are the results?
- (3) How effective is the tree compilation of the DNF formula over simply performing an exact or an approximate computation of the whole formula?

**Experimental Setup** We implemented our system in Java. XQuery lineage queries, translated from the original XPath queries, are evaluated over the XML dataset through the XML::DB API for BaseX, the native XML DBMS. All indexing features of BaseX (in particular, path summaries and text index) were turned on. The resulting DNF formulas are then compiled into trees, potential evaluation strategies are explored and their cost assessed, and the chosen strategy is then run to evaluate the probability of the query. Figure 3 illustrates the architecture of the system:

- (1) The user (XPath) query is translated into the lineage query expressed in XQuery language: Section 3;
- (2, 3) The lineage query is processed by the native DBMS BaseX over the database (PrXML<sup>cie</sup> p-documents stored as ordinary XML documents with *prob* attributes);
- (4) The gathered probabilistic lineage is compiled by ProApprox:

**Lineage preprocessing:** optimizations over the DNF (removing subsumed or invalid clauses, etc.)

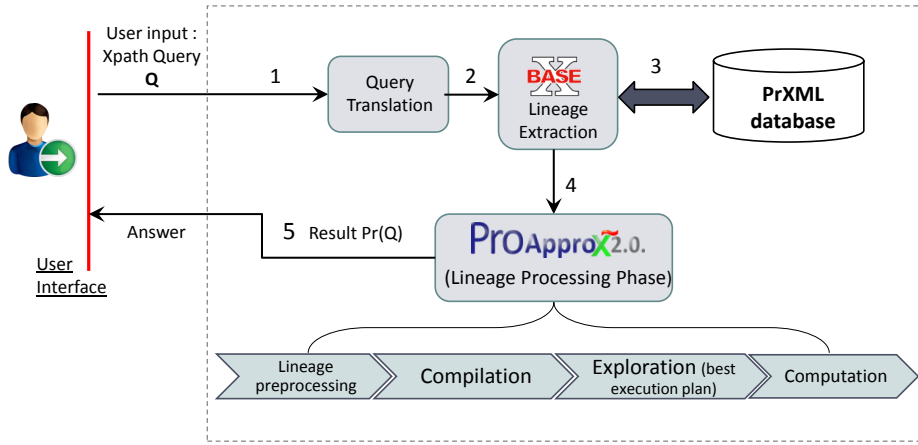


Figure 3: The ProApproX 2.0 architecture.

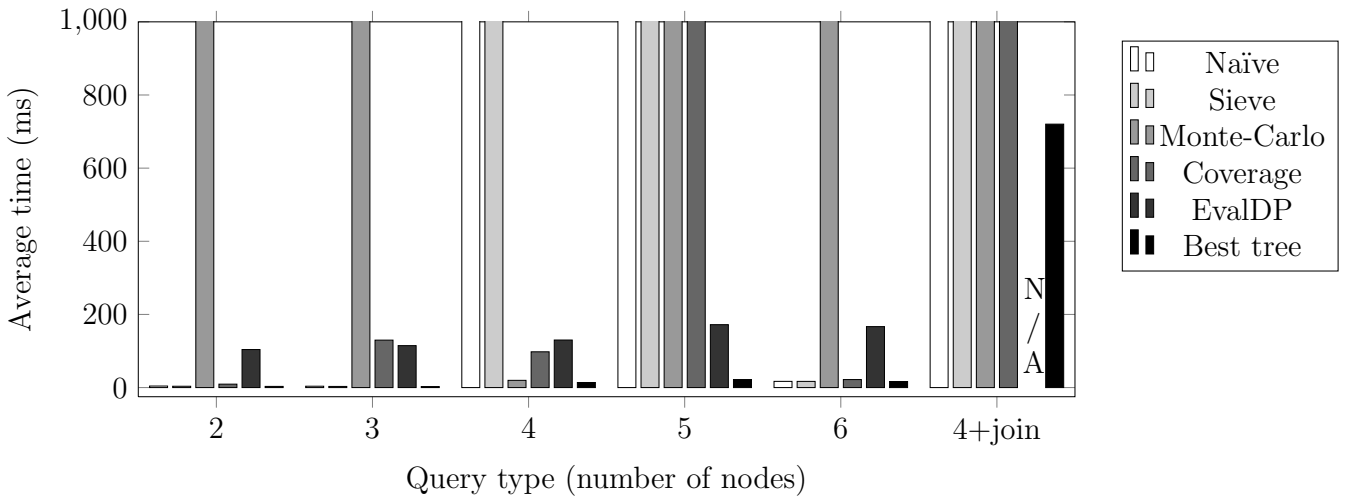


Figure 4: Running time of the different algorithms on the MondialDB dataset [20, 21]

**Compilation:** Section 5

**Exploration:** Section 7

**Computation:** following the algorithms set at the computation nodes of the best evaluation tree

(5) The probability of the user query is then displayed.

Our experiments were run on a desktop PC with an Intel Xeon CPU at 2.40 GHz running 64 bit Linux, with the data accessed by BaseX located on a standard magnetic hard drive.

**Data and Queries** Our evaluation was carried out on three different datasets, two of them previously used in the probabilistic XML literature, the third one a synthetic dataset.

1. We used the *MondialDB* probabilistic XML document with local dependencies and 15 tree-pattern queries that served to evaluate EvalDP [20, 21]. EvalDP is a linear-time bottom-up algorithm for evaluating tree-pattern queries without joins on local dependency p-documents. The queries have between two to six nodes depending on their complexity and selectivity. The p-document, also provided by the authors of [21], has been created by adding random probabilistic choices in a regular XML document. We added an extra query with a join (`/mondial//country[.//name=@name]`) to showcase the fact that ProApproX is able to deal with join and no-join queries indifferently.

2. We also obtained the *Movies* dataset and queries, used by Hollander and Van Keulen in [16], from the authors of this work. [16] explores the feasibility of storing and querying probabilistic XML documents in a probabilistic relational database. The authors considered real-life uncertain data obtained from a probabilistic data integration application [35], that produces a local-dependency p-document. The application integrates movie data from a TV guide<sup>2</sup> with that of the Internet Movie Database<sup>3</sup>. As in [16], the query was ran on five p-documents of different sizes.
3. Finally, we tested the performance of our DNF probability evaluator over *synthetic* data: random DNF formulas with a number of literals  $L$  ranging from 1 to 12,000. Specifically, for a given target size number of literals, we iteratively randomly drew clauses in the following manner, until the target size was reached: first, draw a random clause size uniformly at random between 1 and  $[L/3, L/10, \text{ and } L/100]$ ; second, uniformly draw positive or negative literals from a fixed set of variables (of size  $L/2$ ) avoiding contradicting literals.

**Methodology** Each query is translated to an XQuery lineage query that is evaluated on BaseX to capture the DNF lineage. The time of this operation is what we call *XQuery time*. To compare our approach to baselines, we first run the four main computation algorithms over the entire DNF lineage: sieve and naïve algorithms, Monte-Carlo sampling, and the self-adjusting coverage algorithm. The computation is stopped at one minute if not finished before. Afterwards, we fully compile the DNF in a tree structure. We then use our cost model to determine what the best strategy to evaluate the probability of each leaf. We call this the *full-tree* strategy. Finally, we used the exploration strategy detailed in Section 7 to find, if possible, a tree whose cost is lower than both the cost of evaluating the whole DNF with any of the algorithms and the cost of the full tree strategy. This is the *best-tree* strategy, for which we can distinguish between the *Exploration time* (to decide of the best tree and of the best allocation of error bounds and evaluation algorithms) and the *Evaluation time* itself, when the actual evaluation algorithms are run and the final probability is obtained. The running time of the best-tree strategy, ProApprox’s default behavior, is thus obtained by summing *XQuery time*, *Exploration time*, and *Evaluation time*. In all experiments, we set the parameters for a multiplicative approximation with  $\varepsilon = 0.1$  and a reliability factor  $\delta = 95\%$ .

**Comparison with EvalDP** We compare with EvalDP, using the implementation kindly provided by the authors of [21]. We show the results of our system, baselines and EvalDP on the 16 queries, averaged by type, in Figure 4. The *best tree* strategy of ProApproX has highly competitive execution time, being in all cases under 1s and the fastest probability computation method, sometimes tied with one of the baselines. The best baseline algorithm depends on the query: when the lineage is simple (types 2, 3, or 6) exact naïve or sieve computation methods are extremely fast, and this is the choice that our optimizer selects. On query type 4, Monte-Carlo approximations are fast (this is usually the sign that we have a high lower bound on the probability value, making it easy to derive a multiplicative approximation from an additive approximation). On query type 5 and the join query, interestingly enough, the best-tree strategy is significantly faster than all baselines, meaning that our evaluation tree is useful and that different evaluation strategies need to be applied to different nodes of the tree.

Unlike EvalDP, whose performance is linear in the size of the data (and remains in our experiments relatively uniform for all queries), the performance of ProApproX is obviously related to both the size and the complexity of the lineage DNF. The part that is directly related to the size of the data in ProApproX is the XQuery time recorded by BaseX for lineage query evaluation. But the part that is sensitive to the size of the DNF lineage is the compilation time. As shown in Table 3, the majority of the time is spent in XQuery evaluation for all queries on this dataset. For queries with larger DNFs, the *Compilation time* becomes relatively important. The time required for the Exploration

---

<sup>2</sup><http://www.tvguide.com/>

<sup>3</sup><http://www.imdb.com/>

Query type	Average DNF size	XQuery	Compilation	Exploration	Evaluation
2	6	96.82%	2.47%	0.19%	0.52%
3	5	98.19%	1.38%	0.08%	0.34%
4	43	97.41%	2.01%	0.34%	0.24%
5	252	64.47%	33.76%	1.43%	0.35%
6	6	99.69%	0.29%	0.02%	0.01%
4+join	3656	61.49%	36.12%	1.95%	0.44%

Table 3: Proportion of the time spent on each part of the probabilistic XML query evaluation on the MondialDB dataset for the best tree strategy

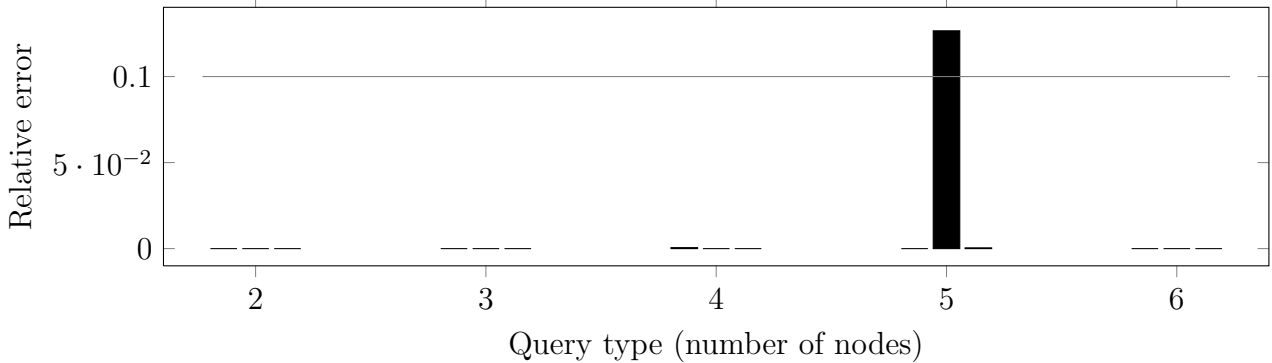


Figure 5: Relative error on the probabilities computed by the algorithm on the MondialDB over each non-join query with respect to the exact probability values ( $\varepsilon = 0.1$ ,  $\delta = 95\%$ )

and Evaluation phases remain negligible in all cases. Finally recall that EvalDP is neither able to deal with join queries, nor with long-distance dependencies, contrarily to ProApproX.

With respect to SPROUT, we recorded one order of magnitude improvements on the largest DNFs. For the last categories of queries, run times for SPROUT were of the order of several seconds, while ProApproX never exceeded 100 milliseconds.

**Accuracy of Results** Obviously, a fast algorithm is of no use if the results obtained are not accurate, something to be wary of in the case of ProApproX since it relies on approximation algorithms. To evaluate the quality of results generated by ProApproX, we based the comparison on the relative error distance between our results and the exact probability computed by the naive algorithm (except in one case where EvalDP returns a wrong answer, probably due to a bug in the code, where we had to compute the exact answer using the naïve algorithm). We left the join query out of this experiment, since it was too costly to compute its probability in an exact manner. Figure 5 is the result on one particular run of the algorithm (since we use sampling, errors tend to be averaged out if we consider multiple runs). The figure shows that for most results ProApproX achieves high accuracy. For 14 out of 15 queries, we record an accuracy far below our relative error margin of 0.1. For one of the 15 queries, we have an error slightly higher than 0.1, but this is in line with what is expected: with a probabilistic guarantee  $\delta$  of 95% and for 15 queries, there is a chance of  $1 - 0.95^{15} = 53.6\%$  that we produce, at least once, a relative accuracy beyond 0.1.

**Performance over the Movies dataset** For the three queries of this dataset, the run time on the five p-documents was recorded to be almost one order of magnitude lower than the performances presented by the authors in [16] using Trio, a relational probabilistic DBMS [25]. In Figure 6, we plot the performance over the join query  $Q_3$ ; points missing in the figure mean runtimes greater than 5s. The score for Trio is plotted on the figure, though it is not directly comparable, as we did not

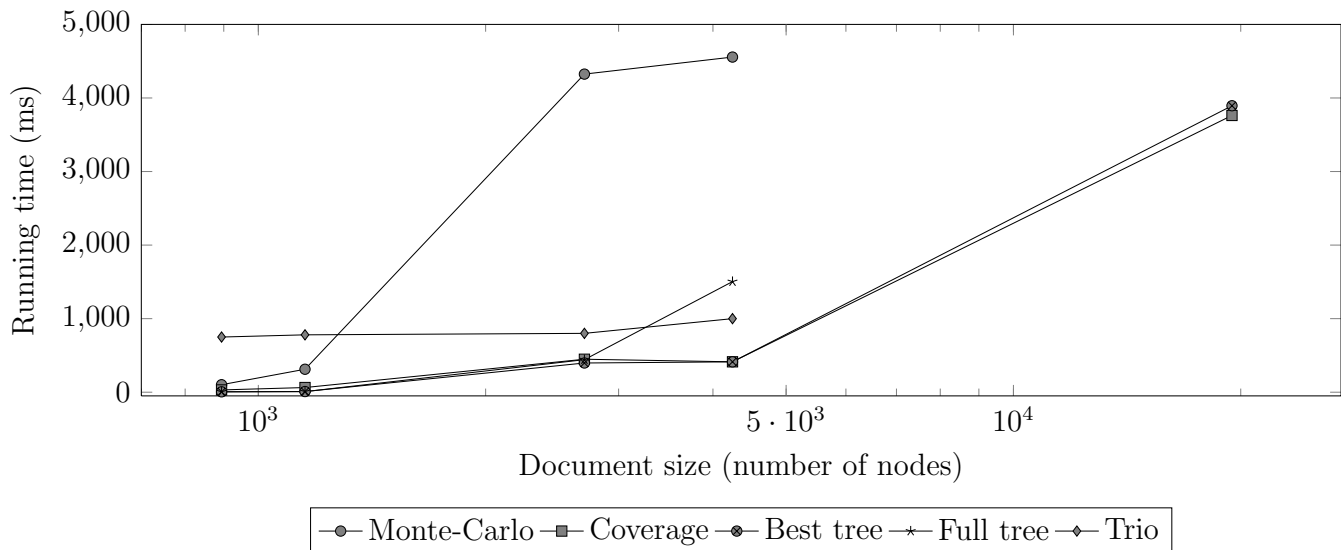


Figure 6: Running time of the different algorithms on query  $Q_3$  of the movie dataset [16]; times greater than 5s are not shown and times for Trio are estimated from [16]

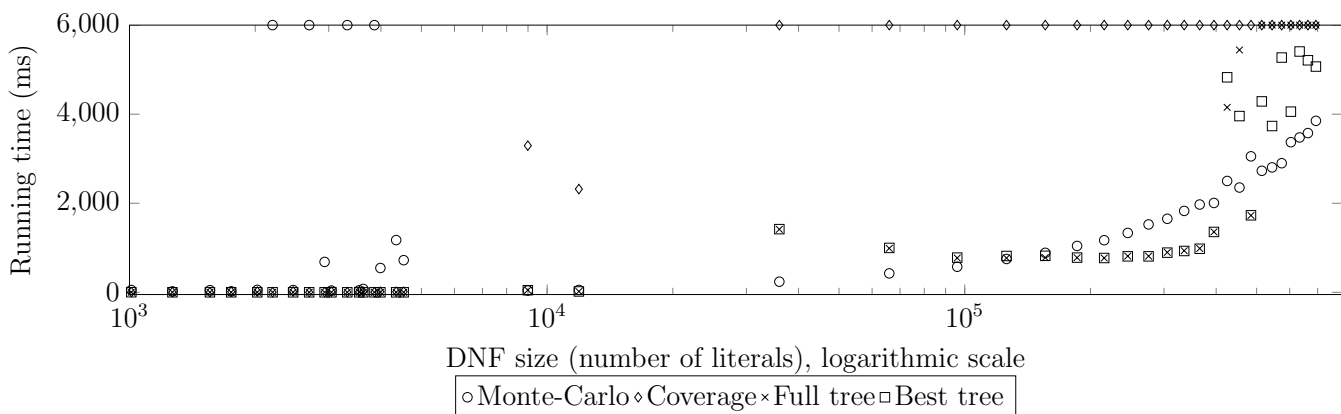


Figure 7: Running time of the different algorithms on the synthetic dataset

reimplement the translation to Trio but took the numbers directly from [16]. Still, it is interesting here to note that Trio has a performance of about 40 seconds on the biggest p-document (19,475 nodes), while ProApprox took around 4 seconds.

Figure 6 also compares the best-tree strategy to the full-tree strategy (always evaluating the leaves of the evaluation tree) and shows that the latter is not always efficient. In this particular example, we see that it is always best to apply the coverage algorithm to the whole formula. We have seen, in Figure 4 that for some other datasets, the coverage algorithm is not efficient. The point here is that by using a cost optimizer, we are able *in all cases* to find the most efficient evaluation strategy (or one close to it).

**Performance over Synthetic Data** Fig. 7 presents the performance of different evaluation methods: the additive Monte-Carlo, the Self-Adjusting coverage algorithm, evaluation over a fully compiled tree (Full tree), and evaluation over the best execution plan (Best tree) that ProApprox chooses for a DNF probability computation. Note that times shown on the top line are actually anywhere beyond. The running time of the self-adjusting coverage algorithm is directly proportional to the size of the formula (i.e., appears as an exponential in a log-linear plot). This is not the case, however, for Monte-Carlo sampling because of the use of the lower bound on the DNF probability: when this value is relatively high (apparently the case for most of the big random DNFs of Fig. 7), Monte-Carlo sampling can be very efficient. On the other hand, when the lower bound is low, the Monte Carlo algorithm can be quite ineffective, as is shown by the irregular variation recorded for this algorithm over DNFs of different sizes. Our best-tree strategy significantly outperforms the self-adjusting coverage algorithm



and gives very low running times (of usually less than 100ms) for DNFs holding until  $10^4$  variables. It also outperforms Monte Carlo for most of the cases, especially for DNFs with low probabilities. However, when the data scales (DNFs with sizes larger than  $\sim 10^{5.6}$ ), Monte Carlo may record better time, which is due to the compilation time needed by ProApprox before deciding on the best strategy. Yet the difference in running time remains reasonable. Interestingly, for very large formulas, the best-tree strategy does a better job than the full-tree strategy. Most importantly, the best tree is *always* under 6s, even though it is sometimes not the quickest (cf. Monte-Carlo), which empirically demonstrates that it is much more robust.

## 9 Related Work

**Probabilistic data management** Managing probabilistic and uncertain data is a topic of much current interest in database research, and there has been a significant amount of work under the relational database scheme. Extensive literature and many systems and prototypes have been introduced, in particular MystiQ [8], Trio [25], Orion (previously known as U-DBMS) [32], and MayBMS [5, 18].

SPROUT [17], the query engine integrated in the probabilistic relational database MayBMS [18], considers query evaluation for tuple independent probabilistic databases (i.e, restricted databases [29]) and Boolean conjunctive queries with inequalities [28] but without self-joins [11]. For tractable queries, the confidence computation is based on OBDDs [27]. If a DNF lineage is captured, it is compiled to a novel kind of decision diagrams called *d-trees*, exploiting negative correlations (inconsistency), independence, and factored representations. The authors pointed out the need for more advanced and more accurate – but still efficient – techniques for estimating lower and upper bounds of probability values.

Later work by R. Fink and D. Olteanu [12] on SPROUT employs novel exact and approximate knowledge compilation techniques that compute probabilities of events for results of arbitrary relational algebra queries. The main component is an incremental compilation algorithm that repeatedly decomposes a propositional formula  $\varphi$  into sub-formulas such that lower and upper bounds of the probability of  $\varphi$  can be efficiently computed from bounds for the sub-formulas. The algorithm is run for a given time budget or until the desired approximation is reached [12]. The authors also addressed the problem of approximating lineage formulas that are hard to evaluate with read read-once formulas, where every variable occurs at most once, or read-once formulas in disjunctive normal form [13].

**Probabilistic XML** Meanwhile, the semistructured model has received less attention. An in-depth study on the compared expressiveness and efficiency of probabilistic XML models was elaborated by Abiteboul, Kimelfeld, Sagiv, and Senellart in [3]. Subsequent research has extended the model with continuous probabilistic distributions [2], constraints [9], or possible trees of unbounded size [6].

Early work on query evaluation over p-documents was introduced by Kimelfeld, Koscharovsky, and Sagiv in [20, 21] producing the EvalDP algorithm that natively processes positive tree-pattern queries, without join. The algorithm is restricted to run only over p-documents with local dependencies, i.e., where dependency can only happen between children of a same parent node. This linear-time bottom-up processing of the tree has first-rate performance, and has been generalized to arbitrary monadic second-order queries through tree automata techniques [6, 10], but is not applicable to queries or data that involve correlations.

In [21], Kimelfeld, Koscharovsky, and Sagiv proposed as a proof of the existence of a randomized polynomial-time approximation scheme an alternative sampling algorithm to obtain a multiplicative approximation of a Boolean query. The probability of query  $Q$  being true in a random instance  $\mathcal{P}$  can be computed [21] by:  $\sum_{i=1}^n \Pr(m_i \triangleleft \mathcal{P}) \times \Pr\left(\bigwedge_{j=1}^{i-1} \neg(m_j \triangleleft \mathcal{P}) \mid m_i \triangleleft \mathcal{P}\right)$ , where  $\Pr(m_i \triangleleft \mathcal{P})$  is the probability that a match  $m_i$  of  $Q$  remains in  $\mathcal{P}$ . The first term is easy to compute by just gathering the probabilities of all events involved in the match; the second term can be approximated by conducting biased draws to considerate the probability that none of the preceding matches exists

in the random document  $\mathcal{P}$ , given that a current match  $m_i$  appears in that same document. We implemented this biased approximation in [31], a demonstration of an early version of ProApproX without any form of cost estimation, and the evaluation leads to very good accuracy. Nevertheless, the convergence guarantee [21] that is obtained from Hoeffding’s inequality requires a running time growing in  $O(n^3 \ln n)$  in the number of matches (i.e., in the number of clauses of the lineage formula).

**DNF probability estimation** We replaced this method by the self-adjusting coverage algorithm, a linear multiplicative estimation of the number of assignments satisfying a DNF, and an extension of the earlier *coverage algorithm* for the *union of sets* problem, proposed by Karp, Luby and Madras [19] (with the background motivation of managing network reliability issues). The DNF counting problem is the special case of the DNF probability problem when all variable probabilities are 0.5, in which case  $\Pr(\varphi) = \frac{\#\varphi}{2^N}$  (with  $\#\varphi$  the number of satisfying assignments of  $\varphi$ ). Both coverage and self-adjusting coverage algorithms for the DNF counting problem can easily be modified for the DNF probability problem, yielding  $(\varepsilon, \delta)$ -approximation algorithms, as shown in a subsequent work by Luby and Velickovic [24], and as deployed in ProApproX.

**Systems** At the practical level, a full-fledged probabilistic semistructured data management system is still missing. Hollander and van Keulen [16] investigate the adequacy of probabilistic relational databases for querying probabilistic semistructured data by mapping XML to relational data. Queries can then be transformed into relational queries and evaluated using a system such as Trio, for which efficient query evaluation algorithms have been developed, both exact and approximate. The downside of the approach is that relational databases do not exploit the specific characteristics of tree-like data encoded into databases, that for instance makes tree-pattern queries over local models tractable [21].

Initial ideas leading to this work were presented in [31, 33] though the specificity of our approach, using a cost estimator to optimize the running time of the (approximate) probability computing, is fully novel.

## 10 Conclusions

We have introduced an original optimizer-like approach to evaluating query results over probabilistic XML. Though some components are specific to this setting (in particular, the translation from XPath queries to XQuery lineage queries), the core of our system aims at solving the very general problem of computing (approximate) probabilities of propositional formulas built out of independent random variables, in an efficient manner. Our formula decomposition technique, and some of the algorithms used, assume that the formula is in DNF, but the technique could probably be extended to arbitrary formulas.

The first principle our approach relies on, as well as the main observation from our experiments, is that *the optimal probability evaluation algorithm to use depends on the characteristics of the formula*: if the formula has few variables, go with the naïve algorithm; if it has few clauses, with the sieve algorithm; if the probability is known to be close to 1, with Monte-Carlo sampling; if the formula was obtained from evaluating a tree-like query over a tree-like structure, use a technique such as EvalDP; if nothing else works, use the self-adjusting coverage algorithm. Our cost model captures this. The second principle is that *different algorithms can be used to evaluate the probability of different parts of a formula*. Our formula decomposition technique, and our exploration of possible evaluation trees, makes use of this.

In ProApproX, the production of the lineage from the original p-document and query, and the evaluation of the probability of this lineage are fully independent. Each of these problem can thus be optimized independently (using, respectively, the XML query optimization of a native XML DBMS, and our formula probability evaluation technique). However, it is likely that exploiting the structure of the query to obtain lineage that is already factorized would result in even more efficient evaluation. This is an obvious direction for future research.

Other potential improvements of interest would be (i) to refine our cost model with non-linear functions of the size of the formulas, to take into account second-order terms in the complexity of the algorithms; (ii) to explore the space of evaluation trees in a more intelligent manner, taking into account the fact that, for instance, it is always better to apply the naïve algorithm on two children of a node than on the node itself – note that this observation is not true in general for approximation algorithms, so this requires non-trivial analyses of the cases; (iii) to keep the same  $\delta$ 's parameter across the whole tree by performing correlated samples in all branches of the tree – this is not trivial either, since we need to keep track of all partial samples performed; (iv) to exploit the fact that most evaluation algorithms scale effortlessly to distribute the probability computation over multi-core or distributed architectures.

## 11 Acknowledgments

This work has been partly supported by the FP7 ARCOMEM and ANR Dataring projects. We would like to thank B. Kimmelfeld, Y. Kosharovskiy, and Y. Sagiv for providing the EvalDP system and dataset [20,21], as well as E. Hollander and M. van Keulen for the Movies dataset [16].

## References

- [1] T. Abdesslem, M. L. Ba, and P. Senellart. A probabilistic XML merging tool. In *Proc. EDBT*, 2011.
- [2] S. Abiteboul, T.-H. H. Chan, E. Kharlamov, W. Nutt, and P. Senellart. Aggregate queries for discrete and continuous probabilistic XML. In *Proc. ICDT*, 2010.
- [3] S. Abiteboul, B. Kimelfeld, Y. Sagiv, and P. Senellart. On the expressiveness of probabilistic XML models. *VLDB J.*, 2009.
- [4] S. Abiteboul and P. Senellart. Querying and updating probabilistic information in XML. In *Proc. EDBT*, 2006.
- [5] L. Antova, C. Koch, and D. Olteanu. Query language support for incomplete information in the MayBMS system. In *Proc. VLDB*, 2007.
- [6] M. Benedikt, E. Kharlamov, D. Olteanu, and P. Senellart. Probabilistic XML via Markov chains. *Proc. VLDB Endowment*, 2010.
- [7] M. Benedikt and C. Koch. XPath leashed. *Proc. ACM*, 2008.
- [8] J. Boulos, N. N. Dalvi, B. Mandhani, S. Mathur, C. Ré, and D. Suciu. MYSTIQ: a system for finding more answers by using probabilities. In *Proc. SIGMOD*, 2005.
- [9] S. Cohen, B. Kimelfeld, and Y. Sagiv. Incorporating constraints in probabilistic XML. *Proc. ACM*, 2009.
- [10] S. Cohen, B. Kimelfeld, and Y. Sagiv. Running tree automata on probabilistic XML. In *Proc. PODS*, 2009.
- [11] N. N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. *VLDB J.*, 2007.
- [12] R. Fink, A. Hogue, D. Olteanu, and S. Rath. Sprout<sup>2</sup>: a squared query engine for uncertain web data. In *Proc. SIGMOD*, 2011.
- [13] R. Fink and D. Olteanu. On the optimal approximation of queries using tractable propositional languages. In *Proc. ICDT*, 2011.
- [14] K. D. Heidtmann. Improved method of inclusion-exclusion applied to  $k$ -out-of- $n$  systems. *Proc. IEEE*, 1982.
- [15] W. Hoeffding. Probability inequalities for sums of bounded random variables. *J. American Statistical Association*, 1963.
- [16] E. Hollander and M. van Keulen. Storing and querying probabilistic XML using a probabilistic relational DBMS. In *Proc. MUD*, 2010.
- [17] J. Huang. Design and implementation of the SPROUT query engine for probabilistic databases. Master's thesis, University of Oxford, 2009.

- [18] J. Huang, L. Antova, C. Koch, and D. Olteanu. MayBMS: a probabilistic database management system. In *Proc. SIGMOD*, 2009.
- [19] R. M. Karp, M. Luby, and N. Madras. Monte-Carlo approximation algorithms for enumeration problems. *Algorithms J.*, 1989.
- [20] B. Kimelfeld, Y. Kosharovsky, and Y. Sagiv. Query efficiency in probabilistic XML models. In *Proc. SIGMOD*, 2008.
- [21] B. Kimelfeld, Y. Kosharovsky, and Y. Sagiv. Query evaluation over probabilistic XML. *VLDB J.*, 2009.
- [22] B. Kimelfeld and Y. Sagiv. Matching twigs in probabilistic XML. In *Proc. VLDB*, 2007.
- [23] B. Kimelfeld and P. Senellart. Probabilistic XML: Models and complexity. In Z. Ma, editor, *Advances in Probabilistic Databases for Uncertain Information Management*. Springer-Verlag, 2012. To appear.
- [24] M. Luby and B. Velickovic. On deterministic approximation of DNF. *Algorithmica J.*, 1996.
- [25] M. Mutsuzaki, M. Theobald, A. de Keijzer, J. Widom, P. Agrawal, O. Benjelloun, A. D. Sarma, R. Murthy, and T. Sugihara. Trio-One: Layering uncertainty and lineage on a conventional DBMS. In *Proc. CIDR*, 2007.
- [26] A. Nierman and H. V. Jagadish. ProTDB: Probabilistic data in XML. In *Proc. VLDB*, 2002.
- [27] D. Olteanu and J. Huang. Using OBDDs for efficient query evaluation on probabilistic databases. In *Proc. SUM*, 2008.
- [28] D. Olteanu and J. Huang. Secondary-storage confidence computation for conjunctive queries with inequalities. In *Proc. SIGMOD*, 2009.
- [29] D. Olteanu, J. Huang, and C. Koch. Sprout: Lazy vs. eager query plans for tuple-independent probabilistic databases. In *ICDE*, 2009.
- [30] P. Senellart and S. Abiteboul. On the complexity of managing probabilistic XML data. In *Proc. PODS*, 2007.
- [31] P. Senellart and A. Souihli. ProApproX: a lightweight approximation query processor over probabilistic trees. In *Proc. SIGMOD*, 2011.
- [32] S. Singh, C. Mayfield, S. Mittal, S. Prabhakar, S. E. Hambrusch, and R. Shah. Orion 2.0: native support for uncertain data. In *Proc. SIGMOD*, 2008.
- [33] A. Souihli. Efficient query evaluation over probabilistic XML with long-distance dependencies. In *Proc. EDBT/ICDT*, 2011.
- [34] L. G. Valiant. The complexity of enumeration and reliability problems. *j. SIAM*, 1979.
- [35] M. van Keulen and A. de Keijzer. Qualitative effects of knowledge rules and user feedback in probabilistic data integration. *VLDB J.*, 2009.
- [36] M. van Keulen, A. de Keijzer, and W. Alink. A probabilistic XML approach to data integration. In *Proc. ICDE*, 2005.