

Parameterized and Fine-Grained Analysis of Query Evaluation Over Bag PDBs

Su Feng ✉

Illinois Institute of Technology, Chicago, USA

Boris Glavic ✉

Illinois Institute of Technology, USA

Aaron Huber ✉

University at Buffalo, USA

Oliver Kennedy ✉

University at Buffalo, USA

Atri Rudra ✉

University at Buffalo, USA

Abstract

The problem of computing the marginal probability of a tuple in the result of a query over set-probabilistic databases (PDBs) is a fundamental problem in set-PDBs. In this work, we study the analog problem for bag semantics: computing a tuple's expected multiplicity exactly and approximately. We are specifically interested in the fine-grained complexity and how it compares to the complexity of deterministic query evaluation algorithms — if these complexities are comparable, it opens the door to practical deployment of probabilistic databases. Unfortunately, our results imply that computing expected multiplicities for Bag-PDBs based on the results produced by such query evaluation algorithms introduces super-linear overhead (under parameterized complexity hardness assumptions/conjectures). We proceed to study approximation of expected multiplicities of result tuples of positive relational algebra queries (\mathcal{RA}^+) over c -TIDBs and for a non-trivial subclass of block-independent databases (BIDBs). We develop a sampling algorithm that computes a $(1 \pm \epsilon)$ -approximation of the expected multiplicity of an output tuple in time linear in the runtime of a comparable deterministic query for any \mathcal{RA}^+ query.

2012 ACM Subject Classification Information systems → Incomplete data

Keywords and phrases PDB, bags, polynomial, boolean formula, etc.

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

1 Introduction

This work explores the problem of computing the expectation of a tuple's multiplicity in a bag TIDB. Our analysis specifically considers a restricted form of bag TIDB, which we call a c -TIDB. A c -TIDB, $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ encodes a bag of uncertain tuples such that each tuple models c disjoint events, where each such set of disjoint events is itself independent of the others. A tuple in \mathcal{D} has a multiplicity of at most c . The set of all worlds is encoded in $\{0, \dots, c\}^D$, which is the set of all vectors of length $|D|$ such that each index corresponds to a distinct $t \in D$ storing its multiplicity. \mathcal{P} is the product distribution over the set of all worlds. A given world $\mathbf{W} = \{0, \dots, c\}^D$ can be interpreted for each $\mathbf{W}[t] = j$ as denoting that tuple t_i appears j times in world \mathbf{W} for $j \in [0, c]$. The resulting product distribution can then be expressed across the n base tuples of the encoding as $p_{i,j} = Pr[\mathbf{W}[t] = j]$, where each distribution is independent for $t \in [n]$. Allowing for $\leq c$ multiplicities across all tuples gives rise to having $\leq (c+1)^n$ possible worlds instead of the usual 2^n possible worlds of the traditional set TIDB. In this work, it is natural to be specifically considering bag query

By def implies disjoint property

important special case a as follows

Do we really need this explanation. Just jump to multiplicity?

the mult. of t in

Better to say that \mathcal{J} -TIDB is the traditional set TIDB.

why is it natural?

23:2 Bag PDB Queries

44 semantics.

45 We can formally state this problem as:

46 ► **Problem 1.1.** Given a c -TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, \mathcal{RA}^+ query Q , and result tuple t ,
 47 compute the expected multiplicity of t : $\mathbb{E}_{\mathcal{D} \sim \mathcal{P}} [Q(\mathcal{D})(t)]$.

Aaron says: I think we use \mathbf{D} to denote something different in one of the proofs. Have to keep an eye open for this to avoid overloading notation.

49 We upperbound the multiplicity of tuples in a c -TIDB since this is what typically seen in
 50 practice. Allowing for unbounded c is an interesting open problem.

51 **Hardness of Set Query Semantics and Bag Query Semantics.** Set query evaluation
 52 semantics over 1-TIDBs have been studied extensively, and the data complexity of the
 53 problem in general has been shown by Dalvi and Suicu to be $\#P$ -hard [13]. For our setting,
 54 there exists a trivial algorithm to compute problem 1.1 for any query over a c -TIDB due to
 55 linearity of expectation. Simply perform the probability computations in a 'sum-of-products'
 56 fashion. This is made more precise when we discuss polynomial equivalence in the following
 57 subsection. Since we can compute problem 1.1 in polynomial time, the interesting question
 58 that we explore deals with hardness of computing expectation using fine-grained analysis
 59 and parameterized complexity.

60 **1.1** One of the main theoretical points in this work is to discern whether or not bag c -TIDB
 61 query semantics are indeed linear in the runtime of an equivalent deterministic query. If this is
 62 true, then this would open up the way for deployment of c -TIDBs in practice. Unfortunately,
 63 we prove that this is not the case. To analyze this question we denote by $T^*(Q, \mathcal{D})$ the
 64 optimal runtime complexity of computing problem 1.1 over c -TIDB \mathcal{D} . Let $\bar{\mathcal{D}}$ denote the
 65 set of tuples in \mathcal{D} , i.e.,

66 ► **Definition 1.2** ($\bar{\mathcal{D}}$). Define $\bar{\mathcal{D}}$ to be the set of tuples appearing across all the possible
 67 worlds of a c -TIDB, formally $\bar{\mathcal{D}} = \{t_i \mid \forall \mathbf{W} \in \{0, \dots, c\}^D, \forall i \in [D] : \mathbf{W}[i] > 0\}$. When a
 68 specific $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ is being referred to, we will use \mathcal{D} to denote the set of tuples.

69 Let $T_{det}^*(Q, \bar{\mathcal{D}})$ be the optimal runtime (with some caveats; discussed in sec. 2.4) of query
 70 Q on a comparable deterministic database $\bar{\mathcal{D}}$ defined next.

71 Table 1 shows our lower bounds for computing problem 1.1 on c -TIDBs.

Our lower bound results. In table 1 we show that depending on what hardness

| Lower bound on $T^*(Q, \mathcal{D})$ | Num. \mathcal{P} s | Hardness Assumption |
|--|----------------------|-------------------------------|
| $\Omega \left((T_{det}^*(Q, \bar{\mathcal{D}}))^{1+\epsilon_0} \right)$ for some $\epsilon_0 > 0$ | Single | Triangle Detection hypothesis |
| $\omega \left((T_{det}^*(Q, \bar{\mathcal{D}}))^{C_0} \right)$ for all $C_0 > 0$ | Multiple | $\#W[0] \neq \#W[1]$ |
| $\Omega \left((T_{det}^*(Q, \bar{\mathcal{D}}))^{c_0 \cdot k} \right)$ for some $c_0 > 0$ | Multiple | Conjecture 3.2 |

Table 1 Our lower bounds for a specific hard query Q parameterized by k . The \mathcal{D} is over the same (family of) $\bar{\mathcal{D}}$ and those with 'Multiple' in the second column need the algorithm to be able to handle multiple \mathcal{P} s (for a given $\bar{\mathcal{D}}$). The last column states the hardness assumptions that imply the lower bounds in the first column (ϵ_0, C_0, c_0 are constants that are independent of k).

72
 73 result/conjecture we assume, we get various emphatic versions of *no* as an answer to our
 74 question. To make some sense of the other lower bounds in Table 1, we note that it is not
 75 too hard to show that $T^*(Q, \mathcal{D}) \leq O \left((T_{det}^*(Q, \bar{\mathcal{D}}))^k \right)$, where k is the largest degree of the

There is no notion of "degree" of a query as far as I know.

on $T^*(Q, \mathcal{D})$ for

of computing expected result tuple our mult.

Does this need to be defined since \mathcal{D} is not a set anymore

Re-state this as saying we'll achieve $c = O(1)$.

I don't think so

can if Problem be solved in time

This should also include c as a param

+ PROPAGATE

poly time

specifically, we ask

Also need to justify why computing expected mult. is imp/ useful
 explain that we are interested in the exponent of poly runtime.

This change needs to be PROPAGATED.

I.e. many of the subsequent $\bar{\mathcal{D}}$ should say our \mathcal{D} is \mathcal{D}
 At this point $T^*(Q, \mathcal{D}) \leq O(T_{det}^*(Q, \bar{\mathcal{D}})^k)$

$k = (d_0, 1)^D ; \mathcal{P}$

Spent a sentence or two on the implication of this result.

$$\begin{aligned} \Phi[\pi_A(Q), \bar{D}, t] &= \sum_{t': \pi_A(t')=t} \Phi[Q, \bar{D}, t'] & \Phi[Q_1 \cup Q_2, \bar{D}, t] &= \Phi[Q_1, \bar{D}, t] + \Phi[Q_2, \bar{D}, t] \\ \Phi[\sigma_\theta(Q), \bar{D}, t] &= \begin{cases} \Phi[Q, \bar{D}, t] & \text{if } \theta(t) \\ 0 & \text{otherwise.} \end{cases} & \Phi[Q_1 \bowtie Q_2, \bar{D}, t] &= \Phi[Q_1, \bar{D}, \pi_{attr(Q_1)}t] \\ & & & \cdot \Phi[Q_2, \bar{D}, \pi_{attr(Q_2)}t] \\ & & \Phi[R, \bar{D}, t] &= X_t \end{aligned}$$

Figure 1 Construction of the lineage (polynomial) for an \mathcal{RA}^+ query over a c -TIDB, where \mathbf{X} consists of all X_t over all R in \bar{D} and t in R . Here $\bar{D}.R$ denotes the instance of relation R in \bar{D} . Please note, after we introduce the reduction to 1-BIDB, the base case will be expressed alternatively.

this is not a standard term I think

not defined

76 query Q (i.e., join width) over all result tuples t (and the parameter that defines our family
77 of hard queries).

78 What our lower bound in the third row says is that one cannot get more than a polynomial
79 improvement over essentially the trivial algorithm for problem 1.1. However, this result
80 assumes a hardness conjecture that is not as well studied as those in the first two rows of the
81 table (see Sec. 3 for more discussion on the hardness assumptions). Further, we note that
82 existing results already imply the claimed lower bounds if we were to replace the $T_{det}^*(Q, \bar{D})$
83 by just $|\bar{D}|$ (indeed these results follow from known lower bound for deterministic query
84 processing). Our contribution is to then identify a family of hard queries where deterministic
85 query processing is 'easy' but computing the expected multiplicities is hard.

time

86 **Our upper bound results.** We introduce an $(1 \pm \epsilon)$ -approximation algorithm that computes
87 problem 1.1 in $O_\epsilon(T_{det}^*(Q, \bar{D}))$. In contrast, known approximation techniques ([38, 30]) in
88 set-PDBs need time $\Omega(|C|^{2k})$ (see Appendix G). Further, we generalize the PDB data model
89 considered by the approximation algorithm to a class of bag-Block Independent Disjoint
90 Databases (see Sec. 2.1.1) (BIDBs).

restate by saying our algo works for a more general notion of bag-PDBs beyond c -TIDBs

91 1.1 Polynomial Equivalence

92 A common encoding of probabilistic databases (e.g., in [28, 27, 5, 2] and many others)
93 relies on annotating tuples with lineages, propositional formulas that describe the set of
94 possible worlds that the tuple appears in. The bag semantics analog is a provenance/lineage
95 polynomial $\Phi[Q, \bar{D}, t]$ [25], a polynomial with non-zero integer coefficients and exponents,
96 over integer variables \mathbf{X} encoding input tuple multiplicities.

97 We drop Q, \bar{D} , and t from $\Phi[Q, \bar{D}, t]$ when they are clear from the context or irrelevant to
98 the discussion. We now specify the problem of computing the expectation of tuple multiplicity
99 in the language of lineage polynomials:

100 **Problem 1.3** (Expected Multiplicity of Lineage Polynomials). *Given an \mathcal{RA}^+ query Q ,
101 c -TIDB \mathcal{D} and result tuple t , compute the expected multiplicity of the polynomial $\Phi[Q, \bar{D}, t]$
102 (i.e., $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}} [\Phi[Q, \bar{D}, t](\mathbf{W})]$, where $\mathbf{W} \in \{0, \dots, c\}$).*

103 We note that computing Problem 1.1 is equivalent to computing Problem 1.3 (see Proposition 2.1).

should be \bar{D}

104 1.2 Our Machinery Techniques

105 **Lower Bound Proof Techniques.** All of our results rely on working with a *reduced* form
106 of the lineage polynomial Φ . In fact, it turns out that for the 1- c -TIDB case, computing
107

1?

Ref to Fig 1 here

All the \bar{D} should be \bar{D} here.

move this to line 152

Q: Do we need to notion of 1-BIDB reduced poly?

23:4 Bag PDB Queries

the expected multiplicity (over bag query semantics) is *exactly* the same as evaluating this reduced polynomial over the probabilities that define the TIDB. This is also true when the query input(s) is a block independent disjoint probabilistic database (with tuple multiplicity of at most 1), which we refer to as a 1-BIDB. For our results to be applicable to c -TIDBs, we introduce the following reduction.

Definition 1.4. Any c -TIDB \mathcal{D} , can be reduced to an equivalent 1-BIDB \mathcal{D}' in the following manner. For each $t_i \in \mathcal{D}$, create a block of $c + 1$ disjoint BIDB tuples in \mathcal{D}' such that each tuple in the newly formed block is mapped to its own boolean variable $X_{i,j}$ for $i \in |D|$ and $j \in [c + 1]$. Then, given $\mathbf{W} \in \{0, \dots, c\}^D$, the equivalent world in \mathcal{D}' will set each variable $X_{i,j} = 1$ for each $\mathbf{W}[i] = j$, while (for $\ell \neq j$) all other $X_{i,\ell} \in \mathbf{X}$ of \mathcal{D}' are set to 0.

Example 1.5. Consider the *Route* relation of fig. 2 and query $Q = \pi_{City_1}(\mathbf{Route})$. The output relation Q is $\{\langle \text{Chicago}, X \rangle, \langle \text{Chicago}, Y \rangle\}$ and can be represented as a c -TIDB $Q' = \{\langle \text{Chicago}, X', 2 \rangle\}$, where the following probabilities are true: $Pr[X' = 0] = Pr[\neg X \wedge \neg Y]$, $Pr[X' = 1] = Pr[(X \vee Y) \wedge (\neg X \vee \neg Y)]$, and $Pr[X' = 2] = Pr[X \wedge Y]$. Q' can then be reduced to a 1-BIDB by creating a block of the following disjoint tuples: $Q'' = \{\langle \text{Chicago}, X'_0 \rangle, \langle \text{Chicago}, X'_1 \rangle, \langle \text{Chicago}, X'_2 \rangle\}$ such that $Pr[X'_i = 1] = Pr[X' = i]$.

Next, we motivate this reduced polynomial. Consider the query Q defined as follows over the bag relations of Fig. 2:

```
SELECT 1 FROM OnTime a, Route r, OnTime b
WHERE a.city = r.city1 AND b.city = r.city2
```

It can be verified that $\Phi(A, B, C, E, X, Y, Z)$ for the sole result tuple (i.e. the count) of Q is $AXB + BYE + BZC$. Now consider the product query $Q^2 = Q \times Q$. The lineage polynomial for Q^2 is given by $\Phi^2(A, B, C, E, X, Y, Z) = A^2X^2B^2 + B^2Y^2E^2 + B^2Z^2C^2 + 2AXB^2YE + 2AXB^2ZC + 2B^2YEZC$.

By exploiting linearity of expectation, further pushing expectation through independent variables and observing that for any $W \in \{0, 1\}$, we have $W^2 = W$, the expectation is $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi^2(\mathbf{W})]$ (where W_A is the random variable corresponding to A , distributed by \mathcal{P})

$$\mathbb{E}[W_A] \mathbb{E}[W_X] \mathbb{E}[W_B] + \mathbb{E}[W_B] \mathbb{E}[W_Y] \mathbb{E}[W_E] + \mathbb{E}[W_B] \mathbb{E}[W_Z] \mathbb{E}[W_C] + 2 \mathbb{E}[W_A] \mathbb{E}[W_X] \mathbb{E}[W_B] \mathbb{E}[W_Y] \mathbb{E}[W_E] + 2 \mathbb{E}[W_A] \mathbb{E}[W_Y] \mathbb{E}[W_B] \mathbb{E}[W_Z] \mathbb{E}[W_C] + 2 \mathbb{E}[W_B] \mathbb{E}[W_Y] \mathbb{E}[W_E] \mathbb{E}[W_Z] \mathbb{E}[W_C].$$

This property leads us to consider a structure related to the lineage polynomial.

Definition 1.6. For any polynomial $\Phi(\mathbf{X})$ corresponding to a c -TIDB (henceforth, c -TIDB-lineage polynomial), define the reduced polynomial $\tilde{\Phi}(\mathbf{X})$ to be the polynomial obtained by setting all exponents $e > 1$ in the standard monomial basis (SMB)¹ form of $\Phi(\mathbf{X})$ to 1.

With $\Phi^2(A, B, C, E, X, Y, Z)$ as an example, we have:

$$\tilde{\Phi}^2(A, B, C, E, X, Y, Z) = AXB + BYE + BZC + 2AXB^2YE + 2AXB^2ZC + 2BYE^2ZC.$$

Note that we have argued that for our specific example the expectation that we want is $\tilde{\Phi}^2(Pr(A = 1), Pr(B = 1), Pr(C = 1), Pr(E = 1), Pr(X = 1), Pr(Y = 1), Pr(Z = 1))$. Lemma 1.7 generalizes the equivalence to all $\mathcal{R.A}^+$ queries on TIDBs (proof in Appendix B.5).

¹ This is the representation, typically used in set-PDBs, where the polynomial is represented as sum of 'pure' products. See Definition 2.2 for a formal definition.

Proposal: Expand the example to have general mult maybe have fewer tuples in example. You can define directly Φ by the \mathcal{P} over $X_{i,j}$

This is hard to follow since 1-BIDB has not been defined.

→ defining / using BIDB

can't do without no longer as if Φ is in \mathcal{P} instead $\mathcal{P}(0,1)$

The entire example needs to be worked to handle mult > 1 .

This is incorrect. You replace each X_i by $\sum_{j=1}^c X_{i,j}$

147 ▶ **Lemma 1.7.** Let \mathcal{D} be a 1-BIDB such that the probability distribution \mathcal{P} over $\mathbf{W} \in$
 148 $\{0,1\}^{|\mathcal{D}|}$ (the set of all worlds) is induced by the disjoint condition and the probability
 149 vector $\mathbf{p} = (p_1, \dots, p_{|\mathcal{D}|})$ where $p_i = Pr(W_i = 1)$. For any 1-BIDB-lineage polynomial
 150 $\Phi(\mathbf{X}) = \Phi[Q, \overline{D}, t](\mathbf{X})$, it holds that $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi(\mathbf{W})] = \tilde{\Phi}(\mathbf{p})$.

This lemma should be for c-TIDBs: mention l.b for 1-TIDB

Sec 1.2 starts here

151 To prove our hardness result we show that for the same Q from the example above, for
 152 an arbitrary ‘product width’ k , the query Q^k is able to encode various hard graph-counting
 153 problems (assuming $O(n)$ tuples rather than the $O(1)$ tuples in Fig. 2). We do so by
 154 considering an arbitrary graph G (analogous to the *Route* relation of Q) and analyzing how
 155 the coefficients in the (univariate) polynomial $\tilde{\Phi}(p, \dots, p)$ relate to counts of subgraphs in G
 156 that are isomorphic to various graphs with k edges. E.g., we exploit the fact that the leading
 157 coefficient in Φ corresponding to Q^k is proportional to the number of k -matchings in G , a
 158 known hard problem in parameterized/fine-grained complexity literature.

159 For an upper bound on approximating the expected count, it is easy to check that if all the
 160 probabilities are constant then $\tilde{\Phi}(p_1, \dots, p_n)$ (i.e. evaluating the original lineage polynomial
 161 over the probability values) is a constant factor approximation. For example, using Q^2 from
 162 above, using p_A to denote $Pr[A = 1]$ (and similarly for the other variables), we can see that

This should be under upper bound techniques

163
$$\Phi^2(\mathbf{p}) = p_A^2 p_X^2 p_B^2 + p_B^2 p_Y^2 p_E^2 + p_B^2 p_Z^2 p_C^2 + 2p_A p_X p_B^2 p_Y p_E + 2p_A p_X p_B^2 p_Z p_C + 2p_B^2 p_Y p_E p_Z p_C$$

 164
$$\leq p_A p_X p_B + p_B p_Y p_E + p_B p_Z p_C + 2p_A p_X p_B p_Y p_E + 2p_A p_X p_B p_Z p_C + 2p_B p_Y p_E p_Z p_C = \tilde{\Phi}(\mathbf{p})$$

166 If we assume that all seven probability values are at least $p_0 > 0$, we get that $\Phi^2(\mathbf{p})$ is in
 167 the range $[(p_0)^3 \cdot \tilde{\Phi}(\mathbf{p}), \tilde{\Phi}(\mathbf{p})]$. To get an $(1 \pm \epsilon)$ -multiplicative approximation we uniformly
 168 sample monomials from the SMB representation of Φ and ‘adjust’ their contribution to $\tilde{\Phi}(\cdot)$.

169 **Upper Bound Techniques.** Our negative results (table 1) indicate that c -TIDBs can not
 170 achieve comparable performance to deterministic databases for exact results (under complexity
 171 assumptions). In fact, under plausible hardness conjectures, one cannot (drastically) improve
 172 upon the trivial algorithm to exactly compute the expected multiplicities for c -TIDBs. A
 173 natural followup is whether we can do better if we are willing to settle for an approximation
 174 to the expected multiplicities. In the remainder of this work, we demonstrate that a $(1 \pm \epsilon)$
 175 (multiplicative) approximation with competitive performance is achievable.

even for c=1

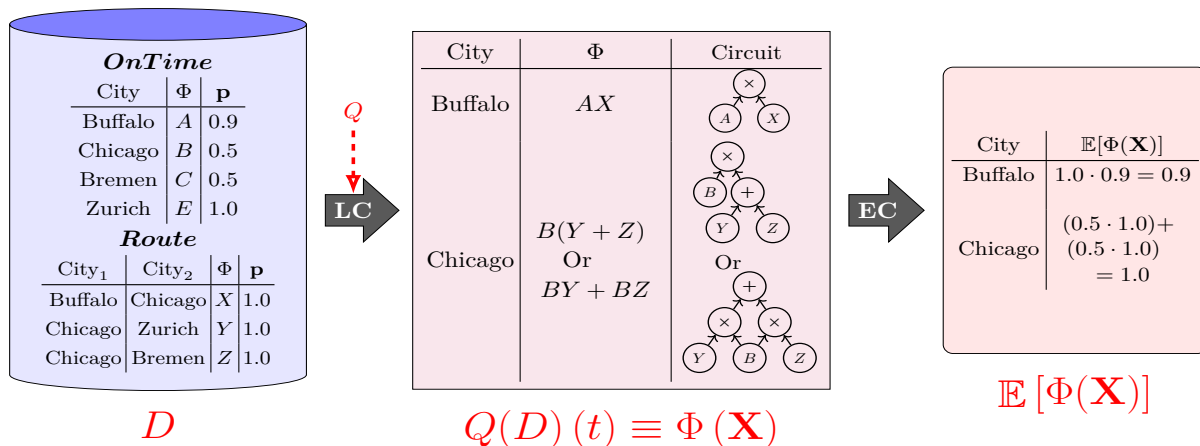


Figure 2 Intensional Query Evaluation Model ($Q = \pi_{\text{City}}(\text{Route} \bowtie_{\text{City}_1 = \text{City}} \text{OnTime})$).

176 We adopt the two-step intensional model of query evaluation used in set-PDBs, as
 177 illustrated in Fig. 2: (i) Lineage Computation (LC): Given input D and Q , output every tuple
 178 t that possibly satisfies Q , annotated with its lineage polynomial ($\Phi(\mathbf{X}) = \Phi[Q, \overline{D}, t](\mathbf{X})$);
 179 (ii) Expectation Computation (EC): Given $\Phi(\mathbf{X})$ for each tuple, compute $\mathbb{E}[\Phi(\mathbf{W})]$. Let
 180 $T_{LC}(Q, D, \mathcal{C})$ denote the runtime of LC when it outputs \mathcal{C} (which is a representation of Φ

D

181 as an arithmetic circuit — more on this representation shortly). Denote by $T_{EC}(\mathcal{C})$ (recall \mathcal{C}
 182 is the output of LC) the runtime of EC, allowing us to formally define our objective:

183 ► **Problem 1.8** (Bag- c -TIDB linear time approximation). Given c -TIDB \mathcal{D} , \mathcal{RA}^+ query
 184 Q , is there a $(1 \pm \epsilon)$ -approximation of $\mathbb{E}_{\mathcal{D} \sim \mathcal{P}_{\bar{\Omega}}} [Q(\mathbf{D})(t)]$ for all result tuples t where $\exists \mathcal{C}$:
 185 $T_{LC}(Q, D, \mathcal{C}) + T_{EC}(\mathcal{C}) \leq O_{\epsilon}(T_{det}^*(Q, D))$?

186 We show in Appendix E.2.1 an $O(T_{det}^*(Q, D))$ algorithm for constructing the lineage
 187 polynomial for all result tuples of an \mathcal{RA}^+ query Q (or more more precisely, a single circuit
 188 \mathcal{C} with one sink per tuple representing the tuple's lineage). A key insight of this paper is
 189 that the representation of \mathcal{C} matters. For example, if we insist that \mathcal{C} represent the lineage
 190 polynomial in SMB, the answer to the above question in general is no, since then we will
 191 need $|\mathcal{C}| \geq \Omega((T_{det}^*(Q, D))^k)$, and hence, just $T_{LC}(Q, D, \mathcal{C})$ will be too large.

192 However, systems can directly emit compact, factorized representations of $\Phi(\mathbf{X})$ (e.g.,
 193 as a consequence of the standard projection push-down optimization [23]). For example,
 194 in Fig. 2, $B(Y + Z)$ is a factorized representation of the SMB-form $BY + BZ$. Accordingly,
 195 this work uses (arithmetic) circuits² as the representation system of $\Phi(\mathbf{X})$.

196 Given that there exists a representation \mathcal{C}^* such that $T_{LC}(Q, D, \mathcal{C}^*) \leq O(T_{det}^*(Q, D))$, we
 197 can now focus on the complexity of EC. We can represent the factorized lineage polynomial
 198 by its corresponding arithmetic circuit \mathcal{C} (whose size we denote by $|\mathcal{C}|$). As we also show in
 199 Appendix E.2.2, this size is also bounded by $T_{det}^*(Q, D)$ (i.e., $|\mathcal{C}^*| \leq O(T_{det}^*(Q, D))$). Thus,
 200 the question of approximation can be reframed as:

201 ► **Problem 1.9** (Problem 1.8 reframed). Given one circuit \mathcal{C} that encodes $\Phi[Q, \bar{D}, t]$ for
 202 all result tuples t (one sink per t) for bag-PDB \mathcal{D} and \mathcal{RA}^+ query Q , does there exist an
 203 algorithm that computes a $(1 \pm \epsilon)$ -approximation of $\mathbb{E}_{\mathcal{D} \sim \mathcal{P}_{\bar{\Omega}}} [Q(\mathbf{D})(t)]$ (for all result tuples t)
 204 in $O(|\mathcal{C}|)$ time?

205 **Old Stuff**

206 A probabilistic database (PDB) \mathcal{D} is a pair $(\bar{\Omega}, \mathcal{P}_{\bar{\Omega}})$, where $\bar{\Omega}$ is a set of deterministic
 207 database instances called possible worlds and $\mathcal{P}_{\bar{\Omega}}$ is a probability distribution over $\bar{\Omega}$. A
 208 tuple independent database (TIDB) (to which we will refer to later) is a PDB such that
 209 each tuple is an independent random event. A commonly studied problem in probabilistic
 210 databases is, given a query Q , PDB \mathcal{D} , and possible query result tuple t , to compute the
 211 tuple's marginal probability of being in the query's result, i.e., computing the expectation
 212 of a Boolean random variable over $\mathcal{P}_{\bar{\Omega}}$ that is 1 for every $D \in \bar{\Omega}$ for which $t \in Q(D)$ and 0
 213 otherwise. In this work, we are interested in bag semantics, where each tuple is associated
 214 with a multiplicity. Following [25], we model bag databases (resp., relations) as functions
 215 from each t to the tuple's multiplicity $D(t) \in \mathbb{N}$ in a possible world D . We refer to such a
 216 probabilistic database as a bag-probabilistic database or bag-PDB for short.

217 The natural generalization of the (set) problem of computing marginal probabilities of
 218 query result tuples to bag semantics is to compute the expectation of a random variable over
 219 $\mathcal{P}_{\bar{\Omega}}$ that is assigned value $Q(D)(t) \in \mathbb{N}$ in world $D \in \bar{\Omega}$, formally $\mathbb{E}_{\mathcal{D} \sim \mathcal{P}_{\bar{\Omega}}} [Q(\mathbf{D})(t)]$.

² An arithmetic circuit is a DAG with variable and/or numeric source nodes and internal, each nodes representing either an addition or multiplication operator.

include ϵ in def?

some more work is needed to come up with this to answer 1.3

Need some discussion on how to handle P1.9

reframed is not as good as P1.9 is not equivalent to P1.8 but it is sufficient

Imp. to clearly state this.