

1 Parameterized and Fine-Grained Analysis of Query 2 Evaluation Over Bag PDBs

3 **Su Feng** ✉

4 Illinois Institute of Technology, Chicago, USA

5 **Boris Glavic** ✉

6 Illinois Institute of Technology, USA

7 **Aaron Huber** ✉

8 University at Buffalo, USA

9 **Oliver Kennedy** ✉

10 University at Buffalo, USA

11 **Atri Rudra** ✉

12 University at Buffalo, USA

13 — Abstract —

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

27 **2012 ACM Subject Classification** Information systems → Incomplete data

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

29 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

30 **1** Introduction

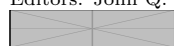
31 This work explores the problem of computing the expectation of a tuple's multiplicity
32 in an important special case of bag TIDB, which we call a c -TIDB. A c -TIDB, $\mathcal{D} =$
33 $(\{0, \dots, c\}^D, \mathcal{P})$ encodes a bag of uncertain tuples such that each tuple in \mathcal{D} has a multiplicity
34 of at most c . D is the set of tuples appearing across all possible worlds, and the set of all
35 worlds is encoded in $\{0, \dots, c\}^D$, which is the set of all vectors of length $n = |D|$ such that
36 each index corresponds to a distinct $t \in D$ storing its multiplicity. \mathcal{P} is a product distribution
37 over the set of all worlds. A given world $\mathbf{M} \in \{0, \dots, c\}^D$ can be interpreted such that, for
38 each $t \in D$, $\mathbf{M}[t]$ is the multiplicity of t in \mathbf{M} . The resulting product distribution can then
39 be encoded as $p_t = Pr[\mathbf{M}[t] = j]$ (for $j \in [c]$), where each t is an independent random event.
40 Allowing for $\leq c$ multiplicities across all tuples gives rise to having $\leq (c + 1)^n$ possible worlds
41 instead of the usual 2^n possible worlds of a 1-TIDB, which (assuming set query semantics),
42 is the same as the traditional set TIDB. In this work, since we are generally considering bag
43 query input, we will only be considering bag query semantics. We denote by $Q(\mathbf{M})(t)$ the



© Aaron Huber, Oliver Kennedy, Atri Rudra, Su Feng, Boris Glavic;
licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:62



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

44 multiplicity of t in query Q over possible world $\mathbf{M} \in \{0, \dots, c\}^D$.

45 We can formally state our problem of computing the expected multiplicity of a result
46 tuple as:

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

49 It is natural to explore computing the expected multiplicity of result tuple as this is the
50 analog for computing the marginal probability of a tuple in a set PDB. In this work we will
51 assume that $c = O(1)$ since this is what typically seen in practice. Allowing for unbounded
52 c is an interesting open problem.

53 **Hardness of Set Query Semantics and Bag Query Semantics.** Set query evaluation
54 semantics over 1-TIDBs have been studied extensively, and the data complexity of the
55 problem in general has been shown by Dalvi and Suicu to be #P-hard [13]. For our setting,
56 there exists a trivial polytime algorithm to compute Problem 1.1 for any query over a c -TIDB
57 due to linearity of expectation by simply computing the expectation over a ‘sum-of-products’
58 representation of the query operations of $Q(\mathcal{D})(t)$. Since we can compute Problem 1.1 in
59 polynomial time, the interesting question that we explore deals with analyzing the hardness
60 of computing expectation using fine-grained analysis and parameterized complexity, where
61 we are interested in the exponent of polynomial runtime.

62 Specifically, in this work we ask if Problem 1.1 can be solved in time linear in the runtime
63 of an equivalent deterministic query. If this is true, then this would open up the way for
64 deployment of c -TIDBs in practice. To analyze this question we denote by $T^*(Q, \mathcal{D})$ the
65 optimal runtime complexity of computing Problem 1.1 over c -TIDB \mathcal{D} .

66 Let $T_{det}(Q, D, c) \equiv Q(D)$ for arbitrary query Q , deterministic database D , and multiplicity
67 bound c . Let $T_{det}^*(Q, D, c) = \min_{Q': Q' \equiv Q} T_{det}(Q', D, c)$ be the optimal runtime (with some
68 caveats; discussed in Sec. 2.3) of query Q on deterministic database D .

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

■ **Table 1** Our lower bounds for a specific hard query Q parameterized by k . For $\mathcal{D} = \{\{0, \dots, c\}^D, \mathcal{P}\}$ those with ‘Multiple’ in the second column need the algorithm to be able to handle multiple \mathcal{P} (for a given 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).

69 **Our lower bound results.** Our question is whether or not it is always true that $T^*(Q, \mathcal{D}) \leq$
70 $T_{det}^*(Q, D, c)$. Unfortunately this is not the case. Table 1 shows our results.

71 Specifically, depending on what hardness result/conjecture we assume, we get various
72 asymptotic versions of *no* as an answer to our question. To make some sense of the other
73 lower bounds in Table 1, we note that it is not too hard to show that $T^*(Q, \mathcal{D}) \leq$
74 $O \left((T_{det}^*(Q, D, c))^k \right)$, where k is the join width (our notion of join width follows from Definition 2.2
75 and Fig. 1.) of the query Q over all result tuples t (and the parameter that defines our family
76 of hard queries).

77 What our lower bound in the third row says is that one cannot get more than a polynomial
78 improvement over essentially the trivial algorithm for Problem 1.1. However, this result

not planned what that is

what is that?

m?

RA T

I am not sure people

multiple probability distributions? a get under p values?

① This should be the query result

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

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

79 assumes a hardness conjecture that is not as well studied as those in the first two rows
80 of the table (see Sec. 3 for more discussion on the hardness assumptions). Further, we
81 note that existing results already imply the claimed lower bounds if we were to replace the
82 $T_{det}^*(Q, D, c)$ by just n (indeed these results follow from known lower bound for deterministic
83 query processing). Our contribution is to then identify a family of hard queries where
84 deterministic query processing is ‘easy’ but computing the expected multiplicities is hard.

85 **Our upper bound results.** We introduce an $(1 \pm \epsilon)$ -approximation algorithm that computes
86 Problem 1.1 in time $O_\epsilon(T_{det}^*(Q, D, c))$. This means, when we are okay with approximation,
87 that we solve Problem 1.1 in time linear in the size of the deterministic query and bag
88 PDBs are deployable in practice. In contrast, known approximation techniques ([38, 30])
89 in set-PDBs need time $\Omega(T_{det}^*(Q, D, c)^{2k})$ (see Appendix G). Further, our approximation
90 algorithm works for a more general notion of bag PDBs beyond c -TIDBs (see Sec. 2.1.1).

Only not absolute

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 (see Fig. 1) $\Phi[Q, D, t]$ [25], a polynomial with non-zero integer coefficients and
96 exponents, over integer variables \mathbf{X} encoding input tuple multiplicities.

97 We drop Q, D , and t from $\Phi[Q, D, t]$ 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 1.2** (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 $\Phi[Q, D, t]$*
102 *(i.e., $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$, where $\mathbf{W} \in \{0, \dots, c\}^D$).*

103 We note that computing Problem 1.1 is equivalent to computing Problem 1.2 (see Proposition 2.5).

104 *↳ input way?*
105 All of our results rely on working with a *reduced* form of the lineage polynomial Φ . In fact,
106 it turns out that for the 1-TIDB case, computing the expected multiplicity (over bag query
107 semantics) is *exactly* the same as evaluating this reduced polynomial over the probabilities
108 that define the 1-TIDB. This is also true when the query input(s) is a block independent
109 disjoint probabilistic database (with tuple multiplicity of at most 1), which we refer to as
110 a 1-BIDB. Next, we motivate this reduced polynomial. Consider the query Q_1 defined as
111 follows over the bag relations of Fig. 2:

112
113
114

```
SELECT 1 FROM T t1, Route r, T t2
WHERE t1.city = r.city1 AND t2.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 $Q_1^2 = Q_1 \times Q_1$. The lineage polynomial for Q_1^2 is given by $\Phi_1^2(A, B, C, E, X, Y, Z)$

$$= A^2X^2B^2 + B^2Y^2E^2 + B^2Z^2C^2 + 2AXB^2YE + 2AXB^2ZC + 2B^2YEZC.$$

116

To compute $\mathbb{E}[\Phi_1^2]$ we can use linearity of expectation and push the expectation through each summand. To keep things simple, let us focus on the monomial $\Phi_1^{(ABX)^2} = A^2X^2B^2$ as the procedure is the same for all other monomials of Φ_1^2 . Let W_X be the random variable corresponding to a lineage variable X . Because the distinct variables in the product are independent, we can push expectation through them yielding $\mathbb{E}[W_A^2W_X^2W_B^2] = \mathbb{E}[W_A^2]\mathbb{E}[W_X^2]\mathbb{E}[W_B^2]$. Since $W_A, W_B \in \{0, 1\}$ we can further derive $\mathbb{E}[W_A]\mathbb{E}[W_X^2]\mathbb{E}[W_B]$ by the fact that for any $W \in \{0, 1\}$, $W^2 = W$. However, we get stuck with $\mathbb{E}[W_X^2]$, since $W_X \in \{0, 1, 2\}$ and for $W_X \leftarrow 2$, $W_X^2 \neq W_X$.

124

Denote the variables of Φ to be $\text{VARS}(\Phi)$. In the c -TIDB setting, $\Phi(\mathbf{X})$ has an equivalent reformulation (Φ_R) that is of use to us. Given $X_t \in \text{VARS}(\Phi)$, by definition $X_t \in \{0, \dots, c\}$. We can replace X_t by $\sum_{j \in [c]} X_{t,j}$ where each $X_{t,j} \in \{0, 1\}$. Then for any $\mathbf{M} \in \{0, \dots, c\}^D$, we set $X_{t,j} = 1$ for $\mathbf{M}_t = j$, while $X_{t,j'} = 0$ for all $j' \neq j \in [c]$. By construction then $\Phi(\mathbf{X}) \equiv \Phi_R(\mathbf{X})$ since for any $X_t \in \text{VARS}(\Phi)$ we have the equality $X_t = j = \sum_{j \in [c]} jX_j$.

128

Considering again our example,

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

$$\begin{aligned} \Phi_{1,R}^{(ABX)^2}(A, X, B) &= \Phi^{(ABX)^2} \left(\sum_{j_1 \in [c]} j_1 A_{j_1}, \sum_{j_2 \in [c]} j_2 X_{j_2}, \sum_{j_3 \in [c]} j_3 B_{j_3} \right) \\ &= \left(\sum_{j_1 \in [c]} j_1 A_{j_1} \right)^2 \left(\sum_{j_2 \in [c]} j_2 X_{j_2} \right)^2 \left(\sum_{j_3 \in [c]} j_3 B_{j_3} \right)^2. \end{aligned}$$

Since the set of multiplicities for tuple t by nature are disjoint we can drop all cross terms and have $\Phi_{1,R}^2 = \sum_{j_1, j_2, j_3 \in [c]} j_1^2 A_{j_1}^2 j_2^2 X_{j_2}^2 j_3^2 B_{j_3}^2$. Computing expectation we get $\mathbb{E}[\Phi_{1,R}^2] = \sum_{j_1, j_2, j_3 \in [c]} j_1^2 j_2^2 j_3^2 \mathbb{E}[W_{A_{j_1}}]\mathbb{E}[W_{X_{j_2}}]\mathbb{E}[W_{B_{j_3}}]$, since we now have that all $W_{X_j} \in \{0, 1\}$. This leads us to consider a structure related to the lineage polynomial.

► **Definition 1.3.** For any polynomial $\Phi((X_t)_{t \in D})$ define the reformulated polynomial $\Phi_R((X_{t,j})_{t \in D, j \in [c]})$ to be the polynomial $\Phi_R = \Phi\left(\left(\sum_{j \in [c]} j \cdot X_{t,j}\right)_{t \in D}\right)$ and ii) define the reduced polynomial $\tilde{\Phi}((X_{t,j})_{t \in D, j \in [c]})$ to be the polynomial resulting from converting Φ_R into the standard monomial basis (SMB),¹ removing all monomials containing the term $X_{t,j}X_{t,j'}$ for $t \in D, j \neq j' \in [c]$, and setting all variable exponents $e > 1$ to 1.

Continuing with the example $\Phi_1^2(A, B, C, E, X_1, X_2, Y, Z)$, to save clutter we i) do not show the full expansion for variables with greatest multiplicity = 1 since e.g. for variable A , the sum of products itself evaluates to $1^2 \cdot A^2 = A$, and ii) for $\sum_{j \in [c]} j^2 \cdot X_j$, we omit the

¹ This is the representation, typically used in set-PDBs, where the polynomial is reresented as sum of 'pure' products. See Definition 2.1 for a formal definition.

146 summands encoding multiplicities > 2 , since the greatest multiplicity of the tuple annotated
 147 with X is 2, likewise those summands will always be evaluated to 0 since the tuple will never
 148 have a multiplicity of > 2 .

$$\begin{aligned}
 & \tilde{\Phi}_1^2(A, B, C, E, X_1, X_2, Y, Z) = \\
 & A \left(\sum_{j \in [c]} j^2 X_j \right) B + BYE + BZC + 2A \left(\sum_{j \in [c]} j^2 X_j \right) BYE + 2A \left(\sum_{j \in [c]} j^2 X_j \right) BZC + 2BYE ZC = \\
 & ABX_1 + AB(2)^2 X_2 + BYE + BZC + 2AX_1 BYE + 2A(2)^2 X_2 BYE + 2AX_1 BZC + 2A(2)^2 X_2 BZC + 2BYE ZC.
 \end{aligned}$$

154 Note that we have argued that for our specific example the expectation that we want is
 155 $\tilde{\Phi}_1^2(\Pr(A=1), \Pr(B=1), \Pr(C=1), \Pr(E=1), \Pr(X_1=1), \Pr(X_2=1), \Pr(Y=1), \Pr(Z=1))$.
 156 Lemma 1.4 generalizes the equivalence to *all* \mathcal{RA}^+ queries on c -TIDBs (proof in Appendix B.5).

157 **► Lemma 1.4.** *For any c -TIDB \mathcal{D} , \mathcal{RA}^+ query Q , and lineage polynomial $\Phi(\mathbf{X}) =$
 158 $\Phi[Q, D, t](\mathbf{X})$, it holds that $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi_R(\mathbf{W})] = \tilde{\Phi}(\mathbf{p})$, where $\mathbf{p} = \left((p_{t,j})_{t \in D, j \in [c]} \right)$.*

159 1.2 Our Techniques

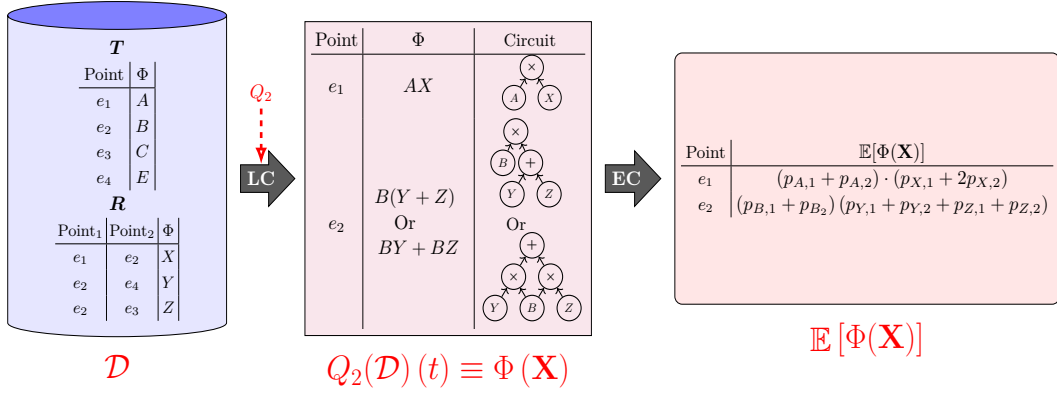
160 **Lower Bound Proof Techniques.** Our main hardness result shows that computing Problem 1.1
 161 is $\#\mathcal{W}[1]$ – hard for 1-TIDB. To prove this result we show that for the same Q_1 from the
 162 example above, for an arbitrary ‘product width’ k , the query Q^k is able to encode various
 163 hard graph-counting problems (assuming $O(n)$ tuples rather than the $O(1)$ tuples in Fig. 2).
 164 We do so by considering an arbitrary graph G (analogous to relation \mathbf{R} of Q) and analyzing
 165 how the coefficients in the (univariate) polynomial $\tilde{\Phi}(p, \dots, p)$ relate to counts of subgraphs
 166 in G that are isomorphic to various graphs with k edges. E.g., we exploit the fact that the
 167 leading coefficient in Φ corresponding to Q^k is proportional to the number of k -matchings in
 168 G , a known hard problem in parameterized/fine-grained complexity literature.

169 **Upper Bound Techniques.** Our negative results (Table 1) indicate that c -TIDBs (even
 170 for $c = 1$) can not achieve comparable performance to deterministic databases for exact
 171 results (under complexity assumptions). In fact, under plausible hardness conjectures, one
 172 cannot (drastically) improve upon the trivial algorithm to exactly compute the expected
 173 multiplicities for 1-TIDBs. A natural followup is whether we can do better if we are willing
 174 to settle for an approximation to the expected multiplicities.

175 We adopt the two-step intensional model of query evaluation used in set-PDBs, as
 176 illustrated in Fig. 2: (i) Lineage Computation (LC): Given input D and Q , output every tuple
 177 t that possibly satisfies Q , annotated with its lineage polynomial ($\Phi(\mathbf{X}) = \Phi[Q, D, t](\mathbf{X})$);
 178 (ii) Expectation Computation (EC): Given $\Phi(\mathbf{X})$ for each tuple, compute $\mathbb{E}[\Phi(\mathbf{W})]$. Let
 179 $T_{LC}(Q, D, \mathcal{C})$ denote the runtime of LC when it outputs \mathcal{C} (which is a representation of Φ as
 180 an arithmetic circuit — more on this representation shortly). Denote by $T_{EC}(\mathcal{C}, \epsilon)$ (recall \mathcal{C}
 181 is the output of LC) the runtime of EC, which we can leverage Definition 1.3 and Lemma 1.4
 182 to address the next formal objective:

183 **► Problem 1.5** (c -TIDB linear time approximation). *Given c -TIDB \mathcal{D} , \mathcal{RA}^+ query Q ,
 184 is there a $(1 \pm \epsilon)$ -approximation of $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[Q(\mathbf{W})(t)]$ for all result tuples t where $\exists \mathcal{C} :$
 185 $T_{LC}(Q, D, \mathcal{C}) + T_{EC}(\mathcal{C}, \epsilon) \leq O_\epsilon(T_{det}^*(Q, D, c))$?*

186 We show in Appendix E.2.1 an $O(T_{det}^*(Q, D, c))$ algorithm for constructing the lineage
 187 polynomial for all result tuples of an \mathcal{RA}^+ query Q (or more precisely, a single circuit



■ **Figure 2** Intensional Query Evaluation Model ($Q_2 = \pi_{\text{Point}}(T \bowtie_{\text{Point}=\text{Point}_1} R)$ and $c = 2$).

188 \mathbf{C} with one sink per tuple representing the tuple’s lineage). A key insight of this paper is
 189 that the representation of \mathbf{C} matters. For example, if we insist that \mathbf{C} represent the lineage
 190 polynomial in SMB, the answer to the above question in general is no, since then we will
 191 need $|\mathbf{C}| \geq \Omega\left((T_{det}^*(Q, D, c))^k\right)$, and hence, just $T_{LC}(Q, D, \mathbf{C})$ 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 \mathbf{C}^* such that $T_{LC}(Q, D, \mathbf{C}^*) \leq O(T_{det}^*(Q, D, c))$,
 197 we can now focus on the complexity of EC. We can represent the factorized lineage polynomial
 198 by its corresponding arithmetic circuit \mathbf{C} (whose size we denote by $|\mathbf{C}|$). As we also show
 199 in Appendix E.2.2, this size is also bounded by $T_{det}^*(Q, D, c)$ (i.e., $|\mathbf{C}^*| \leq O(T_{det}^*(Q, D, c))$).
 200 Thus, the question of approximation can be stated as the following stronger (since Problem 1.5
 201 has access to *all* equivalent \mathbf{C} representing $Q(\mathbf{W})(t)$), but sufficient condition:

202 ► **Problem 1.6.** *Given one circuit \mathbf{C} that encodes $\Phi[Q, D, t]$ for all result tuples t (one sink
 203 per t) for bag-PDB \mathcal{D} and \mathcal{RA}^+ query Q , does there exist an algorithm that computes a
 204 $(1 \pm \epsilon)$ -approximation of $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[Q(\mathbf{W})(t)]$ (for all result tuples t) in $O(|\mathbf{C}|)$ time?*

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

$$209 \quad \Phi_1^2(\mathbf{p}) = p_A^2 p_X^2 p_B^2 + p_B^2 p_Y^2 p_E^2 + p_B^2 p_Z^2 p_C^2 + 2p_A p_X p_B^2 p_Y p_E + 2p_A p_X p_B^2 p_Z p_C + 2p_B^2 p_Y p_E p_Z p_C$$

$$210 \quad \leq p_A p_X p_B + p_B p_Y p_E + p_B p_Z p_C + 2p_A p_X p_B p_Y p_E + 2p_A p_X p_B p_Z p_C + 2p_B p_Y p_E p_Z p_C = \tilde{\Phi}_1^2(\mathbf{p})$$

212 If we assume that all seven probability values are at least $p_0 > 0$, we get that $\Phi_1^2(\mathbf{p})$ is in
 213 the range $[(p_0)^3 \cdot \tilde{\Phi}_1^2(\mathbf{p}), \tilde{\Phi}_1^2(\mathbf{p})]$. In sec. 4 we demonstrate that a $(1 \pm \epsilon)$ (multiplicative)
 214 approximation with competitive performance is achievable. To get an $(1 \pm \epsilon)$ -multiplicative
 215 approximation and solve Problem 1.6, using \mathbf{C} we uniformly sample monomials from the

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.

216 equivalent SMB representation of Φ (without materializing the SMB representation) and
 217 ‘adjust’ their contribution to $\tilde{\Phi}(\cdot)$.

218

219 **Applications.** Recent work in heuristic data cleaning [49, 43, 40, 8, 43] emits a PDB when
 220 insufficient data exists to select the ‘correct’ data repair. Probabilistic data cleaning is a
 221 crucial innovation, as the alternative is to arbitrarily select one repair and ‘hope’ that queries
 222 receive meaningful results. Although PDB queries instead convey the trustworthiness of
 223 results [35], they are impractically slow [18, 17], even in approximation (see Appendix G).
 224 Bags, as we consider, are sufficient for production use, where bag-relational algebra is already
 225 the default for performance reasons. Our results show that bag-PDBs can be competitive,
 226 laying the groundwork for probabilistic functionality in production database engines.

227 **Paper Organization.** We present relevant background and notation in Sec. 2. We then
 228 prove our main hardness results in Sec. 3 and present our approximation algorithm in Sec. 4.
 229 Finally, we discuss related work in Sec. 5 and conclude in Sec. 6. All proofs are in the
 230 appendix.

2 Background and Notation

231

2.1 Polynomial Definition and Terminology

232

233 A polynomial over $\mathbf{X} = (X_1, \dots, X_n)$ with individual degree $B < \infty$ is formally defined as
 234 (where $c_{\mathbf{d}} \in \mathbb{N}$):

$$235 \quad \Phi(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \{0, \dots, B\}^D} c_{\mathbf{d}} \cdot \prod_{t \in D} X_t^{d_t}. \quad (1)$$

236 **Definition 2.1** (Standard Monomial Basis). *The term $\prod_{t \in D} X_t^{d_t}$ in Eq. (1) is a monomial.*
 237 *A polynomial $\Phi(\mathbf{X})$ is in standard monomial basis (SMB) when we keep only the terms with*
 238 *$c_{\mathbf{d}} \neq 0$ from Eq. (1).*

239 Unless otherwise noted, we consider all polynomials to be in SMB representation. When it is
 240 unclear, we use $\text{SMB}(\Phi)$ to denote the SMB form of a polynomial Φ .

241 **Definition 2.2** (Degree). *The degree of polynomial $\Phi(\mathbf{X})$ is the largest $\|\mathbf{d}\|_1$ such that*
 242 *$c_{(d_1, \dots, d_n)} \neq 0$.*

243 As an example, the degree of the polynomial $X^2 + 2XY^2 + Y^2$ is 3. Product terms in lineage
 244 arise only from join operations (Fig. 1), so intuitively, the degree of a lineage polynomial
 245 is analogous to the largest number of joins needed to produce a result tuple. We call a
 246 polynomial $\Phi(\mathbf{X})$ a *c-TIDB-lineage polynomial* (or simply lineage polynomial), if there exists
 247 a \mathcal{RA}^+ query Q , *c*-TIDB \mathcal{D} , and result tuple t such that $\Phi(\mathbf{X}) = \Phi[Q, \mathcal{D}, t](\mathbf{X})$.

2.1.1 c-TIDBs and 1-BIDBs

248

249 An *incomplete database* Ω is a set of deterministic databases ω called possible worlds.

250 A *c*-TIDB \mathcal{D} is a pair $(\{0, \dots, c\}^D, \mathcal{P})$ such that $\{0, \dots, c\}^D$ is an incomplete database
 251 whose set of possible worlds is the $c + 1^n$ tuple/multiplicity combinations across all $t \in D$,
 252 where $|D| = n$, $D = \bigcup_{\mathbf{M} \in \{0, \dots, c\}^D, \mathbf{M}_t \geq 1} t$ is the set of possible tuples across possible worlds,
 253 and \mathcal{P} is a probability distribution over $\{0, \dots, c\}^D$.

254 A block independent database (BIDB) is a related probabilistic data model $\mathcal{D} = (\Omega, \mathcal{P})$
 255 such that the base set of tuples $D = \bigcup_{\omega \in \Omega, t \in \omega} t$ is partitioned into a set of n independent

256 blocks $\{(b_t)_{t \in [n]}\}$ such that the set of tuples $\{(t_j)_{j \in [|b_t|]}\}$ in block b_t are disjoint from one
 257 another. This construction produces the set of possible worlds Ω that consists of all unique
 258 combinations of tuples in D with the constraint that for any $\omega \in \Omega$, no two tuples $t_j, t_{j'}, j \neq j'$
 259 from the same block b_t exist together. A c -BIDB has the further requirement that each block
 260 has a multiplicity of at most c . We present a reduction that is useful in producing our results:

261 ► **Definition 2.3** (*c-TIDB reduction*). Given c -TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, let $\mathcal{D}' =$
 262 (Ω, \mathcal{P}') be the 1-BIDB obtained in the following manner: for each $t \in D$, create block
 263 $b_t = \{(t, jX_{t,j})_{j \in [c]}\}$, such that $X_{t,j} \in \{0, 1\}$. The probability distribution \mathcal{P}' is the one
 264 induced by $\mathbf{p} = ((p_{t,j})_{t \in D, j \in [c]})$ and the BIDB disjoint requirement.

265 For the c -TIDB \mathcal{D} , each $X_t \in [c]$, while in the reduced 1-BIDB \mathcal{D}' , each $X_{t,j} \in \{0, 1\}$.
 266 Hence, in the setting of 1-BIDB, the base case of Fig. 1 now becomes $\Phi[R, D, t] = \sum_{j \in [c]} jX_{t,j}$.
 267 Then given the disjoint requirement and the semantics for constructing the lineage polynomial
 268 over a 1-BIDB, $\Phi[R, D', t]$ is of the same structure as the reformulated polynomial Φ_R of
 269 step i) from Definition 1.3, which then implies that $\tilde{\Phi}$ is the reduced polynomial that results
 270 from step ii) of Definition 1.3, and further that Lemma 1.4 immediately follows for 1-BIDB
 271 polynomials: $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}'}[\Phi(\mathbf{W})] = \tilde{\Phi}(\mathbf{p})$.

272 **Aaron says:** @atri, not sure if \mathcal{P}' should be \mathcal{P}'' (in the above expectation) as discussed
 below. Since $\mathcal{P}' \equiv \mathcal{P}''$, then the proof still holds for Lemma 1.4, but maybe it is
 important to \mathcal{P}'' to drive the point home that we iterate over the all worlds set (as
 opposed to the set of possible worlds) when computing the expectation of a polynomial.
 Or maybe it suffices to note that $\mathcal{P}' \equiv \mathcal{P}''$.

273 Instead of looking only at the possible worlds of \mathcal{D} , one can consider all worlds, including
 274 those that cannot exist due to disjointness. The all worlds set can be modeled by $\mathbf{M} \in$
 275 $\{0, 1\}^{cn}$,³ such that $\mathbf{M}_{t,j} \in \mathbf{M}$ represents whether or not the multiplicity of t is j . We denote
 276 a probability distribution over all $\mathbf{M} \in \{0, 1\}^n$ as \mathcal{P}'' . When \mathcal{P}'' is the one induced from
 277 each $p_{t,j}$ while assigning $Pr[\mathbf{M}] = 0$ for any \mathbf{M} with $\mathbf{M}_{t,j} = \mathbf{M}_{t,j'} = 1$ for $j \neq j'$, we end up
 278 with a bijective mapping from \mathcal{P}' to \mathcal{P}'' , such that each mapping is equivalent, implying the
 279 distributions are equivalent. Appendix B.2 has more details.

280 Let $|\Phi|$ be the number of operators in Φ .

281 ► **Corollary 2.4.** If Φ is a BIDB-lineage polynomial already in SMB, then the expectation of
 282 Φ , i.e., $\mathbb{E}[\Phi] = \tilde{\Phi}(p_1, \dots, p_n)$ can be computed in $O(|\Phi|)$ time.

283 Queries over probabilistic databases are evaluated using the so-called possible world semantics.
 284 Under the possible world semantics, the result of a query Q over an incomplete database
 285 Ω is the set of query answers produced by evaluating Q over each possible world $\omega \in \Omega$:
 286 $\{Q(\omega) : \omega \in \Omega\}$.

287 The result of a query is the pair $(Q(\omega), \mathcal{P}')$ where \mathcal{P}' is a probability distribution that
 288 assigns to each possible query result the sum of the probabilities of the worlds that produce
 289 this answer: $Pr[\omega \in \Omega] = \sum_{\omega' \in \Omega, Q(\omega') = Q(\omega)} Pr[\omega']$.

290 Recalling Fig. 1 again, which defines the lineage polynomial $\Phi[Q, D, t]$ for any \mathcal{RA}^+
 291 query. We now make a meaningful connection between possible world semantics and world
 292 assignments on the lineage polynomial.

³ Here and later, especially in Sec. 4, we will rename the variables as X_1, \dots, X_n , where $n = \sum_{i=1}^{\ell} |b_i|$.

293 ▶ **Proposition 2.5** (Expectation of polynomials). *Given a bag-PDB $\mathcal{D} = (\Omega, \mathcal{P})$, \mathcal{RA}^+ query*
 294 *Q , and lineage polynomial $\Phi[Q, D, t]$ for arbitrary result tuple t , we have (denoting \mathbf{D} as the*
 295 *random variable over Ω): $\mathbb{E}_{\mathbf{D} \sim \mathcal{P}}[Q(\mathbf{D})(t)] = \mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$.*

296 A formal proof of Proposition 2.5 is given in Appendix B.3.⁴ We focus on the problem of
 297 computing $\mathbb{E}_{\mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$ from now on, assume implicit Q, D, t , and drop them from
 298 $\Phi[Q, D, t]$ (i.e., $\Phi(\mathbf{X})$ will denote a polynomial).

299 2.2 Formalizing Problem 1.6

300 We represent lineage polynomials via *arithmetic circuits* [9], a standard way to represent
 301 polynomials over fields (particularly in the field of algebraic complexity) that we use for
 302 polynomials over \mathbb{N} in the obvious way. Since we are particularly using circuits to model
 303 lineage polynomials, we can refer to these circuits as lineage circuits. However, when the
 304 meaning is clear, we will drop the term lineage and only refer to them as circuits.

305 ▶ **Definition 2.6** (Circuit). *A circuit \mathcal{C} is a Directed Acyclic Graph (DAG) whose source*
 306 *gates (in degree of 0) consist of elements in either \mathbb{N} or \mathbf{X} . For each result tuple there exists*
 307 *one sink gate. The internal gates have binary input and are either sum (+) or product (\times)*
 308 *gates. Each gate has the following members: **type**, **partial**, **input**, **degree**, **Lweight**, and*
 309 ***Rweight**, where **type** is the value type $\{+, \times, \text{VAR}, \text{NUM}\}$ and **input** the list of inputs. Source*
 310 *gates have an extra member **val** storing the value. \mathcal{C}_L (\mathcal{C}_R) denotes the left (right) input of \mathcal{C} .*

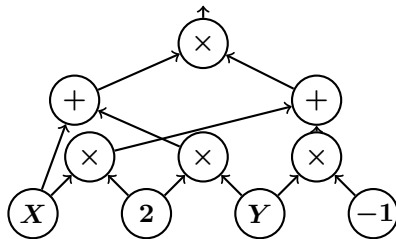
311 **Aaron says:** Does the following matter, i.e., does it point anything out special for our
 312 research?

312 When the underlying DAG is a tree (with edges pointing towards the root), the structure
 313 is an expression tree T . In such a case, the root of T is analogous to the sink of \mathcal{C} . The fields
 314 **partial**, **degree**, **Lweight**, and **Rweight** are used in the proofs of Appendix D.

315 The circuits in Fig. 2 encode their respective polynomials in column Φ . Note that each
 316 circuit \mathcal{C} encodes a tree, with edges pointing towards the root.

317 We next formally define the relationship of
 318 circuits with polynomials. While the definition
 319 assumes one sink for notational convenience, it
 320 easily generalizes to the multiple sinks case.

321 ▶ **Definition 2.7** ($\text{POLY}(\cdot)$). *Denote $\text{POLY}(\mathcal{C})$ to*
 322 *be the function from the sink of circuit \mathcal{C} to*
 323 *its corresponding polynomial (in SMB). $\text{POLY}(\cdot)$*
 324 *is recursively defined on \mathcal{C} as follows, with*
 325 *addition and multiplication following the standard*
 326 *interpretation for polynomials:*



325 ■ **Figure 3** Circuit encoding of $(X + 2Y)(2X - Y)$

$$327 \text{POLY}(\mathcal{C}) = \begin{cases} \text{POLY}(\mathcal{C}_L) + \text{POLY}(\mathcal{C}_R) & \text{if } \mathcal{C}.\text{type} = + \\ \text{POLY}(\mathcal{C}_L) \cdot \text{POLY}(\mathcal{C}_R) & \text{if } \mathcal{C}.\text{type} = \times \\ \mathcal{C}.\text{val} & \text{if } \mathcal{C}.\text{type} = \text{VAR OR NUM}. \end{cases}$$

⁴ Although Proposition 2.5 follows, e.g., as an obvious consequence of [28]’s Theorem 7.1, we are unaware of any formal proof for bag-probabilistic databases.

328 \mathcal{C} need not encode $\Phi(\mathbf{X})$ in the same, default SMB representation. For instance, \mathcal{C} could
 329 encode the factorized representation $(X + 2Y)(2X - Y)$ of $\Phi(\mathbf{X}) = 2X^2 + 3XY - 2Y^2$, as
 330 shown in Fig. 3, while $\text{POLY}(\mathcal{C}) = \Phi(\mathbf{X})$ is always the equivalent SMB representation.

331 **► Definition 2.8** (Circuit Set). $\text{CSet}(\Phi(\mathbf{X}))$ is the set of all possible circuits \mathcal{C} such that
 332 $\text{POLY}(\mathcal{C}) = \Phi(\mathbf{X})$.

333 The circuit of Fig. 3 is an element of $\text{CSet}(2X^2 + 3XY - 2Y^2)$. One can think of
 334 $\text{CSet}(\Phi(\mathbf{X}))$ as the infinite set of circuits where for each element \mathcal{C} , $\text{POLY}(\mathcal{C}) = \Phi(\mathbf{X})$.

335 We are now ready to formally state the final version of Problem 1.6.

336 **► Definition 2.9** (The Expected Result Multiplicity Problem). Let \mathcal{D} be an arbitrary BIDB-
 337 PDB and \mathbf{X} be the set of variables annotating tuples in D_Ω . Fix an \mathcal{RA}^+ query Q and a
 338 result tuple t . The EXPECTED RESULT MULTIPLICITY PROBLEM is defined as follows:

340 **Input:** $\mathcal{C} \in \text{CSet}(\Phi(\mathbf{X}))$ for $\Phi(\mathbf{X}) = \Phi[Q, D, t]$ **Output:** $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$

341 2.3 Relationship to Deterministic Query Runtimes

342 To decouple our results from specific join algorithms, we first abstract the cost of a join.

343 **► Definition 2.10** (Join Cost). Denote by $T_{\text{join}}(R_1, \dots, R_m)$ the runtime of an algorithm
 344 for computing the m -ary join $R_1 \bowtie \dots \bowtie R_m$. We require only that the algorithm must
 345 enumerate its output, i.e., that $T_{\text{join}}(R_1, \dots, R_m) \geq |R_1 \bowtie \dots \bowtie R_m|$.

346 Worst-case optimal join algorithms [37, 36] and query evaluation via factorized databases [39]
 347 (as well as work on FAQs [33]) can be modeled as \mathcal{RA}^+ queries (though the query size is
 348 data dependent). For these algorithms, $T_{\text{join}}(R_1, \dots, R_m)$ is linear in the AGM bound [6].
 349 Our cost model for general query evaluation follows from the join cost:

$$\begin{aligned}
 350 \quad T_{\text{det}}(R, D) &= |D \cdot R| \quad T_{\text{det}}(\sigma Q, D) = T_{\text{det}}(Q, D) \quad T_{\text{det}}(\pi Q, D) = T_{\text{det}}(Q, D) + |Q(D)| \\
 & \quad T_{\text{det}}(Q \cup Q', D) = T_{\text{det}}(Q, D) + T_{\text{det}}(Q', D) + |Q(D)| + |Q'(D)| \\
 351 \quad T_{\text{det}}(Q_1 \bowtie \dots \bowtie Q_m, D) &= T_{\text{det}}(Q_1, D) + \dots + T_{\text{det}}(Q_m, D) + T_{\text{join}}(Q_1(D), \dots, Q_m(D))
 \end{aligned}$$

352 Under this model, an \mathcal{RA}^+ query Q evaluated over database D has runtime $O(T_{\text{det}}(Q, D))$.
 353 We assume that full table scans are used for every base relation access. We can model index
 354 scans by treating an index scan query $\sigma_\theta(R)$ as a base relation.

355 Finally, Lemma E.2 and Lemma E.3 show that for any \mathcal{RA}^+ query Q and D_Ω , there
 356 exists a circuit \mathcal{C}^* such that $T_{LC}(Q, D_\Omega, \mathcal{C}^*)$ and $|\mathcal{C}^*|$ are both $O(T_{\text{det}}(Q, D_\Omega))$. Recall we
 357 assumed these two bounds when we moved from Problem 1.5 to Problem 1.6.

358 3 Hardness of Exact Computation

359 In this section, we will prove the hardness results claimed in Table 1 for a specific (family) of
 360 hard instance (Q, \mathcal{D}) for Problem 1.2 where \mathcal{D} is a 1-TIDB. Note that this implies hardness
 361 for c -TIDBs ($c \geq 1$), BIDBs and general bag-PDB, showing Problem 1.2 cannot be done in
 362 $O(T_{\text{det}}^*(Q, D))$ runtime.

363 3.1 Preliminaries

364 Our hardness results are based on (exactly) counting the number of (not necessarily induced)
 365 subgraphs in G isomorphic to H . Let $\#(G, H)$ denote this quantity. We can think of H

366 as being of constant size and G as growing. In particular, we will consider the problems of
 367 computing the following counts (given G in its adjacency list representation): $\#(G, \mathfrak{A})$ (the
 368 number of triangles), $\#(G, \mathfrak{A}\mathfrak{A}\mathfrak{A})$ (the number of 3-matchings), and the latter's generalization
 369 $\#(G, \mathfrak{A} \cdots \mathfrak{A}^k)$ (the number of k -matchings). We use $T_{match}(k, G)$ to denote the optimal
 370 runtime of computing $\#(G, \mathfrak{A} \cdots \mathfrak{A}^k)$ exactly. Our hardness results in Sec. 3.2 are based on
 371 the following hardness results/conjectures:

372 ► **Theorem 3.1** ([11]). *Given positive integer k and undirected graph $G = (V, E)$ with
 373 no self-loops or parallel edges, $T_{match}(k, G) \geq \omega(f(k) \cdot |E|^c)$ for any function f and fixed
 374 constant c independent of $|E|$ and k (assuming $\#\mathbb{W}[0] \neq \#\mathbb{W}[1]$).*

375 ► **Conjecture 3.2.** *There exists an absolute constant $c_0 > 0$ such that for every $G = (V, E)$,
 376 we have $T_{match}(k, G) \geq \Omega(|E|^{c_0 \cdot k})$ for large enough k .*

377 We note that the above conjecture is somewhat non-standard. In particular, the best known
 378 algorithm to compute $\#(G, \mathfrak{A} \cdots \mathfrak{A}^k)$ takes time $\Omega(|V|^{k/2})$ (i.e. if this is the best algorithm
 379 then $c_0 = \frac{1}{4}$) [11]. What the above conjecture is saying is that one can only hope for a
 380 polynomial improvement over the state of the art algorithm to compute $\#(G, \mathfrak{A} \cdots \mathfrak{A}^k)$.

381 Our hardness result in Section 3.3 is based on the following conjectured hardness result:

382 ► **Conjecture 3.3.** *There exists a constant $\epsilon_0 > 0$ such that given an undirected graph
 383 $G = (V, E)$, computing $\#(G, \mathfrak{A})$ exactly cannot be done in time $o(|E|^{1+\epsilon_0})$.*

384 The so called *Triangle detection hypothesis* (cf. [34]), which states that detecting the presence
 385 of triangles in G takes time $\Omega(|E|^{4/3})$, implies that in Conjecture 3.3 we can take $\epsilon_0 \geq \frac{1}{3}$.

386 All of our hardness results rely on a simple lineage polynomial encoding of the edges
 387 of a graph. To prove our hardness result, consider a graph $G = (V, E)$, where $|E| = m$,
 388 $V = [n]$. Our lineage polynomial has a variable X_i for every i in $[n]$. Consider the polynomial
 389 $\Phi_G(\mathbf{X}) = \sum_{(i,j) \in E} X_i \cdot X_j$. The hard polynomial for our problem will be a suitable power $k \geq 3$
 390 of the polynomial above:

391 ► **Definition 3.4.** *For any graph $G = (V, E)$ and $k \geq 1$, define*

$$392 \quad \Phi_G^k(X_1, \dots, X_n) = \left(\sum_{(i,j) \in E} X_i \cdot X_j \right)^k.$$

393 Returning to Fig. 2, it is easy to see that $\Phi_G^k(\mathbf{X})$ is the lineage polynomial whose structure
 394 mirrors the query Q_2 from Sec. 1. Let us alias

```
395 SELECT 1 FROM T t1, R r, T t2
396 WHERE t1.city = r.city1 AND t2.city = r.city2
397
398
```

399 as R_i for each $i \in [k]$. The query Q^k then becomes

```
400 SELECT COUNT(*) FROM R1 JOIN R2 JOIN...JOIN Rk
401
402
```

403 Further, the c -TIDB instance of Fig. 2 generalizes to one compatible to Definition 3.4 as
 404 follows. Relation T has n tuples corresponding to each vertex for i in $[n]$, each with probability
 405 p_i and R has tuples corresponding to the edges E (each with probability of 1).⁵ In other

⁵ Technically, $\Phi_G^k(\mathbf{X})$ should have variables corresponding to tuples in R as well, but since they always are present with probability 1, we drop those. Our argument also works when all the tuples in R also are present with probability p but to simplify notation we assign probability 1 to edges.

406 words, for this instance D contains the set of n unary tuples in T (which corresponds to
 407 V) and m binary tuples in R (which corresponds to E). Note that this implies that Φ_G^k is
 408 indeed a c -TIDB-lineage polynomial.

409 **Aaron says:** Can the proofs generalize to 2-TIDB, as the new updated Fig. 2 now is?

410 Next, we note that the runtime for answering Q^k on deterministic database D , as defined
 411 above, is $O(m)$ (i.e. deterministic query processing is ‘easy’ for this query):

412 ► **Lemma 3.5.** *Let Q^k and D be as defined above. Then $T_{det}(Q^k, D)$ is $O(km)$.*

413 **Aaron says:** Should the above be $T_{det}()$ or $T_{det}^*()$?

414 3.2 Multiple Distinct p Values

415 We are now ready to present our main hardness result.

416 **Aaron says:** Note that Definition 1.3 has been changed, where we compute $\tilde{\Phi}$ from Φ_R .

417 ► **Theorem 3.6.** *Let p_0, \dots, p_{2k} be $2k + 1$ distinct values in $(0, 1]$. Then computing
 418 $\tilde{\Phi}_G^k(p_i, \dots, p_i)$ (over all $i \in [2k + 1]$ for arbitrary $G = (V, E)$ needs time $\Omega(T_{match}(k, G))$,
 419 assuming $T_{match}(k, G) \geq \omega(|E|)$.*

420 Note that the second row of Table 1 follows from Proposition 2.5, Theorem 3.6, Lemma 3.5,
 421 and Theorem 3.1 while the third row is proved by Proposition 2.5, Theorem 3.6, Lemma 3.5,
 422 and Conjecture 3.2. Since Conjecture 3.2 is non-standard, the latter hardness result should
 423 be interpreted as follows. Any substantial polynomial improvement for Problem 1.2 (over the
 424 trivial algorithm that converts Φ into SMB and then uses Corollary 2.4 for EC) would lead
 425 to an improvement over the state of the art *upper* bounds on $T_{match}(k, G)$. Finally, note
 426 that Theorem 3.6 needs one to be able to compute the expected multiplicities over $(2k + 1)$
 427 distinct values of p_i , each of which corresponds to distinct \mathcal{P} (for the same D), which explain
 428 the ‘Multiple’ entry in the second column in the second and third row in Table 1. Next, we
 429 argue how to get rid of this latter requirement.

430 3.3 Single p value

431 While Theorem 3.6 shows that computing $\tilde{\Phi}(p, \dots, p)$ for multiple values of p in general is
 432 hard it does not rule out the possibility that one can compute this value exactly for a *fixed*
 433 value of p . Indeed, it is easy to check that one can compute $\tilde{\Phi}(p, \dots, p)$ exactly in linear time
 434 for $p \in \{0, 1\}$. Next we show that these two are the only possibilities:

435 ► **Theorem 3.7.** *Fix $p \in (0, 1)$. Then assuming Conjecture 3.3 is true, any algorithm that
 436 computes $\tilde{\Phi}_G^3(p, \dots, p)$ for arbitrary $G = (V, E)$ exactly has to run in time $\Omega(|E|^{1+\epsilon_0})$, where
 437 ϵ_0 is as defined in Conjecture 3.3.*

438 Note that Proposition 2.5 and Theorem 3.7 above imply the hardness result in the first
 439 row of Table 1. We note that Theorem 3.1 and Conjecture 3.2 (and the lower bounds in the
 440 second and third row of Table 1) need k to be large enough (in particular, we need a family
 441 of hard queries). But the above Theorem 3.7 (and the lower bound in first row of Table 1)
 442 holds for $k = 3$ (and hence for a fixed query).

4 $1 \pm \epsilon$ Approximation Algorithm

Aaron says: If we get rid of the problem statements, then we need to remember to get rid of references to particular problem statements, as in below.

In Sec. 3, we showed that the answer to Problem 1.6 is no. With this result, we now design an approximation algorithm for our problem that runs in $O(|\mathcal{C}|)$ for a very broad class of circuits (see the discussion after Lemma 4.8 for more). The following approximation algorithm applies to BIDD lineage polynomials (over \mathcal{RA}^+ queries), though our bounds are more meaningful for a non-trivial subclass of queries over BIDDs that contains all queries on TIDDs, as well as the queries of the PDBench benchmark [1]. All proofs and pseudocode can be found in Appendix D.

Aaron says: We are going to have to rework γ in this section, as well as the proof for our result.

4.1 Preliminaries and some more notation

We now introduce definitions and notation related to circuits and polynomials that we will need to state our upper bound results. First we introduce the expansion $E(\mathcal{C})$ of circuit \mathcal{C} which is used in our algorithm for sampling monomials (part of our approximation algorithm).

► **Definition 4.1** ($E(\mathcal{C})$). For a circuit \mathcal{C} , we define $E(\mathcal{C})$ as a list of tuples (\mathbf{v}, c) , where \mathbf{v} is a set of variables and $c \in \mathbb{N}$. $E(\mathcal{C})$ has the following recursive definition (\circ is list concatenation).

$$E(\mathcal{C}) = \begin{cases} E(\mathcal{C}_L) \circ E(\mathcal{C}_R) & \text{if } \mathcal{C}.type = + \\ \{(\mathbf{v}_L \cup \mathbf{v}_R, c_L \cdot c_R) \mid (\mathbf{v}_L, c_L) \in E(\mathcal{C}_L), (\mathbf{v}_R, c_R) \in E(\mathcal{C}_R)\} & \text{if } \mathcal{C}.type = \times \\ List[(\emptyset, \mathcal{C}.val)] & \text{if } \mathcal{C}.type = NUM \\ List[(\{\mathcal{C}.val\}, 1)] & \text{if } \mathcal{C}.type = VAR. \end{cases}$$

Later on, we will denote the monomial composed of the variables in \mathbf{v} as \mathbf{v}_m . As an example of $E(\mathcal{C})$, consider \mathcal{C} illustrated in Fig. 3. $E(\mathcal{C})$ is then $[(X, 2), (XY, -1), (XY, 4), (Y, -2)]$. This helps us redefine $\tilde{\Phi}$ (see Eq. (2)) in a way that makes our algorithm more transparent.

► **Definition 4.2** ($|\mathcal{C}|$). For any circuit \mathcal{C} , the corresponding positive circuit, denoted $|\mathcal{C}|$, is obtained from \mathcal{C} as follows. For each leaf node ℓ of \mathcal{C} where $\ell.type$ is NUM , update $\ell.value$ to $|\ell.value|$.

We will overload notation and use $|\mathcal{C}|(\mathbf{X})$ to mean $POLY(|\mathcal{C}|)$. Conveniently, $|\mathcal{C}|(1, \dots, 1)$ gives us $\sum_{(\mathbf{v}, c) \in E(\mathcal{C})} |c|$.

► **Definition 4.3** ($SIZE(\cdot)$, $DEPTH(\cdot)$). The functions $SIZE$ and $DEPTH$ output the number of gates and levels respectively for input \mathcal{C} .

► **Definition 4.4** ($DEG(\cdot)$).⁶ $DEG(\mathcal{C})$ is defined recursively as follows:

$$DEG(\mathcal{C}) = \begin{cases} \max(DEG(\mathcal{C}_L), DEG(\mathcal{C}_R)) & \text{if } \mathcal{C}.type = + \\ DEG(\mathcal{C}_L) + DEG(\mathcal{C}_R) + 1 & \text{if } \mathcal{C}.type = \times \\ 1 & \text{if } \mathcal{C}.type = VAR \\ 0 & \text{otherwise.} \end{cases}$$

⁶ Note that the degree of $POLY(|\mathcal{C}|)$ is always upper bounded by $DEG(\mathcal{C})$ and the latter can be strictly larger (e.g. consider the case when \mathcal{C} multiplies two copies of the constant 1— here we have $deg(\mathcal{C}) = 1$ but degree of $POLY(|\mathcal{C}|)$ is 0).

23:14 Bag PDB Queries

472 Next, we use the following notation for the complexity of multiplying integers:

473 ► **Definition 4.5** ($\overline{\mathcal{M}}(\cdot, \cdot)$). ⁷ In a RAM model of word size of W -bits, $\overline{\mathcal{M}}(M, W)$ denotes
 474 the complexity of multiplying two integers represented with M -bits. (We will assume that for
 475 input of size N , $W = O(\log N)$.)

476 Finally, to get linear runtime results, we will need to define another parameter modeling
 477 the (weighted) number of monomials in $\mathbf{E}(\mathcal{C})$ that need to be ‘canceled’ when monomials
 478 with dependent variables are removed (??). Let $\text{ISIND}(\cdot)$ be a boolean function returning
 479 true if monomial \mathbf{v}_m is composed of independent variables and false otherwise; further, let $\mathbb{1}_\theta$
 480 also be a boolean function returning true if θ evaluates to true.

481 ► **Definition 4.6** (Parameter γ). Given a BIDB circuit \mathcal{C} define

Aaron says: Technically, \mathbf{v} is a set of variables rather than a monomial. Perhaps we don’t need the $\text{VAR}(\cdot)$ function and can replace it with a function that returns the monomial represented by a set of variables. **FIXED:** need to propagate this to the appendix (\mathbf{v}_m)

Aaron says: To add, this is an issue on line 1073, 1117 of app C.

484
$$\gamma(\mathcal{C}) = \frac{\sum_{(v,c) \in \mathbf{E}(\mathcal{C})} |c| \cdot \mathbb{1}_{\neg \text{ISIND}(\mathbf{v}_m)}}{|\mathcal{C}|(1, \dots, 1)}.$$

485 4.2 Our main result

486 **Algorithm Idea.** Our approximation algorithm (APPROXIMATE $\tilde{\Phi}$ pseudo code in Appendix D.1)
 487 is based on the following observation. Given a lineage polynomial $\Phi(\mathbf{X}) = \text{POLY}(\mathcal{C})$ for circuit
 488 \mathcal{C} over BIDB , we have:

489
$$\tilde{\Phi}(p_1, \dots, p_n) = \sum_{(v,c) \in \mathbf{E}(\mathcal{C})} \mathbb{1}_{\text{ISIND}(\mathbf{v}_m)} \cdot |c| \cdot \prod_{X_i \in \mathbf{v}} p_i. \quad (2)$$

490 Given the above, the algorithm is a sampling based algorithm for the above sum: we
 491 sample (via SAMPLEMONOMIAL) $(\mathbf{v}, \mathbf{c}) \in \mathbf{E}(\mathcal{C})$ with probability proportional to $|c|$ and
 492 compute $Y = \mathbb{1}_{\text{ISIND}(\mathbf{v}_m)} \cdot \prod_{X_i \in \mathbf{v}} p_i$. Repeating the sampling appropriate number of times and
 493 computing the average of Y gives us our final estimate. ONEPASS is used to compute the
 494 sampling probabilities needed in SAMPLEMONOMIAL (details are in Appendix D).

495 **Runtime analysis.** We can argue the following runtime for the algorithm outlined above:

496 ► **Theorem 4.7.** Let \mathcal{C} be an arbitrary BIDB circuit and define $\Phi(\mathbf{X}) = \text{POLY}(\mathcal{C})$ and let
 497 $k = \text{DEG}(\mathcal{C})$. Let $\gamma = \gamma(\mathcal{C})$. Further let it be the case that $p_i \geq p_0$ for all $i \in [n]$. Then an
 498 estimate \mathcal{E} of $\tilde{\Phi}(p_1, \dots, p_n)$ satisfying

499
$$\Pr \left(\left| \mathcal{E} - \tilde{\Phi}(p_1, \dots, p_n) \right| > \epsilon' \cdot \tilde{\Phi}(p_1, \dots, p_n) \right) \leq \delta \quad (3)$$

⁷ We note that when doing arithmetic operations on the RAM model for input of size N , we have that $\overline{\mathcal{M}}(O(\log N), O(\log N)) = O(1)$. More generally we have $\overline{\mathcal{M}}(N, O(\log N)) = O(N \log N \log \log N)$.

500 can be computed in time

$$501 \quad O\left(\left(\text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta} \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C})}{(\epsilon')^2 \cdot (1-\gamma)^2 \cdot p_0^{2k}}\right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right). \quad (4)$$

502 In particular, if $p_0 > 0$ and $\gamma < 1$ are absolute constants then the above runtime simplifies to
 503 $O_k\left(\left(\frac{1}{(\epsilon')^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta}\right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right).$

504 The restriction on γ is satisfied by any TIDB (where $\gamma = 0$) as well as for all three queries
 505 of the PDBench BIDB benchmark (see Appendix D.10 for experimental results).

506 We briefly connect the runtime in Eq. (4) to the algorithm outline earlier (where we
 507 ignore the dependence on $\overline{\mathcal{M}}(\cdot, \cdot)$, which is needed to handle the cost of arithmetic operations
 508 over integers). The $\text{SIZE}(\mathcal{C})$ comes from the time take to run ONEPASS once (ONEPASS
 509 essentially computes $|\mathcal{C}|(1, \dots, 1)$ using the natural circuit evaluation algorithm on \mathcal{C}). We
 510 make $\frac{\log \frac{1}{\delta}}{(\epsilon')^2 \cdot (1-\gamma)^2 \cdot p_0^{2k}}$ many calls to SAMPLEMONOMIAL (each of which essentially traces $O(k)$
 511 random sink to source paths in \mathcal{C} all of which by definition have length at most $\text{DEPTH}(\mathcal{C})$).

512 Finally, we address the $\overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))$ term in the runtime.

513 **► Lemma 4.8.** *For any BIDB circuit \mathcal{C} with $\text{DEG}(\mathcal{C}) = k$, we have $|\mathcal{C}|(1, \dots, 1) \leq 2^{2^k \cdot \text{DEPTH}(\mathcal{C})}$.
 514 Further, if \mathcal{C} is a tree, then we have $|\mathcal{C}|(1, \dots, 1) \leq \text{SIZE}(\mathcal{C})^{O(k)}$.*

515 Note that the above implies that with the assumption $p_0 > 0$ and $\gamma < 1$ are absolute
 516 constants from Theorem 4.7, then the runtime there simplifies to $O_k\left(\frac{1}{(\epsilon')^2} \cdot \text{SIZE}(\mathcal{C})^2 \cdot \log \frac{1}{\delta}\right)$
 517 for general circuits \mathcal{C} . If \mathcal{C} is a tree, then the runtime simplifies to $O_k\left(\frac{1}{(\epsilon')^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta}\right)$,
 518 which then answers Problem 1.6 in yes for such circuits.

519 Finally, note that by Proposition E.1 and Lemma E.2 for any \mathcal{RA}^+ query Q , there exists a
 520 circuit \mathcal{C}^* for $\Phi[Q, D, t]$ such that $\text{DEPTH}(\mathcal{C}^*) \leq O_{|Q|}(\log n)$ and $\text{SIZE}(\mathcal{C}) \leq O_k(T_{det}^*(Q, D_\Omega))$.
 521 Using this along with Lemma 4.8, Theorem 4.7 and the fact that $n \leq T_{det}^*(Q, D_\Omega)$, we answer
 522 Problem 1.5 in the affirmative as follows:

523 **► Corollary 4.9.** *Let Q be an \mathcal{RA}^+ query and \mathcal{D} be an BIDB with $p_0 > 0$ and $\gamma < 1$ (where
 524 p_0, γ as in Theorem 4.7) are absolute constants. Let $\Phi(\mathbf{X}) = \Phi[Q, D, t]$ for any result tuple
 525 t with $\text{deg}(\Phi) = k$. Then one can compute an approximation satisfying Eq. (3) in time
 526 $O_{k, |Q|, \epsilon', \delta}(T_{det}^*(Q, D_\Omega))$ (given Q, D_Ω and p_i for each $i \in [n]$ that defines \mathcal{P}).*

527 If we want to approximate the expected multiplicities of all $Z = O(n^k)$ result tuples t
 528 simultaneously, we just need to run the above result with δ replaced by $\frac{\delta}{Z}$. Note this increases
 529 the runtime by only a logarithmic factor.

530 **5 Related Work**

531 **Probabilistic Databases** (PDBs) have been studied predominantly for set semantics.
 532 Approaches for probabilistic query processing (i.e., computing marginal probabilities of
 533 tuples), fall into two broad categories. *Intensional* (or *grounded*) query evaluation computes
 534 the *lineage* of a tuple and then the probability of the lineage formula. It has been shown
 535 that computing the marginal probability of a tuple is #P-hard [46] (by reduction from
 536 weighted model counting). The second category, *extensional* query evaluation, is in PTIME,
 537 but is limited to certain classes of queries. Dalvi et al. [14] and Olteanu et al. [21] proved
 538 dichotomies for UCQs and two classes of queries with negation, respectively. Amarilli et al.
 539 investigated tractable classes of databases for more complex queries [3]. Another line of work

540 studies which structural properties of lineage formulas lead to tractable cases [31, 41, 44]. In
541 this paper we focus on intensional query evaluation with polynomials.

542 Many data models have been proposed for encoding PDBs more compactly than as sets of
543 possible worlds. These include tuple-independent databases [47] (TIDBs), block-independent
544 databases (BIDBs) [42], and *PC-tables* [26]. Fink et al. [19] study aggregate queries over
545 a probabilistic version of the extension of K-relations for aggregate queries proposed in [4]
546 (*pvc-tables*) that supports bags, and has runtime complexity linear in the size of the lineage.
547 However, this lineage is encoded as a tree; the size (and thus the runtime) are still superlinear
548 in $T_{det}^*(Q, D_\Omega)$. The runtime bound is also limited to a specific class of (hierarchical) queries,
549 suggesting the possibility of a generalization of [14]’s dichotomy result to bag-PDBs.

550 Several techniques for approximating tuple probabilities have been proposed in related
551 work [20, 15, 38, 12], relying on Monte Carlo sampling, e.g., [12], or a branch-and-bound
552 paradigm [38]. Our approximation algorithm is also based on sampling.

553 **Compressed Encodings** are used for Boolean formulas (e.g, various types of circuits
554 including OBDDs [29]) and polynomials (e.g., factorizations [39]) some of which have been
555 utilized for probabilistic query processing, e.g., [29]. Compact representations for which
556 probabilities can be computed in linear time include OBDDs, SDDs, d-DNNF, and FBDD.
557 [16] studies circuits for absorptive semirings while [45] studies circuits that include negation
558 (expressed as the monus operation). Algebraic Decision Diagrams [7] (ADDs) generalize
559 BDDs to variables with more than two values. Chen et al. [10] introduced the generalized
560 disjunctive normal form. Appendix H covers more related work on fine-grained complexity.

561 **6 Conclusions and Future Work**

562 We have studied the problem of calculating the expected multiplicity of a query result tuple,
563 a problem that has a practical application in probabilistic databases over multisets. We show
564 that under various parameterized complexity hardness results/conjectures computing the
565 expected multiplicities exactly is not possible in time linear in the corresponding deterministic
566 query processing time. We prove that it is possible to approximate the expectation of a
567 lineage polynomial in linear time in the deterministic query processing over TIDBs and BIDBs
568 (assuming that there are few cancellations). Interesting directions for future work include
569 development of a dichotomy for bag PDBs. While we can handle higher moments (this
570 follows fairly easily from our existing results— see Appendix F), more general approximations
571 are an interesting area for exploration, including those for more general data models.

572 **References**

-
- 573 1 pdbench. <http://pdbench.sourceforge.net/>. Accessed: 2020-12-15.
- 574 2 Parag Agrawal, Omar Benjelloun, Anish Das Sarma, Chris Hayworth, Shubha U. Nabar,
575 Tomoe Sugihara, and Jennifer Widom. Trio: A system for data, uncertainty, and lineage. In
576 *VLDB*, pages 1151–1154, 2006.
- 577 3 Antoine Amarilli, Pierre Bourhis, and Pierre Senellart. Probabilities and provenance via tree
578 decompositions. *PODS*, 2015.
- 579 4 Yael Amsterdamer, Daniel Deutch, and Val Tannen. Provenance for aggregate queries. In
580 *PODS*, pages 153–164, 2011.
- 581 5 Lyublena Antova, Thomas Jansen, Christoph Koch, and Dan Olteanu. Fast and simple
582 relational processing of uncertain data.
- 583 6 Albert Atserias, Martin Grohe, and Dániel Marx. Size bounds and query plans for relational
584 joins. *SIAM J. Comput.*, 42(4):1737–1767, 2013. doi:10.1137/110859440.

- 585 7 R. Iris Bahar, Erica A. Frohm, Charles M. Gaona, Gary D. Hachtel, Enrico Macii, Abelardo
586 Pardo, and Fabio Somenzi. Algebraic decision diagrams and their applications. In *IEEE CAD*,
587 1993.
- 588 8 George Beskales, Ihab F. Ilyas, and Lukasz Golab. Sampling the repairs of functional
589 dependency violations under hard constraints. *Proc. VLDB Endow.*, 3(1):197–207, 2010.
- 590 9 Peter Bürgisser, Michael Clausen, and Mohammad Amin Shokrollahi. *Algebraic complexity*
591 *theory*, volume 315. Springer, 1997.
- 592 10 Hubie Chen and Martin Grohe. Constraint satisfaction with succinctly specified relations. *J.*
593 *Comput. Syst. Sci.*, 76(8):847–860, 2010.
- 594 11 Radu Curticapean. Counting matchings of size k is $w[1]$ -hard. In *ICALP*, volume 7965, pages
595 352–363, 2013.
- 596 12 N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. *VLDB*, 16(4):544,
597 2007.
- 598 13 Nilesh Dalvi and Dan Suciu. The dichotomy of conjunctive queries on probabilistic structures.
599 In *PODS*, pages 293–302, 2007.
- 600 14 Nilesh Dalvi and Dan Suciu. The dichotomy of probabilistic inference for unions of conjunctive
601 queries. *JACM*, 59(6):30, 2012.
- 602 15 Maarten Van den Heuvel, Peter Ivanov, Wolfgang Gatterbauer, Floris Geerts, and Martin
603 Theobald. Anytime approximation in probabilistic databases via scaled dissociations. In
604 *SIGMOD*, pages 1295–1312, 2019.
- 605 16 Daniel Deutch, Tova Milo, Sudeepa Roy, and Val Tannen. Circuits for datalog provenance. In
606 *ICDT*, pages 201–212, 2014.
- 607 17 Su Feng, Boris Glavic, Aaron Huber, and Oliver Kennedy. Efficient uncertainty tracking for
608 complex queries with attribute-level bounds. In *SIGMOD*, 2021.
- 609 18 Su Feng, Aaron Huber, Boris Glavic, and Oliver Kennedy. Uncertainty annotated databases -
610 a lightweight approach for approximating certain answers. In *SIGMOD*, 2019.
- 611 19 Robert Fink, Larisa Han, and Dan Olteanu. Aggregation in probabilistic databases via
612 knowledge compilation. *PVLDB*, 5(5):490–501, 2012.
- 613 20 Robert Fink, Jiewen Huang, and Dan Olteanu. Anytime approximation in probabilistic
614 databases. *VLDBJ*, 22(6):823–848, 2013.
- 615 21 Robert Fink and Dan Olteanu. Dichotomies for queries with negation in probabilistic databases.
616 *TODS*, 41(1):4:1–4:47, 2016.
- 617 22 Jörg Flum and Martin Grohe. *Parameterized Complexity Theory*. Texts in Theoretical
618 Computer Science. An EATCS Series. Springer, 2006. doi:10.1007/3-540-29953-X.
- 619 23 Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database systems - the*
620 *complete book (2. ed.)*. Pearson Education, 2009.
- 621 24 George Grätzer. *Universal algebra*. Springer Science & Business Media, 2008.
- 622 25 Todd J. Green, Gregory Karvounarakis, and Val Tannen. Provenance semirings. In *PODS*,
623 pages 31–40, 2007.
- 624 26 Todd J Green and Val Tannen. Models for incomplete and probabilistic information. In *EDBT*,
625 pages 278–296. 2006.
- 626 27 T. Imielinski and W. Lipski. Incomplete information in relational databases. 1989.
- 627 28 Tomasz Imieliński and Witold Lipski Jr. Incomplete information in relational databases.
628 *JACM*, 31(4):761–791, 1984.
- 629 29 Abhay Kumar Jha and Dan Suciu. Probabilistic databases with markovviews. *PVLDB*,
630 5(11):1160–1171, 2012.
- 631 30 Richard M. Karp, Michael Luby, and Neal Madras. Monte-carlo approximation algorithms for
632 enumeration problems. *J. Algorithms*, 10(3):429–448, 1989.
- 633 31 Batya Kenig, Avigdor Gal, and Ofer Strichman. A new class of lineage expressions over
634 probabilistic databases computable in p -time. In *SUM*, volume 8078, pages 219–232, 2013.
- 635 32 Oliver Kennedy and Christoph Koch. Pip: A database system for great and small expectations.
636 In *ICDE*, 2010.

- 637 33 Mahmoud Abo Khamis, Hung Q. Ngo, and Atri Rudra. Faq: Questions asked frequently. In
638 *PODS*, pages 13–28, 2016.
- 639 34 Tsvi Kopelowitz and Virginia Vassilevska Williams. Towards optimal set-disjointness and
640 set-intersection data structures. In *ICALP*, volume 168, pages 74:1–74:16, 2020.
- 641 35 Poonam Kumari, Said Achmiz, and Oliver Kennedy. Communicating data quality in on-demand
642 curation. In *QDB*, 2016.
- 643 36 Hung Q. Ngo. Worst-case optimal join algorithms: Techniques, results, and open problems.
644 In *PODS*, 2018.
- 645 37 Hung Q. Ngo, Christopher Ré, and Atri Rudra. Skew strikes back: new developments in the
646 theory of join algorithms. *SIGMOD Rec.*, 42(4):5–16, 2013.
- 647 38 Dan Olteanu, Jiewen Huang, and Christoph Koch. Approximate confidence computation in
648 probabilistic databases. In *ICDE*, pages 145–156, 2010.
- 649 39 Dan Olteanu and Maximilian Schleich. Factorized databases. *SIGMOD Rec.*, 45(2):5–16, 2016.
- 650 40 Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. Holoclean: Holistic data
651 repairs with probabilistic inference. *Proc. VLDB Endow.*, 10(11):1190–1201, 2017.
- 652 41 Sudeepa Roy, Vittorio Perduca, and Val Tannen. Faster query answering in probabilistic
653 databases using read-once functions. In *ICDT*, 2011.
- 654 42 C. Ré and D. Suciu. Materialized views in probabilistic databases: for information exchange
655 and query optimization. In *VLDB*, pages 51–62, 2007.
- 656 43 Christopher De Sa, Alexander Ratner, Christopher Ré, Jaeho Shin, Feiran Wang, Sen Wu, and
657 Ce Zhang. Incremental knowledge base construction using deepdive. *VLDB J.*, 26(1):81–105,
658 2017.
- 659 44 Prithviraj Sen, Amol Deshpande, and Lise Getoor. Read-once functions and query evaluation
660 in probabilistic databases. *PVLDB*, 3(1):1068–1079, 2010.
- 661 45 Pierre Senellart. Provenance and probabilities in relational databases. *SIGMOD Record*,
662 46(4):5–15, 2018.
- 663 46 Leslie G. Valiant. The complexity of enumeration and reliability problems. *SIAM J. Comput.*,
664 8(3):410–421, 1979.
- 665 47 Guy Van den Broeck and Dan Suciu. Query processing on probabilistic data: A survey. 2017.
- 666 48 Virginia Vassilevska Williams. Some open problems in fine-grained complexity. *SIGACT News*,
667 49(4):29–35, 2018. doi:10.1145/3300150.3300158.
- 668 49 Ying Yang, Niccolò Meneghetti, Ronny Fehling, Zhen Hua Liu, Dieter Gawlick, and Oliver
669 Kennedy. Lenses: An on-demand approach to etl. *PVLDB*, 8(12):1578–1589, 2015.

670 **7 Acknowledgements**

671 We thank Virginia Williams for showing us Eq. (20), which greatly simplified our earlier
672 proof of Lemma 3.8, and for graciously allowing us to use it.

673 **A** Generalizing Beyond Set Inputs

674 **A.1** TIDBs

675 In our definition of TIDBs (Sec. 2.1.1), we assumed a model of TIDBs where each input
 676 tuple is assigned a probability p of having multiplicity 1. That is, we assumed inputs to be
 677 sets, but interpret queries under bag semantics. Other sensible generalizations of TIDBs
 678 from set semantics to bag semantics also exist.

679 One very natural such generalization is to assign each input tuple t a multiplicity m_t and
 680 probability p : the tuple has probability p to exist with multiplicity m_t , and otherwise has
 681 multiplicity 0. If the maximal multiplicity of all input tuples in the TIDB is bounded by
 682 some constant, then a generalization of our hardness results and approximation algorithm
 683 can be achieved by changing the construction of lineage polynomials (in Fig. 1) as follows
 684 (all other cases remain the same as in fig. 1):

$$685 \quad \Phi[R, D_\Omega, t] = \begin{cases} m_t X_t & \text{if } D_\Omega.R(t) = m_t \\ 0 & \text{otherwise.} \end{cases}$$

687 That is the variable representing a tuple is multiplied by m_t to encode the tuple's multiplicity
 688 m_t . We note that our lower bounds still hold for this model since we only need $m_t = 1$ for all
 689 tuples t . Further, it can be argued that our proofs (as is) for approximation algorithms also
 690 work for this model. The only change is that since we now allow $m_t > 1$ some of the constants
 691 in the runtime analysis of our algorithms change but the overall asymptotic runtime bound
 692 remains the same.

693 Yet another option would be to assign each tuple a probability distribution over multiplicities.
 694 It seems very unlikely that our results would extend to a model that allows arbitrary
 695 probability distributions over multiplicities (our current proof techniques definitely break
 696 down). However, we would like to note that the special case of a Poisson binomial distribution
 697 (sum of independent but not necessarily identical Bernoulli trials) over multiplicities can be
 698 handled as follows: we add an additional identifier attribute to each relation in the database.
 699 For a tuple t with maximal multiplicity m_t , we create m_t copies of t with different identifiers.
 700 To answer a query over this encoding, we first project away the identifier attribute (note that
 701 as per Fig. 1, in Φ this would add up all the variables corresponding to the same tuple t).

702 **A.2** BIDBs

703 The approach described above works for BIDBs as well if we define the bag version of BIDBs
 704 to associate each tuple t a multiplicity m_t . Recall that we associate each tuple in a block
 705 with a unique variable. Thus, the modified lineage polynomial construction shown above can
 706 be applied for BIDBs too (and our approximation results also hold).

707 **B** Missing details from Section 2

708 **B.1** \mathcal{K} -relations and $\mathbb{N}[\mathbf{X}]$ -encoded PDBs

709 We can use \mathcal{K} -relations to model bags. A \mathcal{K} -relation [25] is a relation whose tuples
 710 are annotated with elements from a commutative semiring $\mathcal{K} = \{K, \oplus_{\mathcal{K}}, \otimes_{\mathcal{K}}, 0_{\mathcal{K}}, 1_{\mathcal{K}}\}$. A
 711 commutative semiring is a structure with a domain K and associative and commutative
 712 binary operations $\oplus_{\mathcal{K}}$ and $\otimes_{\mathcal{K}}$ such that $\otimes_{\mathcal{K}}$ distributes over $\oplus_{\mathcal{K}}$, $0_{\mathcal{K}}$ is the identity of $\oplus_{\mathcal{K}}$,

713 $\mathbb{1}_{\mathcal{K}}$ is the identity of $\otimes_{\mathcal{K}}$, and $\mathbb{0}_{\mathcal{K}}$ annihilates all elements of K when combined by $\otimes_{\mathcal{K}}$. Let
 714 \mathcal{U} be a countable domain of values. Formally, an n -ary \mathcal{K} -relation R over \mathcal{U} is a function
 715 $R : \mathcal{U}^n \rightarrow K$ with finite support $\text{supp}(R) = \{t \mid R(t) \neq \mathbb{0}_{\mathcal{K}}\}$. A \mathcal{K} -database is defined
 716 similarly, where we view the \mathcal{K} -database (relation) as a function mapping tuples to their
 717 respective annotations. \mathcal{RA}^+ query semantics over \mathcal{K} -relations are analogous to the lineage
 718 construction semantics of Fig. 1, with the exception of replacing $+$ with $\oplus_{\mathcal{K}}$ and \cdot with $\otimes_{\mathcal{K}}$.

719 Consider the semiring $\mathbb{N} = \{\mathbb{N}, +, \times, 0, 1\}$ of natural numbers. \mathbb{N} -databases model bag
 720 semantics by annotating each tuple with its multiplicity. A probabilistic \mathbb{N} -database (\mathbb{N} -PDB)
 721 is a PDB where each possible world is an \mathbb{N} -database. We study the problem of computing
 722 statistical moments for query results over such databases. Given an \mathbb{N} -PDB $\mathcal{D} = (\Omega, \mathcal{P})$,
 723 (\mathcal{RA}^+) query Q , and possible result tuple t , we sum $Q(D)(t) \cdot \mathcal{P}(D)$ for all $D \in \Omega$ to compute
 724 the expected multiplicity of t . Intuitively, the expectation of $Q(D)(t)$ is the number of
 725 duplicates of t we expect to find in result of query Q .

726 Let $\mathbb{N}[\mathbf{X}]$ denote the set of polynomials over variables $\mathbf{X} = (X_1, \dots, X_n)$ with natural
 727 number coefficients and exponents. Consider now the semiring (abusing notation) $\mathbb{N}[\mathbf{X}] =$
 728 $\{\mathbb{N}[\mathbf{X}], +, \cdot, 0, 1\}$ whose domain is $\mathbb{N}[\mathbf{X}]$, with the standard addition and multiplication of
 729 polynomials. We define an $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ as the tuple $(D_{\mathbb{N}[\mathbf{X}]}, \mathcal{P})$, where $\mathbb{N}[\mathbf{X}]$ -
 730 database $D_{\mathbb{N}[\mathbf{X}]}$ is paired with the probability distribution \mathcal{P} across the set of possible worlds
 731 represented by $D_{\mathbb{N}[\mathbf{X}]}$, i.e. the one induced from $\mathcal{P}_{\mathbb{N}[\mathbf{X}]}$, the probability distribution over \mathbf{X} .
 732 Note that the notation is slightly abused since the first element of the pair is an encoded
 733 set of possible worlds, i.e. $D_{\mathbb{N}[\mathbf{X}]}$ is the deterministic bounding database. We denote by
 734 $\Phi[Q, D_{\mathbb{N}[\mathbf{X}]}, t]$ the annotation of tuple t in the result of $Q(D_{\mathbb{N}[\mathbf{X}]})(t)$, and as before, interpret
 735 it as a function $\Phi[Q, D_{\mathbb{N}[\mathbf{X}]}, t] : \{0, 1\}^{|\mathbf{X}|} \rightarrow \mathbb{N}$ from vectors of variable assignments to the
 736 corresponding value of the annotating polynomial. $\mathbb{N}[\mathbf{X}]$ -encoded PDBs and a function Mod
 737 (which transforms an $\mathbb{N}[\mathbf{X}]$ -encoded PDB to an equivalent \mathbb{N} -PDB) are both formalized next.

738 To justify the use of $\mathbb{N}[\mathbf{X}]$ -databases, we need to show that we can encode any \mathbb{N} -PDB in
 739 this way and that the query semantics over this representation coincides with query semantics
 740 over its respective \mathbb{N} -PDB. For that it will be opportune to define representation systems for
 741 \mathbb{N} -PDBs.

Boris says:
cite

742 **► Definition B.1 (Representation System).** A representation system for \mathbb{N} -PDBs is a tuple
 743 (\mathcal{M}, Mod) where \mathcal{M} is a set of representations and Mod associates with each $M \in \mathcal{M}$ an
 744 \mathbb{N} -PDB \mathcal{D} . We say that a representation system is closed under a class of queries \mathcal{Q} if for
 745 any query $Q \in \mathcal{Q}$ and $M \in \mathcal{M}$ we have:

$$746 \quad Mod(Q(M)) = Q(Mod(M))$$

747 A representation system is complete if for every \mathbb{N} -PDB \mathcal{D} there exists $M \in \mathcal{M}$ such
 748 that:

$$749 \quad Mod(M) = \mathcal{D}$$

750 As mentioned above we will use $\mathbb{N}[\mathbf{X}]$ -databases paired with a probability distribution
 751 as a representation system, referring to such databases as $\mathbb{N}[\mathbf{X}]$ -encoded PDBs. Given
 752 $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$, one can think of the of \mathcal{P} as the probability distribution across
 753 all worlds $\{0, 1\}^n$. Denote a particular world to be \mathbf{w} . For convenience let $\psi_{\mathbf{w}} : \mathcal{D}_{\mathbb{N}[\mathbf{X}]} \rightarrow \mathcal{D}_{\mathbb{N}}$
 754 be a function that computes the corresponding \mathbb{N} -PDB upon assigning all values $w_i \in \mathbf{w}$ to
 755 $X_i \in \mathbf{X}$ of $D_{\mathbb{N}[\mathbf{X}]}$. Note the one-to-one correspondence between elements $\mathbf{w} \in \{0, 1\}^n$ to the
 756 worlds encoded by $D_{\mathbb{N}[\mathbf{X}]}$ when \mathbf{w} is assigned to \mathbf{X} (assuming a domain of $\{0, 1\}$ for each X_i).
 757 We can think of $\psi_{\mathbf{w}}(D_{\mathbb{N}[\mathbf{X}]})(t)$ as the semiring homomorphism $\mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$ that applies the
 assignment \mathbf{w} to all variables \mathbf{X} of a polynomial and evaluates the resulting expression in \mathbb{N} .

Boris says:
explain
connection to
homomorphism
lifting in
 \mathcal{K} -relations

759 ► **Definition B.2** (*Mod*($\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$)). Given an $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$, we compute its
 760 equivalent \mathbb{N} -PDB $\mathcal{D}_{\mathbb{N}} = \text{Mod}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = (\Omega, \mathcal{P}')$ as:

$$761 \quad \Omega = \{\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) \mid \mathbf{w} \in \{0, 1\}^n\}$$

$$762 \quad \forall D \in \Omega : Pr(D) = \sum_{\mathbf{w} \in \{0, 1\}^n : \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w})$$

764 For instance, consider a $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ consisting of a single tuple $t_1 = (1)$ annotated with
 765 $X_1 + X_2$ with probability distribution $Pr([0, 0]) = 0$, $Pr([0, 1]) = 0$, $Pr([1, 0]) = 0.3$ and
 766 $Pr([1, 1]) = 0.7$. This $\mathbb{N}[\mathbf{X}]$ -encoded PDB encodes two possible worlds (with non-zero
 767 probability) that we denote using their world vectors.

$$768 \quad D_{[0,1]}(t_1) = 1 \quad \text{and} \quad D_{[1,1]}(t_1) = 2$$

769 Importantly, as the following proposition shows, any finite \mathbb{N} -PDB can be encoded as an
 770 $\mathbb{N}[\mathbf{X}]$ -encoded PDB and $\mathbb{N}[\mathbf{X}]$ -encoded PDBs are closed under \mathcal{RA}^+ [25].

771 ► **Proposition B.3.** $\mathbb{N}[\mathbf{X}]$ -encoded PDBs are a complete representation system for \mathbb{N} -PDBs
 772 that is closed under \mathcal{RA}^+ queries.

773 **Proof.** To prove that $\mathbb{N}[\mathbf{X}]$ -encoded PDBs are complete consider the following construction
 774 that for any \mathbb{N} -PDB $\mathcal{D} = (\Omega, \mathcal{P})$ produces an $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]} = (\mathcal{D}_{\mathbb{N}[\mathbf{X}]}, \mathcal{P}')$ such
 775 that $\text{Mod}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = \mathcal{D}$. Let $\Omega = \{D_1, \dots, D_{|\Omega|}\}$. For each world D_i we create a corresponding
 776 variable X_i . In $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ we assign each tuple t the polynomial:

$$777 \quad D_{\mathbb{N}[\mathbf{X}]}(t) = \sum_{i=1}^{|\Omega|} D_i(t) \cdot X_i$$

778 The probability distribution \mathcal{P}' assigns all world vectors zero probability except for $|\Omega|$ world
 779 vectors (representing the possible worlds) \mathbf{w}_i . All elements of \mathbf{w}_i are zero except for the
 780 position corresponding to variables X_i which is set to 1. Unfolding definitions it is trivial
 781 to show that $\text{Mod}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = \mathcal{D}$. Thus, $\mathbb{N}[\mathbf{X}]$ -encoded PDBs are a complete representation
 782 system.

783 Since $\mathbb{N}[\mathbf{X}]$ is the free object in the variety of semirings, Birkhoff's HSP theorem implies
 784 that any assignment $\mathbf{X} \rightarrow \mathbb{N}$, which includes as a special case the assignments $\psi_{\mathbf{w}}$ used here,
 785 uniquely extends to the semiring homomorphism alluded to above, $\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) (t) : \mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$.
 786 For a polynomial $\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) (t)$ substitutes variables based on \mathbf{w} and then evaluates the
 787 resulting expression in \mathbb{N} . For instance, consider the polynomial $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}(t) = \Phi = X + Y$
 788 and assignment $\mathbf{w} := X = 0, Y = 1$. We get $\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) (t) = 0 + 1 = 1$. Closure under
 789 \mathcal{RA}^+ queries follows from this and from [25]'s Proposition 3.5, which states that semiring
 790 homomorphisms commute with queries over \mathcal{K} -relations. ◀

791 B.2 TIDBs and BIDBs in the $\mathbb{N}[\mathbf{X}]$ -encoded PDB model

792 Two important subclasses of $\mathbb{N}[\mathbf{X}]$ -encoded PDBs that are of interest to us are the bag
 793 versions of tuple-independent databases (TIDBs) and block-independent databases (BIDBs).
 794 Under set semantics, a TIDB is a deterministic database D where each tuple t is assigned
 795 a probability p_t . The set of possible worlds represented by a TIDB D is all subsets of D .
 796 The probability of each world is the product of the probabilities of all tuples that exist with
 797 one minus the probability of all tuples of D that are not part of this world, i.e., tuples are
 798 treated as independent random events. In a BIDB, we also assign each tuple a probability,

799 but additionally partition D into blocks. The possible worlds of a BIDB D are all subsets of
 800 D that contain at most one tuple from each block. Note then that the tuples sharing the
 801 same block are disjoint, and the sum of the probabilities of all the tuples in the same block b
 802 is at most 1. The probability of such a world is the product of the probabilities of all tuples
 803 present in the world. For bag TIDBs and BIDBs, we define the probability of a tuple to be
 804 the probability that the tuple exists with multiplicity at least 1.

805 In this work, we define TIDBs and BIDBs as subclasses of $\mathbb{N}[\mathbf{X}]$ -encoded PDBs defined
 806 over variables \mathbf{X} (Definition B.2) where \mathbf{X} can be partitioned into blocks that satisfy the
 807 conditions of a TIDB or BIDB (stated formally in Sec. 2.1.1). In this work, we consider
 808 one further deviation from the standard: We use bag semantics for queries. Even though
 809 tuples cannot occur more than once in the input TIDB or BIDB, they can occur with a
 810 multiplicity larger than one in the result of a query. Since in TIDBs and BIDBs, there is a
 811 one-to-one correspondence between tuples in the database and variables, we can interpret a
 812 vector $\mathbf{w} \in \{0, 1\}^n$ as denoting which tuples exist in the possible world $\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ (the ones
 813 where $w_i = 1$). For BIDBs specifically, note that at most one of the bits corresponding to
 814 tuples in each block will be set (i.e., for any pair of bits $w_j, w_{j'}$ that are part of the same
 815 block $b_i \supseteq \{t_{i,j}, t_{i,j'}\}$, at most one of them will be set). Denote the vector \mathbf{p} to be a vector
 816 whose elements are the individual probabilities p_i of each tuple t_i . Given PDB \mathcal{D} \mathcal{P} is the
 817 distribution induced by \mathbf{p} , which we will denote $\mathcal{P}(\mathbf{p})$.

$$818 \quad \mathbb{E}_{\mathbf{w} \sim \mathcal{P}(\mathbf{p})} [\Phi(\mathbf{W})] = \sum_{\substack{\mathbf{w} \in \{0,1\}^n \\ s.t. w_j, w_{j'} = 1 \rightarrow \exists b_i \supseteq \{t_{i,j}, t_{i,j'}\}}} \Phi(\mathbf{w}) \prod_{\substack{j \in [n] \\ s.t. w_j = 1}} p_j \prod_{\substack{j \in [n] \\ s.t. w_j = 0}} (1 - p_j) \quad (5)$$

820 Recall that tuple blocks in a TIDB always have size 1, so the outer summation of eq. (5) is
 821 over the full set of vectors.

Boris says:

822 Oliver's
 823 conjecture:
 824 Bag-TIDBs +
 825 Q can express
 826 any finite
 827 bag-PDB: A
 828 well-known
 829 result for set
 830 semantics
 831 PDBs is
 832 that while
 833 not all finite
 834 PDBs can
 835 be encoded
 836 as TIDBs,
 837 any finite
 838 PDB can
 839 be encoded
 840 using a TIDB
 841 and a query.
 842 An analog
 843 result holds
 844 in our case:
 845 any finite
 846 N-PDB can
 847 be encoded
 848 as a bag
 849 TIDB and a
 850 query (WHAT
 851 CLASS? ADD
 852 PROOF)

823 B.3 Proof of Proposition 2.5

Proof. We need to prove for \mathbb{N} -PDB $\mathcal{D} = (\Omega, \mathcal{P})$ and $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]} = (D'_{\mathbb{N}[\mathbf{X}]}, \mathcal{P}')$
 824 where $Mod(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = \mathcal{D}$ that $\mathbb{E}_{\mathcal{D} \sim \mathcal{P}} [Q(\mathcal{D})(t)] = \mathbb{E}_{\mathbf{w} \sim \mathcal{P}'} [\Phi[Q, D_{\mathbb{N}[\mathbf{X}]}], t](\mathbf{w})]$ By expanding
 825 $\Phi[Q, D_{\mathbb{N}[\mathbf{X}]}], t]$ and the expectation we have:

$$826 \quad \mathbb{E}_{\mathbf{w} \sim \mathcal{P}'} [\Phi(\mathbf{W})] = \sum_{\mathbf{w} \in \{0,1\}^n} Pr(\mathbf{w}) \cdot Q(D_{\mathbb{N}[\mathbf{X}]})(t)(\mathbf{w})$$

827 From $Mod(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = \mathcal{D}$, we have that the range of $\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ is Ω , so

$$828 \quad = \sum_{D \in \Omega} \sum_{\mathbf{w} \in \{0,1\}^n: \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w}) \cdot Q(D_{\mathbb{N}[\mathbf{X}]})(t)(\mathbf{w})$$

829 The inner sum is only over \mathbf{w} where $\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D$ (i.e., $Q(D_{\mathbb{N}[\mathbf{X}]})(t)(\mathbf{w}) = Q(D)(t)$)

$$830 \quad = \sum_{D \in \Omega} \sum_{\mathbf{w} \in \{0,1\}^n: \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w}) \cdot Q(D)(t)$$

831 By distributivity of $+$ over \times

$$832 \quad = \sum_{D \in \Omega} Q(D)(t) \sum_{\mathbf{w} \in \{0,1\}^n: \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w})$$

833 From the definition of \mathcal{P} in definition B.2, given $Mod(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = \mathcal{D}$, we get

$$834 \quad = \sum_{D \in \Omega} Q(D)(t) \cdot Pr(D) = \mathbb{E}_{\mathcal{D} \sim \mathcal{P}} [Q(\mathcal{D})(t)]$$

◀

842 B.4 Proposition B.4

843 Note the following fact:

844 ► **Proposition B.4.** *For any BIDB-lineage polynomial $\Phi(X_1, \dots, X_n)$ and all \mathbf{w} such that*
 845 *$Pr[\mathbf{W} = \mathbf{w}] > 0$, it holds that $\Phi(\mathbf{w}) = \tilde{\Phi}(\mathbf{w})$.*

846 **Proof.** Note that any Φ in factorized form is equivalent to its SMB expansion. For each
 847 term in the expanded form, further note that for all $b \in \{0, 1\}$ and all $e \geq 1$, $b^e = b$.
 848 Finally, note that there are exactly three cases where the expectation of a monomial term
 849 $\mathbb{E} \left[c_{\mathbf{d}} \prod_{i=1}^n X_i^{d_i} \right]$ is zero: (i) when $c_{\mathbf{d}} = 0$, (ii) when $p_i = 0$ for some i where $\mathbf{d}_i \geq 1$,
 850 and (iii) when X_i and X_j are in the same block for some i, j where $\mathbf{d}_i, \mathbf{d}_j \geq 1$. ◀

851 B.5 Proof for Lemma 1.4

852 **Proof.** Let Φ be a polynomial of n variables with highest degree = B , defined as follows:

$$853 \quad \Phi(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \{0, \dots, B\}^n} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n X_i^{d_i}.$$

854 Let the boolean function $\text{ISIND}(\cdot)$ take \mathbf{d} as input and return true if there does not exist any
 855 dependent variables in \mathbf{d} , i.e., $\nexists b, i \neq j \mid d_{b,i}, d_{b,j} \geq 1$.⁸ Then in expectation we have

$$856 \quad \mathbb{E}_{\mathbf{W}} [\Phi(\mathbf{W})] = \mathbb{E}_{\mathbf{W}} \left[\sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n W_i^{d_i} + \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \neg \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n W_i^{d_i} \right] \quad (6)$$

$$857 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[\prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n W_i^{d_i} \right] + \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \neg \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[\prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n W_i^{d_i} \right] \quad (7)$$

$$858 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[\prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n W_i^{d_i} \right] \quad (8)$$

$$859 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n \mathbb{E}_{\mathbf{W}} [W_i^{d_i}] \quad (9)$$

$$860 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n \mathbb{E}_{\mathbf{W}} [W_i] \quad (10)$$

$$861 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n p_i \quad (11)$$

$$862 \quad = \tilde{\Phi}(p_1, \dots, p_n). \quad (12)$$

863

⁸ This BIDB notation is used and discussed in sec. 2.1.1

Eq. (6) is the result of substituting in the definition of Φ given above. Then we arrive at eq. (7) by linearity of expectation. Next, eq. (8) is the result of the independence constraint of BIBDs, specifically that any monomial composed of dependent variables, i.e., variables from the same block b , has a probability of 0. Eq. (9) is obtained by the fact that all variables in each monomial are independent, which allows for the expectation to be pushed through the product. In eq. (10), since $W_i \in \{0, 1\}$ it is the case that for any exponent $e \geq 1$, $W_i^e = W_i$. Next, in eq. (11) the expectation of a tuple is indeed its probability.

Finally, it can be verified that Eq. (12) follows since eq. (11) satisfies the construction of $\tilde{\Phi}(p_1, \dots, p_n)$ in Definition 1.3. ◀

873 B.6 Proof For Corollary 2.4

Proof. Note that Lemma 1.4 shows that $\mathbb{E}[\Phi] = \tilde{\Phi}(p_1, \dots, p_n)$. Therefore, if Φ is already in SMB form, one only needs to compute $\tilde{\Phi}(p_1, \dots, p_n)$ ignoring exponent terms (note that such a polynomial is $\tilde{\Phi}(p_1, \dots, p_n)$), which indeed has $O(|\Phi|)$ computations. ◀

877 C Missing details from Section 3

878 C.1 Lemma C.1

► **Lemma C.1.** *Assuming that each $v \in V$ has degree ≥ 1 ,⁹ the PDB relations encoding the edges for Φ_G^k of Definition 3.4 can be computed in $O(m)$ time.*

Proof of Lemma C.1. Only two relations need be constructed, one for the set V and one for the set E . By a simple linear scan, each can be constructed in time $O(m+n)$. Given that the degree of each $v \in V$ is at least 1, we have that $m \geq \Omega(n)$, and this yields the claimed runtime. ◀

885 C.2 Proof of Lemma 3.5

Proof. By the recursive definition of $T_{det}(\cdot, \cdot)$ (see Sec. 2.3), we have the following equation for our hard query Q when $k = 1$, (we denote this as Q^1).

$$888 \quad T_{det}(Q^1, D) = |D.V| + |D.E| + |D.V| + T_{join}(D.V, D.E, D.V).$$

We argue that $T_{join}(D.V, D.E, D.V)$ is at most $O(m)$ by noting that there exists an algorithm that computes $D.V \bowtie D.E \bowtie D.V$ in the same runtime¹⁰. Then by the assumption of Lemma C.1 (each $v \in V$ has degree ≥ 1), the sum of the first three terms is $O(m)$. We then obtain that $T_{det}(Q^1, D) = O(m) + O(m) = O(m)$. For $Q^k = Q_1^1 \times \dots \times Q_k^1$, we have the recurrence $T_{det}(Q^k, D) = T_{det}(Q_1^1, D) + \dots + T_{det}(Q_k^1, D) + T_{join}(Q_1^1, \dots, Q_k^1)$. Since Q^1 outputs a count, computing the join $Q_1^1 \bowtie \dots \bowtie Q_k^1$ is just multiplying k numbers, which takes $O(k)$ time. Thus, we have

$$896 \quad T_{det}(Q^k, D) \leq k \cdot O(m) + O(k) \leq O(km),$$

897 as desired. ◀

⁹ This is WLOG, since any vertex with degree 0 can be dropped without affecting the result of our hard query.

¹⁰ Indeed the trivial algorithm that computes the obvious pair-wise joins has the claimed runtime. That is, we first compute $D.V \bowtie D.E$, which takes $O(m)$ (assuming $D.V$ is stored in hash map) since tuples in $D.V$ can only filter tuples in $D.E$. The resulting subset of tuples in $D.E$ are then again joined (on the right) with $D.V$, which by the same argument as before also takes $O(m)$ time, as desired.

898 C.3 Lemma C.2

899 The following lemma reduces the problem of counting k -matchings in a graph to our problem
900 (and proves Theorem 3.6):

901 ► **Lemma C.2.** *Let p_0, \dots, p_{2k} be distinct values in $(0, 1]$. Then given the values $\tilde{\Phi}_G^k(p_i, \dots, p_i)$
902 for $0 \leq i \leq 2k$, the number of k -matchings in G can be computed in $O(k^3)$ time.*

903 C.4 Proof of Lemma C.2

904 **Proof.** We first argue that $\tilde{\Phi}_G^k(p, \dots, p) = \sum_{i=0}^{2k} c_i \cdot p^i$. First, since $\Phi_G(\mathbf{X})$ has degree 2, it
905 follows that $\Phi_G^k(\mathbf{X})$ has degree $2k$. By definition, $\tilde{\Phi}_G^k(\mathbf{X})$ sets every exponent $e > 1$ to $e = 1$,
906 which means that $\text{DEG}(\tilde{\Phi}_G^k) \leq \text{DEG}(\Phi_G^k) = 2k$. Thus, if we think of p as a variable, then
907 $\tilde{\Phi}_G^k(p, \dots, p)$ is a univariate polynomial of degree at most $\text{DEG}(\tilde{\Phi}_G^k) \leq 2k$. Thus, we can write

$$908 \quad \tilde{\Phi}_G^k(p, \dots, p) = \sum_{i=0}^{2k} c_i p^i$$

909 We note that c_i is *exactly* the number of monomials in the SMB expansion of $\Phi_G^k(\mathbf{X})$ composed
910 of i distinct variables.¹¹

911 Given that we then have $2k + 1$ distinct values of $\tilde{\Phi}_G^k(p, \dots, p)$ for $0 \leq i \leq 2k$, it follows
912 that we have a linear system of the form $\mathbf{M} \cdot \mathbf{c} = \mathbf{b}$ where the i th row of \mathbf{M} is $(p_i^0 \dots p_i^{2k})$,
913 \mathbf{c} is the coefficient vector (c_0, \dots, c_{2k}) , and \mathbf{b} is the vector such that $\mathbf{b}[i] = \tilde{\Phi}_G^k(p_i, \dots, p_i)$.
914 In other words, matrix \mathbf{M} is the Vandermonde matrix, from which it follows that we have
915 a matrix with full rank (the p_i 's are distinct), and we can solve the linear system in $O(k^3)$
916 time (e.g., using Gaussian Elimination) to determine \mathbf{c} exactly. Thus, after $O(k^3)$ work, we
917 know \mathbf{c} and in particular, c_{2k} exactly.

918 Next, we show why we can compute $\#(G, \mathfrak{I} \dots \mathfrak{I}^k)$ from c_{2k} in $O(1)$ additional time. We
919 claim that c_{2k} is $k! \cdot \#(G, \mathfrak{I} \dots \mathfrak{I}^k)$. This can be seen intuitively by looking at the expansion
920 of the original factorized representation

$$921 \quad \Phi_G^k(\mathbf{X}) = \sum_{(i_1, j_1), \dots, (i_k, j_k) \in E} X_{i_1} X_{j_1} \dots X_{i_k} X_{j_k},$$

922 where a unique k -matching in the multi-set of product terms can be selected $\prod_{i=1}^k i = k!$
923 times. Indeed, note that each k -matching $(i_1, j_1) \dots (i_k, j_k)$ in G corresponds to the monomial
924 $\prod_{\ell=1}^k X_{i_\ell} X_{j_\ell}$ in $\Phi_G^k(\mathbf{X})$, with distinct indexes, and this implies that each distinct k -matching
925 appears the exact number of permutations that exist for its particular set of k edges, or $k!$.

926 Since, as noted earlier, c_{2k} represents the number of monomials with $2k$ distinct variables,
927 then it must be that c_{2k} is the overall number of k -matchings. And since we have $k!$ copies
928 of each distinct k -matching, it follows that $c_{2k} = k! \cdot \#(G, \mathfrak{I} \dots \mathfrak{I}^k)$. Thus, simply dividing
929 c_{2k} by $k!$ gives us $\#(G, \mathfrak{I} \dots \mathfrak{I}^k)$, as needed. ◀

930 C.5 Proof of Theorem 3.6

931 **Proof.** For the sake of contradiction, assume we can solve our problem in $o(T_{\text{match}}(k, G))$
932 time. Given a graph G by Lemma C.1 we can compute the PDB encoding in $O(m)$ time. Then

¹¹ Since $\tilde{\Phi}_G^k(\mathbf{X})$ does not have any monomial with degree < 2 , it is the case that $c_0 = c_1 = 0$ but for the sake of simplicity we will ignore this observation.

933 after we run our algorithm on $\tilde{\Phi}_G^k$, we get $\tilde{\Phi}_G^k(p_i, \dots, p_i)$ for every $0 \leq i \leq 2k$ in additional
 934 $O(k) \cdot o(T_{match}(k, G))$ time. Lemma C.2 then computes the number of k -matchings in G in
 935 $O(k^3)$ time. Adding the runtime of all of these steps, we have an algorithm for computing
 936 the number of k -matchings that runs in time

$$937 \quad O(m) + O(k) \cdot o(T_{match}(k, G)) + O(k^3) \quad (13)$$

$$938 \quad \leq o(T_{match}(k, G)). \quad (14)$$

940 We obtain Eq. (14) from the facts that k is fixed (related to m) and the assumption that
 941 $T_{match}(k, G) \geq \omega(m)$. Thus we obtain the contradiction that we can achieve a runtime
 942 $o(T_{match}(k, G))$ that is better than the optimal time $T_{match}(k, G)$ required to compute
 943 k -matchings. \blacktriangleleft

944 C.6 Subgraph Notation and $O(1)$ Closed Formulas

945 We need all the possible edge patterns in an arbitrary G with at most three distinct edges.
 946 We have already seen \mathfrak{A} , \mathfrak{B} and \mathfrak{C} , so we define the remaining patterns:

- 947 ■ Single Edge (\mathfrak{D})
- 948 ■ 2-path (\mathfrak{E})
- 949 ■ 2-matching (\mathfrak{F})
- 950 ■ 3-star (\mathfrak{G})—this is the graph that results when all three edges share exactly one common
 951 endpoint. The remaining endpoint for each edge is disconnected from any endpoint of
 952 the remaining two edges.
- 953 ■ Disjoint Two-Path (\mathfrak{H} , \mathfrak{I})—this subgraph consists of a two-path and a remaining disjoint
 954 edge.

955 For any graph G , the following formulas for $\#(G, H)$ compute their respective patterns
 956 exactly in $O(m)$ time, with d_i representing the degree of vertex i (proofs are in Appendix C.7):

$$957 \quad \#(G, \mathfrak{D}) = m, \quad (15)$$

$$958 \quad \#(G, \mathfrak{E}) = \sum_{i \in V} \binom{d_i}{2} \quad (16)$$

$$959 \quad \#(G, \mathfrak{F}) = \sum_{(i,j) \in E} \frac{m - d_i - d_j + 1}{2} \quad (17)$$

$$960 \quad \#(G, \mathfrak{G}) = \sum_{i \in V} \binom{d_i}{3} \quad (18)$$

$$961 \quad \#(G, \mathfrak{H}) + 3\#(G, \mathfrak{I}) = \sum_{(i,j) \in E} \binom{m - d_i - d_j + 1}{2} \quad (19)$$

$$962 \quad \#(G, \mathfrak{B}) + 3\#(G, \mathfrak{A}) = \sum_{(i,j) \in E} (d_i - 1) \cdot (d_j - 1) \quad (20)$$

964 C.7 Proofs of Eq. (15)-Eq. (20)

965 The proofs for Eq. (15), Eq. (16) and Eq. (18) are immediate.

966 **Proof of Eq. (17).** For edge (i, j) connecting arbitrary vertices i and j , finding all other
 967 edges in G disjoint to (i, j) is equivalent to finding all edges that are not connected to either

968 vertex i or j . The number of such edges is $m - d_i - d_j + 1$, where we add 1 since edge (i, j)
 969 is removed twice when subtracting both d_i and d_j . Since the summation is iterating over
 970 all edges such that a pair $((i, j), (k, \ell))$ will also be counted as $((k, \ell), (i, j))$, division by 2
 971 then eliminates this double counting. Note that m and d_i for all $i \in V$ can be computed in
 972 one pass over the set of edges by simply maintaining counts for each quantity. Finally, the
 973 summation is also one traversal through the set of edges where each operation is either a
 974 lookup ($O(1)$ time) or an addition operation (also $O(1)$) time. ◀

975 **Proof of Eq. (19).** Eq. (19) is true for similar reasons. For edge (i, j) , it is necessary to find
 976 two additional edges, disjoint or connected. As in our argument for Eq. (17), once the number
 977 of edges disjoint to (i, j) have been computed, then we only need to consider all possible
 978 combinations of two edges from the set of disjoint edges, since it doesn't matter if the two
 979 edges are connected or not. Note, the factor 3 of \mathfrak{iii} is necessary to account for the triple
 980 counting of 3-matchings, since it is indistinguishable to the closed form expression which of
 981 the remaining edges are either disjoint or connected to each of the edges in the *initial* set of
 982 edges disjoint to the edge under consideration. Observe that the disjoint case will be counted
 983 3 times since each edge of a 3-path is visited once, and the same 3-path counted in each
 984 visitation. For the latter case however, it is true that since the two path in \mathfrak{ia} is connected,
 985 there will be no multiple counting by the fact that the summation automatically disconnects
 986 the current edge, meaning that a two matching at the current vertex under consideration
 987 will not be counted. Thus, \mathfrak{ia} will only be counted once, precisely when the single disjoint
 988 edge is visited in the summation. The sum over all such edge combinations is precisely then
 989 $\#(G, \mathfrak{ia}) + 3\#(G, \mathfrak{iii})$. Note that all factorials can be computed in $O(m)$ time, and then
 990 each combination $\binom{n}{2}$ can be performed with constant time operations, yielding the claimed
 991 $O(m)$ run time. ◀

992 **Proof of Eq. (20).** To compute $\#(G, \mathfrak{ii})$, note that for an arbitrary edge (i, j) , a 3-path
 993 exists for edge pair (i, ℓ) and (j, k) where i, j, k, ℓ are distinct. Further, the quantity $(d_i -$
 994 $1) \cdot (d_j - 1)$ represents the number of 3-edge subgraphs with middle edge (i, j) and outer
 995 edges $(i, \ell), (j, k)$ such that $\ell \neq j$ and $k \neq i$. When $k = \ell$, the resulting subgraph is a triangle,
 996 and when $k \neq \ell$, the subgraph is a 3-path. Summing over all edges (i, j) gives Eq. (20) by
 997 observing that each triangle is counted thrice, while each 3-path is counted just once. For
 998 reasons similar to Eq. (17), all d_i can be computed in $O(m)$ time and each summand can
 999 then be computed in $O(1)$ time, yielding an overall $O(m)$ run time. ◀

1000 C.8 Tools to prove Theorem 3.7

1001 Note that $\tilde{\Phi}_G^3(p, \dots, p)$ as a polynomial in p has degree at most six. Next, we figure out the
 1002 exact coefficients since this would be useful in our arguments:

1003 ▶ **Lemma C.3.** *For any p , we have:*

$$1004 \quad \tilde{\Phi}_G^3(p, \dots, p) = \#(G, \mathfrak{i})p^2 + 6\#(G, \mathfrak{ia})p^3 + 6\#(G, \mathfrak{iii})p^4 + 6\#(G, \mathfrak{iaa})p^3 \\ 1005 \quad + 6\#(G, \mathfrak{iaa})p^4 + 6\#(G, \mathfrak{iiia})p^4 + 6\#(G, \mathfrak{iaia})p^5 + 6\#(G, \mathfrak{iiiaa})p^6. \quad (21)$$

1007 C.8.1 Proof for Lemma C.3

1008 **Proof.** By definition we have that

$$1009 \quad \Phi_G^3(\mathbf{X}) = \sum_{(i_1, j_1), (i_2, j_2), (i_3, j_3) \in E} \prod_{\ell=1}^3 X_{i_\ell} X_{j_\ell}.$$

1010 Hence $\tilde{\Phi}_G^3(\mathbf{X})$ has degree six. Note that the monomial $\prod_{\ell=1}^3 X_{i_\ell} X_{j_\ell}$ will contribute to the
 1011 coefficient of p^ν in $\tilde{\Phi}_G^3(\mathbf{X})$, where ν is the number of distinct variables in the monomial. Let
 1012 $e_1 = (i_1, j_1)$, $e_2 = (i_2, j_2)$, and $e_3 = (i_3, j_3)$. We compute $\tilde{\Phi}_G^3(\mathbf{X})$ by considering each of the
 1013 three forms that the triple (e_1, e_2, e_3) can take.

1014 CASE 1: $e_1 = e_2 = e_3$ (all edges are the same). When we have that $e_1 = e_2 = e_3$, then
 1015 the monomial corresponds to $\#(G, \mathfrak{I})$. There are exactly m such triples, each with a p^2
 1016 factor in $\tilde{\Phi}_G^3(p, \dots, p)$.

1017 CASE 2: This case occurs when there are two distinct edges of the three, call them e and
 1018 e' . When there are two distinct edges, there is then the occurrence when 2 variables in the
 1019 triple (e_1, e_2, e_3) are bound to e . There are three combinations for this occurrence in $\Phi_G^3(\mathbf{X})$.
 1020 Analogously, there are three such occurrences in $\Phi_G^3(\mathbf{X})$ when there is only one occurrence of
 1021 e , i.e. 2 of the variables in (e_1, e_2, e_3) are e' . This implies that all $3 + 3 = 6$ combinations of
 1022 two distinct edges e and e' contribute to the same monomial in $\tilde{\Phi}_G^3$. Since $e \neq e'$, this case
 1023 produces the following edge patterns: \mathfrak{A} , \mathfrak{B} , which contribute $6p^3$ and $6p^4$ respectively to
 1024 $\tilde{\Phi}_G^3(p, \dots, p)$.

1025 CASE 3: All e_1, e_2 and e_3 are distinct. For this case, we have $3! = 6$ permutations
 1026 of (e_1, e_2, e_3) , each of which contribute to the same monomial in the SMB representation
 1027 of $\Phi_G^3(\mathbf{X})$. This case consists of the following edge patterns: \mathfrak{C} , \mathfrak{D} , \mathfrak{E} , \mathfrak{F} , \mathfrak{G} , \mathfrak{H} , which
 1028 contribute $6p^3, 6p^4, 6p^4, 6p^5$ and $6p^6$ respectively to $\tilde{\Phi}_G^3(p, \dots, p)$. \blacktriangleleft

1029 Since p is fixed, Lemma C.3 gives us one linear equation in $\#(G, \mathfrak{A})$ and $\#(G, \mathfrak{H})$ (we
 1030 can handle the other counts due to equations (15)-(20)). However, we need to generate
 1031 one more independent linear equation in these two variables. Towards this end we generate
 1032 another graph related to G :

1033 **► Definition C.4.** For $\ell \geq 1$, let graph $G^{(\ell)}$ be a graph generated from an arbitrary graph
 1034 G , by replacing every edge e of G with an ℓ -path, such that all inner vertexes of an ℓ -path
 1035 replacement edge are disjoint from all other vertexes.¹²

1036 We will prove Theorem 3.7 by the following reduction:

1037 **► Theorem C.5.** Fix $p \in (0, 1)$. Let G be a graph on m edges. If we can compute $\tilde{\Phi}_G^3(p, \dots, p)$
 1038 exactly in $T(m)$ time, then we can exactly compute $\#(G, \mathfrak{A})$ in $O(T(m) + m)$ time.

1039 For clarity, we repeat the notion of $\#(G, H)$ to mean the count of subgraphs in G isomorphic
 1040 to H . The following lemmas relate these counts in $G^{(2)}$ to $G^{(1)}$ (G). The lemmas are used
 1041 to prove Lemma C.8.

1042 **► Lemma C.6.** The 3-matchings in graph $G^{(2)}$ satisfy the identity:

$$1043 \quad \#(G^{(2)}, \mathfrak{H}) = 8 \cdot \#(G^{(1)}, \mathfrak{H}) + 6 \cdot \#(G^{(1)}, \mathfrak{I} \mathfrak{A}) \\ 1044 \quad \quad \quad + 4 \cdot \#(G^{(1)}, \mathfrak{D}) + 4 \cdot \#(G^{(1)}, \mathfrak{E}) + 2 \cdot \#(G^{(1)}, \mathfrak{A}).$$

1046 **► Lemma C.7.** For $\ell > 1$ and any graph $G^{(\ell)}$, $\#(G^{(\ell)}, \mathfrak{A}) = 0$.

1047 Finally, the following result immediately implies Theorem C.5:

1048 **► Lemma C.8.** Fix $p \in (0, 1)$. Given $\tilde{\Phi}_{G^{(\ell)}}^3(p, \dots, p)$ for $\ell \in [2]$, we can compute in $O(m)$
 1049 time a vector $\mathbf{b} \in \mathbb{R}^3$ such that

$$1050 \quad \begin{pmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{pmatrix} \cdot \begin{pmatrix} \#(G, \mathfrak{A}) \\ \#(G, \mathfrak{H}) \end{pmatrix} = \mathbf{b},$$

¹²Note that $G \equiv G^{(1)}$.

1051 allowing us to compute $\#(G, \mathfrak{A})$ and $\#(G, \mathfrak{B})$ in $O(1)$ time.

1052 C.9 Proofs for Lemma C.6, Lemma C.7, and Lemma C.8

1053 Before proceeding, let us introduce a few more helpful definitions.

1054 ► **Definition C.9** ($E^{(\ell)}$). For $\ell > 1$, we use $E^{(\ell)}$ to denote the set of edges in $G^{(\ell)}$. For
 1055 any graph $G^{(\ell)}$, its edges are denoted by the a pair (e, b) , such that $b \in \{0, \dots, \ell - 1\}$ where
 1056 $(e, 0), \dots, (e, \ell - 1)$ is the ℓ -path that replaces the edge e for $e \in E^{(1)}$.

1057 ► **Definition C.10** ($E_S^{(\ell)}$). Given an arbitrary subgraph $S^{(1)}$ of $G^{(1)}$, let $E_S^{(1)}$ denote the set
 1058 of edges in $S^{(1)}$. Define then $E_S^{(\ell)}$ for $\ell > 1$ as the set of edges in the generated subgraph $S^{(\ell)}$
 1059 (i.e. when we apply Definition C.4 to S to generate $S^{(\ell)}$).

1060 For example, consider $S^{(1)}$ with edges $E_S^{(1)} = \{e_1\}$. Then the edge set of $S^{(2)}$ is defined
 1061 as $E_S^{(2)} = \{(e_1, 0), (e_1, 1)\}$.

1062 ► **Definition C.11** ($\binom{E}{t}$ and $\binom{E}{\leq t}$). Let $\binom{E}{t}$ denote the set of subsets in E with exactly t
 1063 edges. In a similar manner, $\binom{E}{\leq t}$ is used to mean the subsets of E with t or fewer edges.

1064 The following function f_ℓ is a mapping from every 3-edge shape in $G^{(\ell)}$ to its ‘projection’
 1065 in $G^{(1)}$.

1066 ► **Definition C.12.** Let $f_\ell : \binom{E^{(\ell)}}{3} \rightarrow \binom{E^{(1)}}{\leq 3}$ be defined as follows. For any element $s \in \binom{E^{(\ell)}}{3}$
 1067 such that $s = \{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}$, define:

$$1068 \quad f_\ell(\{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}) = \{e_1, e_2, e_3\}.$$

1069 ► **Definition C.13** (f_ℓ^{-1}). For an arbitrary subgraph $S^{(1)}$ of $G^{(1)}$ with at most $m \leq 3$ edges,
 1070 the inverse function $f_\ell^{-1} : \binom{E^{(1)}}{\leq 3} \rightarrow 2^{\binom{E^{(\ell)}}{3}}$ takes $E_S^{(1)}$ and outputs the set of all elements
 1071 $s \in \binom{E_S^{(\ell)}}{3}$ such that $f_\ell(s) = E_S^{(1)}$.

1072 Note, importantly, that when we discuss f_ℓ^{-1} , that each edge present in $E_S^{(1)}$ must have an
 1073 edge in $s \in f_\ell^{-1}(E_S^{(1)})$ that projects down to it. In particular, if $|E_S^{(1)}| = 3$, then it must be the
 1074 case that each $s \in f_\ell^{-1}(E_S^{(1)})$ consists of the following set of edges: $\{(e_i, b), (e_j, b'), (e_m, b'')\}$,
 1075 where i, j and m are distinct.

1076 We are now ready to prove the structural lemmas. To prove the structural lemmas, we
 1077 will count the number of occurrences of \mathfrak{A} and \mathfrak{B} in $G^{(\ell)}$ we count for each $S \in \binom{E_1}{\leq 3}$, how
 1078 many \mathfrak{B} and \mathfrak{A} subgraphs appear in $f_\ell^{-1}(E_S^{(1)})$.

1079 C.9.1 Proof of Lemma C.6

1080 **Proof.** For each subset $E_S^{(1)} \in \binom{E_1}{\leq 3}$, we count the number of 3-matchings in the 3-edge
 1081 subgraphs of $G^{(2)}$ in $f_2^{-1}(E_S^{(1)})$. We first consider the case of $E_S^{(1)} \in \binom{E_1}{3}$, where $E_S^{(1)}$
 1082 is composed of the edges e_1, e_2, e_3 and $f_2^{-1}(E_S^{(1)})$ is the set of all 3-edge subsets $s \in$
 1083 $\{(e_1, 0), (e_1, 1), (e_2, 0), (e_2, 1), (e_3, 0), (e_3, 1)\}$ such that $f_\ell(s) = \{e_1, e_2, e_3\}$. The size of
 1084 the output is denoted $|f_2^{-1}(E^{(1)})|$. For the case where each set of edges of the form
 1085 $\{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}$ for $b_i \in [2], i \in [3]$ is present, we have $|f_2^{-1}(E^{(1)})| = 8$. We count
 1086 the number of 3-matchings from the set $f_2^{-1}(E_S^{(1)})$.

1087 We do a case analysis based on the subgraph $S^{(1)}$ induced by $E_S^{(1)}$.

23:30 Bag PDB Queries

1088 ■ 3-matching (⚡⚡⚡)

1089 When $S^{(1)}$ is isomorphic to ⚡⚡⚡, it is the case that edges in $E_S^{(2)}$ are *not* disjoint only for the
 1090 pairs $(e_i, 0), (e_i, 1)$ for $i \in \{1, 2, 3\}$. By definition, each set of edges in $f_2^{-1}(E_S^{(1)})$ is a three
 1091 matching and $|f_2^{-1}(E_S^{(1)})| = 8$ possible 3-matchings.

1092 ■ Disjoint Two-Path (⚡ ⚡)

1093 For $S^{(1)}$ isomorphic to ⚡ ⚡ edges e_2, e_3 form a 2-path with e_1 being disjoint. This means
 1094 that in $S^{(2)}$ edges $(e_2, 0), (e_2, 1), (e_3, 0), (e_3, 1)$ form a 4-path while $(e_1, 0), (e_1, 1)$ is its own
 1095 disjoint 2-path. We can pick either $(e_1, 0)$ or $(e_1, 1)$ for the first edge in the 3-matching, while
 1096 it is necessary to have a 2-matching from path $(e_2, 0), \dots, (e_3, 1)$. Note that the 4-path allows
 1097 for three possible 2-matchings, specifically,

$$1098 \quad \{(e_2, 0), (e_3, 0)\}, \{(e_2, 0), (e_3, 1)\}, \{(e_2, 1), (e_3, 1)\}.$$

1099 Since these two selections can be made independently, $|f_2^{-1}(E_S^{(1)})| = 2 \cdot 3 = 6$ *distinct*
 1100 3-matchings in $f_2^{-1}(E_S^{(1)})$.

1101 ■ 3-star (⚡)

1102 When $S^{(1)}$ is isomorphic to ⚡, the inner edges $(e_i, 1)$ of $S^{(2)}$ are all connected, and the
 1103 outer edges $(e_i, 0)$ are all disjoint. Note that for a valid 3-matching it must be the case that
 1104 at most one inner edge can be part of the set of disjoint edges. For the case of when exactly
 1105 one inner edge is chosen, there exist 3 possibilities, based on which inner edge is chosen.
 1106 Note that if $(e_i, 1)$ is chosen, the matching has to choose $(e_j, 0)$ for $j \neq i$ and $(e_{j'}, 0)$ for
 1107 $j' \neq i, j' \neq j$. The remaining possible 3-matching occurs when all 3 outer edges are chosen,
 1108 and $|f_2^{-1}(E_S^{(1)})| = 4$.

1109 ■ 3-path (⚡⚡)

1110 When $S^{(1)}$ is isomorphic to ⚡⚡ it is the case that all edges beginning with e_1 and ending with e_3
 1111 are successively connected. This means that the edges of $E_S^{(2)}$ form a 6-path. For a 3-matching
 1112 to exist in $f_2^{-1}(E_S^{(1)})$, we cannot pick both $(e_i, 0)$ and $(e_i, 1)$ or both $(e_i, 1)$ and $(e_j, 0)$ where
 1113 $j = i + 1$. There are four such possibilities: $\{(e_1, 0), (e_2, 0), (e_3, 0)\}, \{(e_1, 0), (e_2, 0), (e_3, 1)\},$
 1114 $\{(e_1, 0), (e_2, 1), (e_3, 1)\}, \{(e_1, 1), (e_2, 1), (e_3, 1)\}$ and $|f_2^{-1}(E_S^{(1)})| = 4$.

1115 ■ Triangle (⚡)

1116 For $S^{(1)}$ isomorphic to ⚡, note that it is the case that the edges in $E_S^{(2)}$ are connected in a
 1117 successive manner, but this time in a cycle, such that $(e_1, 0)$ and $(e_3, 1)$ are also connected.
 1118 While this is similar to the discussion of the three path above, the first and last edges are
 1119 not disjoint. This rules out both subsets of $(e_1, 0), (e_2, 0), (e_3, 1)$ and $(e_1, 0), (e_2, 1), (e_3, 1)$, so
 1120 that $|f_2^{-1}(E_S^{(1)})| = 2$.

1121 Let us now consider when $E_S^{(1)} \in \binom{E_1}{\leq 2}$, i.e. fixed subgraphs among

1122 ■ 2-matching (⚡⚡), 2-path (⚡), 1 edge (⚡)

1123 When $|E_S^{(1)}| = 2$, we can only pick one from each of two pairs, $\{(e_1, 0), (e_1, 1)\}$ and
 1124 $\{(e_2, 0), (e_2, 1)\}$. The third edge choice in $E_S^{(2)}$ will break the disjoint property of a 3-
 1125 matching. Thus, a 3-matching cannot exist in $f_2^{-1}(E_S^{(1)})$. A similar argument holds for
 1126 $|E_S^{(1)}| = 1$, where the output of f_2^{-1} is $\{\emptyset\}$ since there are not enough edges in the input to
 1127 produce any other output.

1128 Observe that all of the arguments above focused solely on the property of subgraph $S^{(1)}$
 1129 being isomorphic. In other words, all $E_S^{(1)}$ of a given “shape” yield the same number of
 1130 3-matchings in $f_2^{-1}(E_S^{(1)})$, and this is why we get the required identity using the above case
 1131 analysis. ◀

1132 C.9.2 Proof of Lemma C.7

1133 **Proof.** The number of triangles in $G^{(\ell)}$ for $\ell \geq 2$ will always be 0 for the simple fact that all
 1134 cycles in $G^{(\ell)}$ will have at least six edges. ◀

1135 C.9.3 Proof of Lemma C.8

1136 **Proof.** The proof consists of two parts. First we need to show that a vector \mathbf{b} satisfying the
 1137 linear system exists and further can be computed in $O(m)$ time. Second we need to show
 1138 that $\#(G, \mathfrak{A}), \#(G, \mathfrak{B})$ can indeed be computed in time $O(1)$.

1139 The lemma claims that for $\mathbf{M} = \begin{pmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{pmatrix}$, $\mathbf{x} = \begin{pmatrix} \#(G, \mathfrak{A}) \\ \#(G, \mathfrak{B}) \end{pmatrix}$
 1140 satisfies the linear system $\mathbf{M} \cdot \mathbf{x} = \mathbf{b}$.

1141 To prove the first step, we use Lemma C.3 to derive the following equality (dropping the
 1142 superscript and referring to $G^{(1)}$ as G):

$$1143 \begin{aligned} & \#(G, \mathfrak{I})p^2 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{B})p^4 + 6\#(G, \mathfrak{C})p^3 + 6\#(G, \mathfrak{D})p^4 \\ 1144 & + 6\#(G, \mathfrak{E})p^4 + 6\#(G, \mathfrak{F})p^5 + 6\#(G, \mathfrak{G})p^6 = \tilde{\Phi}_G^3(p, \dots, p) \end{aligned} \quad (22)$$

$$1145 \begin{aligned} & \#(G, \mathfrak{A}) + \#(G, \mathfrak{B})p + \#(G, \mathfrak{C})p^2 + \#(G, \mathfrak{D})p^3 \\ 1146 & = \frac{\tilde{\Phi}_G^3(p, \dots, p)}{6p^3} - \frac{\#(G, \mathfrak{I})}{6p} - \#(G, \mathfrak{A}) - \#(G, \mathfrak{B})p - \#(G, \mathfrak{D})p \end{aligned} \quad (23)$$

$$1147 \begin{aligned} & \#(G, \mathfrak{A})(1 - 3p) - \#(G, \mathfrak{D})(3p^2 - p^3) = \\ 1148 & \frac{\tilde{\Phi}_G^3(p, \dots, p)}{6p^3} - \frac{\#(G, \mathfrak{I})}{6p} - \#(G, \mathfrak{A}) - \#(G, \mathfrak{B})p - \#(G, \mathfrak{D})p \\ 1149 & - [\#(G, \mathfrak{E})p + 3\#(G, \mathfrak{C})p] - [\#(G, \mathfrak{F})p^2 + 3\#(G, \mathfrak{G})p^2] \end{aligned} \quad (24)$$

1151 Eq. (22) is the result of Lemma C.3. We obtain the remaining equations through standard
 1152 algebraic manipulations.

1153 Note that the LHS of Eq. (24) is obtained using eq. (19) and eq. (20) and is indeed the
 1154 product $\mathbf{M}[1] \cdot \mathbf{x}[1]$. Further note that this product is equal to the RHS of Eq. (24), where
 1155 every term is computable in $O(m)$ time (by equations (15)-(20)). We set $\mathbf{b}[1]$ to the RHS of
 1156 Eq. (24).

1157 We follow the same process in deriving an equality for $G^{(2)}$. Replacing occurrences of
 1158 G with $G^{(2)}$, we obtain an equation (below) of the form of eq. (24) for $G^{(2)}$. Substituting

1159 identities from lemma C.6 and Lemma C.7 we obtain

$$\begin{aligned}
1160 \quad & 0 - (8\#(G, \mathfrak{iii}) + 6\#(G, \mathfrak{i} \mathfrak{A}) + 4\#(G, \mathfrak{AA}) + 4\#(G, \mathfrak{ii}) + 2\#(G, \mathfrak{A})) (3p^2 - p^3) = \\
1161 \quad & \frac{\tilde{\Phi}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{i})}{6p} - \#(G^{(2)}, \mathfrak{A}) - \#(G^{(2)}, \mathfrak{ii})p - \#(G^{(2)}, \mathfrak{AA})p \\
1162 \quad & - \left[\#(G^{(2)}, \mathfrak{i} \mathfrak{A})p^2 + 3\#(G^{(2)}, \mathfrak{iii})p^2 \right] - \left[\#(G^{(2)}, \mathfrak{ii})p + 3\#(G^{(2)}, \mathfrak{A})p \right] \\
& \qquad \qquad \qquad (25)
\end{aligned}$$

$$\begin{aligned}
1163 \quad & (10\#(G, \mathfrak{A}) + 10G\mathfrak{iii})(3p^2 - p^3) = \\
1164 \quad & \frac{\tilde{\Phi}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{i})}{6p} - \#(G^{(2)}, \mathfrak{A}) - \#(G^{(2)}, \mathfrak{ii})p - \#(G^{(2)}, \mathfrak{AA})p \\
1165 \quad & - \left[\#(G^{(2)}, \mathfrak{ii})p + 3\#(G^{(2)}, \mathfrak{A})p \right] - \left[\#(G^{(2)}, \mathfrak{i} \mathfrak{A})p^2 - 3\#(G^{(2)}, \mathfrak{iii})p^2 \right] \\
1166 \quad & + (4\#(G, \mathfrak{AA}) + [6\#(G, \mathfrak{i} \mathfrak{A}) + 18\#(G, \mathfrak{iii})] + [4\#(G, \mathfrak{ii}) + 12\#(G, \mathfrak{A})]) (3p^2 - p^3) \\
1167 \quad & \qquad \qquad \qquad (26)
\end{aligned}$$

1168 The steps to obtaining eq. (26) are analogous to the derivation immediately preceding. As in
1169 the previous derivation, note that the LHS of Eq. (26) is the same as $\mathbf{M}[2] \cdot \mathbf{x}[2]$. The RHS
1170 of Eq. (26) has terms all computable (by equations (15)-(20)) in $O(m)$ time. Setting $\mathbf{b}[2]$ to
1171 the RHS then completes the proof of step 1.

1172 Note that if \mathbf{M} has full rank then one can compute $\#(G, \mathfrak{A})$ and $\#(G, \mathfrak{iii})$ in $O(1)$
1173 using Gaussian elimination.

1174 To show that \mathbf{M} indeed has full rank, we show in what follows that $\text{Det}(\mathbf{M}) \neq 0$ for
1175 every $p \in (0, 1)$. $\text{Det}(\mathbf{M}) =$

$$\begin{aligned}
1176 \quad & \begin{vmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{vmatrix} = (1 - 3p) \cdot 10(3p^2 - p^3) + 10(3p^2 - p^3) \cdot (3p^2 - p^3) \\
1177 \quad & = 10(3p^2 - p^3) \cdot (1 - 3p + 3p^2 - p^3) = 10(3p^2 - p^3) \cdot (-p^3 + 3p^2 - 3p + 1) \\
1178 \quad & = 10p^2(3 - p) \cdot (1 - p)^3 \\
1179 \quad & \qquad \qquad \qquad (27)
\end{aligned}$$

1180 From Eq. (27) it can easily be seen that the roots of $\text{Det}(\mathbf{M})$ are 0, 1, and 3. Hence there
1181 are no roots in $(0, 1)$ and Lemma C.8 follows. \blacktriangleleft

1182 C.10 Proof of Theorem C.5

1183 **Proof.** We can compute $G^{(2)}$ from $G^{(1)}$ in $O(m)$ time. Additionally, if in time $O(T(m))$, we
1184 have $\tilde{\Phi}_{G^{(\ell)}}^3(p, \dots, p)$ for $\ell \in [2]$, then the theorem follows by Lemma C.8. \blacktriangleleft

1185 In other words, if Theorem C.5 holds, then so must Theorem 3.7.

1186 C.11 Proof of Theorem 3.7

1187 **Proof.** For the sake of contradiction, assume that for any G , we can compute $\tilde{\Phi}_G^3(p, \dots, p)$ in
1188 $o(m^{1+\epsilon_0})$ time. Let G be the input graph. Then by Theorem C.5 we can compute $\#(G, \mathfrak{A})$
1189 in further time $o(m^{1+\epsilon_0}) + O(m)$. Thus, the overall, reduction takes $o(m^{1+\epsilon_0}) + O(m) =$
1190 $o(m^{1+\epsilon_0})$ time, which violates Conjecture 3.3. \blacktriangleleft

1191 D Missing Details from Section 4

1192 In the following definitions and examples, we use the following polynomial as an example:

$$1193 \quad \Phi(X, Y) = 2X^2 + 3XY - 2Y^2. \qquad (28)$$

1194 ► **Definition D.1** (Pure Expansion). *The pure expansion of a polynomial Φ is formed by*
 1195 *computing all product of sums occurring in Φ , without combining like monomials. The pure*
 1196 *expansion of Φ generalizes Definition 2.1 by allowing monomials $m_i = m_j$ for $i \neq j$.*

1197 Note that similar in spirit to ??, $\mathbf{E}(\mathcal{C})$ Definition 4.1 reduces all variable exponents $e > 1$ to
 1198 $e = 1$. Further, it is true that $\mathbf{E}(\mathcal{C})$ is the pure expansion of \mathcal{C} .

1199 ► **Example D.2** (Example of Pure Expansion). *Consider the factorized representation $(X +$
 1200 $2Y)(2X - Y)$ of the polynomial in Eq. (28). Its circuit \mathcal{C} is illustrated in Fig. 3. The*
 1201 *pure expansion of the product is $2X^2 - XY + 4XY - 2Y^2$. As an additional example of*
 1202 *Definition 4.1, $\mathbf{E}(\mathcal{C}) = [(X, 2), (XY, -1), (XY, 4), (Y, -2)]$.*

1203 $\mathbf{E}(\mathcal{C})$ effectively¹³ encodes the *reduced* form of $\text{POLY}(\mathcal{C})$, decoupling each monomial into a
 1204 set of variables \mathbf{v} and a real coefficient c . However, unlike the constraint on the input Φ to
 1205 compute $\tilde{\Phi}$, the input circuit \mathcal{C} does not need to be in SMB/SOP form.

1206 ► **Example D.3** (Example for Definition 4.2). *Using the same factorization from Example D.2,*
 1207 *$\text{POLY}(|\mathcal{C}|) = (X + 2Y)(2X + Y) = 2X^2 + XY + 4XY + 2Y^2 = 2X^2 + 5XY + 2Y^2$. Note that*
 1208 *this is not the same as the polynomial from Eq. (28). As an example of the slight abuse of*
 1209 *notation we alluded to, $\text{POLY}(|\mathcal{C}|(1, \dots, 1)) = 2(1)^2 + 5(1)(1) + 2(1)^2 = 9$.*

1210 **Aaron says:** Verify whether we need pure expansion or not.

1211 ► **Definition D.4** (Subcircuit). *A subcircuit of a circuit \mathcal{C} is a circuit \mathcal{S} such that \mathcal{S} is a DAG*
 1212 *subgraph of the DAG representing \mathcal{C} . The sink of \mathcal{S} has exactly one gate g .*

1213 The following results assume input circuit \mathcal{C} computed from an arbitrary \mathcal{RA}^+ query Q
 1214 and arbitrary BIDB \mathcal{D} . We refer to \mathcal{C} as a BIDB circuit.

1215 **Aaron says:** Verify that the proof for Theorem D.5 doesn't rely on properties of \mathcal{RA}^+
 or BIDB.

1216 ► **Theorem D.5.** *Let \mathcal{C} be an arbitrary BIDB circuit and define $\Phi(\mathbf{X}) = \text{POLY}(\mathcal{C})$ and let*
 1217 *$k = \text{DEG}(\mathcal{C})$. Then an estimate \mathcal{E} of $\tilde{\Phi}(p_1, \dots, p_n)$ can be computed in time*

$$1218 \quad O\left(\left(\text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta} \cdot |\mathcal{C}|^2(1, \dots, 1) \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C})}{(\epsilon)^2 \cdot \tilde{\Phi}^2(p_1, \dots, p_n)}\right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right)$$

1219 *such that*

$$1220 \quad \Pr\left(\left|\mathcal{E} - \tilde{\Phi}(p_1, \dots, p_n)\right| > \epsilon \cdot \tilde{\Phi}(p_1, \dots, p_n)\right) \leq \delta. \quad (29)$$

1221 **Atri says:** Aaron: Just copied over from S4. The above text might need smoothening
 into the appendix.

1222 The slight abuse of notation seen in $|\mathcal{C}|(1, \dots, 1)$ is explained after Definition 4.2 and
 1223 an example is given in Example D.3. The only difference in the use of this notation in
 1224 Theorem D.5 is that we include an additional exponent to square the quantity.

1225 D.1 Proof of Theorem D.5

1226 We prove Theorem D.5 constructively by presenting an algorithm $\text{APPROXIMATE}\tilde{\Phi}$ (Algorithm 1)
 1227 which has the desired runtime and computes an approximation with the desired approximation

¹³The minor difference here is that $\mathbf{E}(\mathcal{C})$ encodes the *reduced* form over the SOP pure expansion of the compressed representation, as opposed to the SMB representation

Algorithm 1 APPROXIMATE $\tilde{\Phi}(\mathcal{C}, \mathbf{p}, \delta, \epsilon)$

Input: \mathcal{C} : Circuit**Input:** $\mathbf{p} = (p_1, \dots, p_n) \in [0, 1]^N$ **Input:** $\delta \in [0, 1]$ **Input:** $\epsilon \in [0, 1]$ **Output:** $\text{acc} \in \mathbb{R}$

```

1:  $\text{acc} \leftarrow 0$ 
2:  $N \leftarrow \left\lceil \frac{2 \log \frac{2}{\delta}}{\epsilon^2} \right\rceil$ 
3:  $(\mathcal{C}_{\text{mod}}, \text{size}) \leftarrow \text{ONEPASS}(\mathcal{C})$  ▷ ONEPASS is Algorithm 2
4: for  $i \in 1$  to  $N$  do ▷ Perform the required number of samples
5:    $(\mathbf{M}, \text{sgn}_i) \leftarrow \text{SAMPLEMONOMIAL}(\mathcal{C}_{\text{mod}})$  ▷ SAMPLEMONOMIAL is Algorithm 3. Note
   that  $\text{sgn}_i$  is the sign of the monomial's coefficient and not the coefficient itself
6:   if  $\mathbf{M}$  has at most one variable from each block then
7:      $\mathbf{Y}_i \leftarrow \prod_{X_j \in \mathbf{M}} p_j$  ▷  $\mathbf{M}$  is the sampled monomial's set of variables (cref. appendix D.8)
8:      $\mathbf{Y}_i \leftarrow \mathbf{Y}_i \times \text{sgn}_i$ 
9:      $\text{acc} \leftarrow \text{acc} + \mathbf{Y}_i$  ▷ Store the sum over all samples
10:  end if
11: end for
12:  $\text{acc} \leftarrow \text{acc} \times \frac{\text{size}}{N}$ 
13: return  $\text{acc}$ 

```

1228 guarantee. Algorithm APPROXIMATE $\tilde{\Phi}$ uses Algorithm ONEPASS to compute weights on the
1229 edges of a circuits. These weights are then used to sample a set of monomials of $\Phi(\mathcal{C})$ from
1230 the circuit \mathcal{C} by traversing the circuit using the weights to ensure that monomials are sampled
1231 with an appropriate probability. The correctness of APPROXIMATE $\tilde{\Phi}$ relies on the correctness
1232 (and runtime behavior) of auxiliary algorithms ONEPASS and SAMPLEMONOMIAL that we
1233 state in the following lemmas (and prove later in this part of the appendix).

► **Lemma D.6.** *The ONEPASS function completes in time:*

$$O(\text{SIZE}(\mathcal{C}) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}(1, \dots, 1)|), \log \text{SIZE}(\mathcal{C})))$$

1234 ONEPASS guarantees two post-conditions: First, for each subcircuit \mathcal{S} of \mathcal{C} , we have that
1235 $\mathcal{S}.\text{partial}$ is set to $|\mathcal{S}|(1, \dots, 1)$. Second, when $\mathcal{S}.\text{type} = +$, $\mathcal{S}.\text{Lweight} = \frac{|\mathcal{S}|(1, \dots, 1)}{|\mathcal{S}(1, \dots, 1)}$ and
1236 likewise for $\mathcal{S}.\text{Rweight}$.

1237 To prove correctness of Algorithm 1, we only use the following fact that follows from the above
1238 lemma: for the modified circuit $(\mathcal{C}_{\text{mod}})$ output by ONEPASS, $\mathcal{C}_{\text{mod}}.\text{partial} = |\mathcal{C}|(1, \dots, 1)$.

► **Lemma D.7.** *The function SAMPLEMONOMIAL completes in time*

$$O(\log k \cdot k \cdot \text{DEPTH}(\mathcal{C}) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log \text{SIZE}(\mathcal{C})))$$

1239 where $k = \text{DEG}(\mathcal{C})$. The function returns every $(\mathbf{v}, \text{sign}(\mathbf{c}))$ for $(\mathbf{v}, \mathbf{c}) \in E(\mathcal{C})$ with probability
1240 $\frac{|\mathbf{c}|}{|\mathcal{C}(1, \dots, 1)|}$.

1241 With the above two lemmas, we are ready to argue the following result:

1242 ► **Theorem D.8.** *For any \mathcal{C} with*

1243 **Aaron says:** *Pretty sure this is $\text{DEG}(|\mathcal{C}|)$. Have to read on to be sure.*

1244 $DEG(\text{poly}(|\mathcal{C}|)) = k$, algorithm 1 outputs an estimate acc of $\tilde{\Phi}(p_1, \dots, p_n)$ such that

$$1245 \Pr \left(\left| \text{acc} - \tilde{\Phi}(p_1, \dots, p_n) \right| > \epsilon \cdot |\mathcal{C}|(1, \dots, 1) \right) \leq \delta,$$

1246 in $O \left(\left(\text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta}}{\epsilon^2} \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C}) \right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log \text{SIZE}(\mathcal{C})) \right)$ time.

1247 Before proving Theorem D.8, we use it to argue the claimed runtime of our main result,
1248 Theorem D.5.

1249 **Proof of Theorem D.5.** Set $\mathcal{E} = \text{APPROXIMATE}\tilde{\Phi}(\mathcal{C}, (p_1, \dots, p_n), \delta, \epsilon')$, where

$$1250 \epsilon' = \epsilon \cdot \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)},$$

1251 which achieves the claimed error bound on \mathcal{E} (acc) trivially due to the assignment to ϵ' and
1252 theorem D.8, since $\epsilon' \cdot |\mathcal{C}|(1, \dots, 1) = \epsilon \cdot \frac{\tilde{\Phi}(1, \dots, 1)}{|\mathcal{C}|(1, \dots, 1)} \cdot |\mathcal{C}|(1, \dots, 1) = \epsilon \cdot \tilde{\Phi}(1, \dots, 1)$.

1253 The claim on the runtime follows from Theorem D.8 since

$$1254 \frac{1}{(\epsilon')^2} \cdot \log \left(\frac{1}{\delta} \right) = \frac{\log \frac{1}{\delta}}{\epsilon^2 \left(\frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)} \right)^2}$$

$$1255 = \frac{\log \frac{1}{\delta} \cdot |\mathcal{C}|^2(1, \dots, 1)}{\epsilon^2 \cdot \tilde{\Phi}^2(p_1, \dots, p_n)}.$$

1257 Let us now prove Theorem D.8: ◀

1259 D.2 Proof of Theorem D.8

1260 **Proof.** Consider now the random variables Y_1, \dots, Y_N , where each Y_i is the value of Y_i in
1261 algorithm 1 after line 8 is executed. Overloading $\text{ISIND}(\cdot)$ to receive monomial input (recall
1262 \mathbf{v}_m is the monomial composed of the variables in the set \mathbf{v}), we have

$$1263 Y_i = \mathbb{1}_{(\text{ISIND}(\mathbf{v}_m))} \cdot \prod_{X_i \in \text{VAR}(\mathbf{v})} p_i,$$

1264 where the indicator variable handles the check in Line 6 Then for random variable Y_i , it is
1265 the case that

$$1266 \mathbb{E}[Y_i] = \sum_{(\mathbf{v}, \mathbf{c}) \in \mathbb{E}(\mathcal{C})} \frac{\mathbb{1}_{(\text{ISIND}(\mathbf{v}_m))} \cdot \mathbf{c} \cdot \prod_{X_i \in \text{VAR}(\mathbf{v})} p_i}{|\mathcal{C}|(1, \dots, 1)}$$

$$1267 = \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)},$$

1269 where in the first equality we use the fact that $\text{sgn}_i \cdot |\mathbf{c}| = \mathbf{c}$ and the second equality follows
1270 from Eq. (2) with X_i substituted by p_i .

1271 Let $\bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_i$. It is also true that

$$1272 \mathbb{E}[\bar{Y}] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[Y_i] = \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)}.$$

1273 Hoeffding's inequality states that if we know that each Y_i (which are all independent)
1274 always lie in the intervals $[a_i, b_i]$, then it is true that

$$1275 \Pr(|\bar{Y} - \mathbb{E}[\bar{Y}]| \geq \epsilon) \leq 2 \exp\left(-\frac{2N^2\epsilon^2}{\sum_{i=1}^N (b_i - a_i)^2}\right).$$

1276 Line 5 shows that sgn_i has a value in $\{-1, 1\}$ that is multiplied with $O(k)$ $p_i \in [0, 1]$,
1277 which implies the range for each Y_i is $[-1, 1]$. Using Hoeffding's inequality, we then get:

$$1278 \Pr(|\bar{Y} - \mathbb{E}[\bar{Y}]| \geq \epsilon) \leq 2 \exp\left(-\frac{2N^2\epsilon^2}{2^2N}\right) = 2 \exp\left(-\frac{N\epsilon^2}{2}\right) \leq \delta,$$

1279 where the last inequality dictates our choice of N in Line 2.

1280 For the claimed probability bound of $\Pr\left(|\text{acc} - \tilde{\Phi}(p_1, \dots, p_n)| > \epsilon \cdot |\mathcal{C}|(1, \dots, 1)\right) \leq \delta$,
1281 note that in the algorithm, acc is exactly $\bar{Y} \cdot |\mathcal{C}|(1, \dots, 1)$. Multiplying the rest of the terms
1282 by the additional factor $|\mathcal{C}|(1, \dots, 1)$ yields the said bound.

1283 This concludes the proof for the first claim of theorem D.8. Next, we prove the claim on
1284 the runtime.

1285 Run-time Analysis

1286 The runtime of the algorithm is dominated first by Line 3 (which by Lemma D.6 takes time
1287 $O(\text{SIZE}(\mathcal{C}) \cdot \bar{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C}))))$) and then by N iterations of the loop in
1288 Line 4. Each iteration's run time is dominated by the call to `SAMPLEMONOMIAL` in Line 5
1289 (which by Lemma D.7 takes $O(\log k \cdot k \cdot \text{DEPTH}(\mathcal{C}) \cdot \bar{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C}))))$)
1290 and the check Line 6, which by the subsequent argument takes $O(k \log k)$ time. We sort
1291 the $O(k)$ variables by their block IDs and then check if there is a duplicate block ID or not.
1292 Combining all the times discussed here gives us the desired overall runtime. ◀

1293 D.3 Proof of Theorem 4.7

1294 **Proof.** The result follows by first noting that by definition of γ , we have

$$1295 \tilde{\Phi}(1, \dots, 1) = (1 - \gamma) \cdot |\mathcal{C}|(1, \dots, 1).$$

1296 Further, since each $p_i \geq p_0$ and $\Phi(\mathbf{X})$ (and hence $\tilde{\Phi}(\mathbf{X})$) has degree at most k , we have that

$$1297 \tilde{\Phi}(1, \dots, 1) \geq p_0^k \cdot \tilde{\Phi}(1, \dots, 1).$$

1298 The above two inequalities implies $\tilde{\Phi}(1, \dots, 1) \geq p_0^k \cdot (1 - \gamma) \cdot |\mathcal{C}|(1, \dots, 1)$. Applying this
1299 bound in the runtime bound in Theorem D.5 gives the first claimed runtime. The final
1300 runtime of $O_k\left(\frac{1}{\epsilon^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta} \cdot \bar{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right)$ follows by noting
1301 that $\text{DEPTH}(\mathcal{C}) \leq \text{SIZE}(\mathcal{C})$ and absorbing all factors that just depend on k . ◀

1302 D.4 Proof of Lemma 4.8

1303 We will prove Lemma 4.8 by considering the two cases separately. We start by considering
1304 the case when \mathcal{C} is a tree:

1305 ▶ **Lemma D.9.** *Let \mathcal{C} be a tree (i.e. the sub-circuits corresponding to two children of a node
1306 in \mathcal{C} are completely disjoint). Then we have*

$$1307 |\mathcal{C}|(1, \dots, 1) \leq (\text{SIZE}(\mathcal{C}))^{\text{DEG}(\mathcal{C})+1}.$$

1308 **Proof of Lemma D.9.** For notational simplicity define $N = \text{SIZE}(\mathcal{C})$ and $k = \text{DEG}(\mathcal{C})$. We
 1309 use induction on $\text{DEPTH}(\mathcal{C})$ to show that $|\mathcal{C}|(1, \dots, 1) \leq N^{k+1}$. For the base case, we have
 1310 that $\text{DEPTH}(\mathcal{C}) = 0$, and there can only be one node which must contain a coefficient or
 1311 constant. In this case, $|\mathcal{C}|(1, \dots, 1) = 1$, and $\text{SIZE}(\mathcal{C}) = 1$, and by Definition 4.4 it is the
 1312 case that $0 \leq k = \text{DEG}(\mathcal{C}) \leq 1$, and it is true that $|\mathcal{C}|(1, \dots, 1) = 1 \leq N^{k+1} = 1^{k+1} = 1$ for
 1313 $k \in \{0, 1\}$.

1314 Assume for $\ell > 0$ an arbitrary circuit \mathcal{C} of $\text{DEPTH}(\mathcal{C}) \leq \ell$ that it is true that $|\mathcal{C}|(1, \dots, 1) \leq$
 1315 N^{k+1} .

1316 For the inductive step we consider a circuit \mathcal{C} such that $\text{DEPTH}(\mathcal{C}) = \ell + 1$. The sink can
 1317 only be either a \times or $+$ gate. Let k_L, k_R denote $\text{DEG}(\mathcal{C}_L)$ and $\text{DEG}(\mathcal{C}_R)$ respectively. Consider
 1318 when sink node is \times . Then note that

$$\begin{aligned} 1319 \quad |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) \cdot |\mathcal{C}_R|(1, \dots, 1) \\ 1320 &\leq (N-1)^{k_L+1} \cdot (N-1)^{k_R+1} \\ 1321 &= (N-1)^{k+1} \\ 1322 &\leq N^{k+1}. \end{aligned} \tag{30}$$

1324 In the above the first inequality follows from the inductive hypothesis (and the fact that the
 1325 size of either subtree is at most $N-1$) and Eq. (30) follows by definition 4.4 which states
 1326 that for $k = \text{DEG}(\mathcal{C})$ we have $k = k_L + k_R + 1$.

1327 For the case when the sink gate is a $+$ gate, then for $N_L = \text{SIZE}(\mathcal{C}_L)$ and $N_R = \text{SIZE}(\mathcal{C}_R)$
 1328 we have

$$\begin{aligned} 1329 \quad |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) + |\mathcal{C}_R|(1, \dots, 1) \\ 1330 &\leq N_L^{k+1} + N_R^{k+1} \\ 1331 &\leq (N-1)^{k+1} \\ 1332 &\leq N^{k+1}. \end{aligned} \tag{31}$$

1334 In the above, the first inequality follows from the inductive hypothesis and definition 4.4
 1335 (which implies the fact that $k_L, k_R \leq k$). Note that the RHS of this inequality is maximized
 1336 when the base and exponent of one of the terms is maximized. The second inequality follows
 1337 from this fact as well as the fact that since \mathcal{C} is a tree we have $N_L + N_R = N-1$ and, lastly,
 1338 the fact that $k \geq 0$. This completes the proof.

1339 The upper bound in Lemma 4.8 for the general case is a simple variant of the above proof
 1340 (but we present a proof sketch of the bound below for completeness):

1341 **► Lemma D.10.** *Let \mathcal{C} be a (general) circuit. Then we have*

$$1342 \quad |\mathcal{C}|(1, \dots, 1) \leq 2^{2^{\text{DEG}(\mathcal{C})} \cdot \text{DEPTH}(\mathcal{C})}.$$

1343 **Proof Sketch of Lemma D.10.** We use the same notation as in the proof of Lemma D.9 and
 1344 further define $d = \text{DEPTH}(\mathcal{C})$. We will prove by induction on $\text{DEPTH}(\mathcal{C})$ that $|\mathcal{C}|(1, \dots, 1) \leq$
 1345 $2^{2^k \cdot d}$. The base case argument is similar to that in the proof of Lemma D.9. In the inductive
 1346 case we have that $d_L, d_R \leq d-1$.

1347 For the case when the sink node is \times , we get that

$$\begin{aligned} 1348 \quad |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) \times |\mathcal{C}_R|(1, \dots, 1) \\ 1349 &\leq 2^{2^{k_L} \cdot d_L} \times 2^{2^{k_R} \cdot d_R} \\ 1350 &\leq 2^{2 \cdot 2^{k-1} \cdot (d-1)} \end{aligned}$$

$$\leq 2^{2^k d}.$$

In the above the first inequality follows from inductive hypothesis while the second inequality follows from the fact that $k_L, k_R \leq k - 1$ and $d_L, d_R \leq d - 1$, where we substitute the upperbound into every respective term.

Now consider the case when the sink node is $+$, we get that

$$\begin{aligned} |\mathbf{C}|(1, \dots, 1) &= |\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1) \\ &\leq 2^{2^{k_L} \cdot d_L} + 2^{2^{k_R} \cdot d_R} \\ &\leq 2 \cdot 2^{2^k (d-1)} \\ &\leq 2^{2^k d}. \end{aligned}$$

In the above the first inequality follows from the inductive hypothesis while the second inequality follows from the facts that $k_L, k_R \leq k$ and $d_L, d_R \leq d - 1$. The final inequality follows from the fact that $k \geq 0$. \blacktriangleleft

D.5 OnePass Remarks

Please note that it is *assumed* that the original call to ONEPASS consists of a call on an input circuit \mathbf{C} such that the values of members `partial`, `Lweight` and `Rweight` have been initialized to Null across all gates.

The evaluation of $|\mathbf{C}|(1, \dots, 1)$ can be defined recursively, as follows (where \mathbf{C}_L and \mathbf{C}_R are the ‘left’ and ‘right’ inputs of \mathbf{C} if they exist):

$$|\mathbf{C}|(1, \dots, 1) = \begin{cases} |\mathbf{C}_L|(1, \dots, 1) \cdot |\mathbf{C}_R|(1, \dots, 1) & \text{if } \mathbf{C.type} = \times \\ |\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1) & \text{if } \mathbf{C.type} = + \\ |\mathbf{C.val}| & \text{if } \mathbf{C.type} = \text{NUM} \\ 1 & \text{if } \mathbf{C.type} = \text{VAR}. \end{cases} \quad (32)$$

It turns out that for proof of Lemma D.7, we need to argue that when $\mathbf{C.type} = +$, we indeed have

$$\mathbf{C.Lweight} \leftarrow \frac{|\mathbf{C}_L|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)}; \quad (33)$$

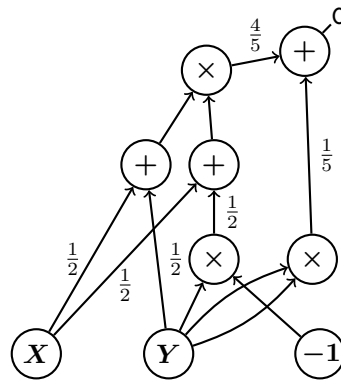
$$\mathbf{C.Rweight} \leftarrow \frac{|\mathbf{C}_R|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)} \quad (34)$$

D.6 OnePass Example

► **Example D.11.** Let T encode the expression $(X + Y)(X - Y) + Y^2$. After one pass, Algorithm 2 would have computed the following weight distribution. For the two inputs of the sink gate \mathbf{C} , $\mathbf{C.Lweight} = \frac{4}{5}$ and $\mathbf{C.Rweight} = \frac{1}{5}$. Similarly, for \mathbf{S} denoting the left input of \mathbf{C}_L , $\mathbf{S.Lweight} = \mathbf{S.Rweight} = \frac{1}{2}$. This is depicted in Fig. 4.

D.7 Proof of OnePass (Lemma D.6)

Proof. We prove the correct computation of `partial`, `Lweight`, `Rweight` values on \mathbf{C} by induction over the number of iterations in the topological order `TOPORD` (line 1) of the input circuit \mathbf{C} . `TOPORD` follows the standard definition of a topological ordering over the DAG structure of \mathbf{C} .



■ **Figure 4** Weights computed by ONEPASS in Example D.11.

■ **Algorithm 2** ONEPASS (C)

Input: C : Circuit

Output: C : Annotated Circuit

Output: $\text{sum} \in \mathbb{N}$

```

1: for  $g$  in TOPORD ( $C$ ) do                                     ▷ TOPORD ( $\cdot$ ) is the topological order of  $C$ 
2:   if  $g.\text{type} = \text{VAR}$  then
3:      $g.\text{partial} \leftarrow 1$ 
4:   else if  $g.\text{type} = \text{NUM}$  then
5:      $g.\text{partial} \leftarrow |g.\text{val}|$ 
6:   else if  $g.\text{type} = \times$  then
7:      $g.\text{partial} \leftarrow g_L.\text{partial} \times g_R.\text{partial}$ 
8:   else
9:      $g.\text{partial} \leftarrow g_L.\text{partial} + g_R.\text{partial}$ 
10:     $g.L\text{weight} \leftarrow \frac{g_L.\text{partial}}{g.\text{partial}}$ 
11:     $g.R\text{weight} \leftarrow \frac{g_R.\text{partial}}{g.\text{partial}}$ 
12:  end if
13:   $\text{sum} \leftarrow g.\text{partial}$ 
14: end for
15: return ( $\text{sum}, C$ )

```

1388 For the base case, we have only one gate, which by definition is a source gate and must be
 1389 either VAR or NUM. In this case, as per eq. (32), lines 3 and 5 correctly compute $C.\text{partial}$
 1390 as 1.

1391 For the inductive hypothesis, assume that ONEPASS correctly computes $S.\text{partial}$,
 1392 $S.L\text{weight}$, and $S.R\text{weight}$ for all gates g in C with $k \geq 0$ iterations over TOPORD.

1393 **Aaron says:** Notes above: Algo uses Reduce, but we don't use that anymore. The figure needs to change to a circuit.

1394 We now prove for $k+1$ iterations that ONEPASS correctly computes the partial , $L\text{weight}$,
 1395 and $R\text{weight}$ values for each gate g_i in C for $i \in [k+1]$. The g_{k+1} must be in the last
 1396 ordering of all gates g_i . When $\text{SIZE}(C) > 1$, if g_{k+1} is a leaf node, we are back to the base
 1397 case. Otherwise g_{k+1} is an internal node which requires binary input.

1398 When $g_{k+1}.\text{type} = +$, then by line 9 $g_{k+1}.\text{partial} = g_{k+1_L}.\text{partial} + g_{k+1_R}.\text{partial}$,
 1399 a correct computation, as per eq. (32). Further, lines 10 and 11 compute $g_{k+1}.L\text{weight} =$

Algorithm 3 SAMPLEMONOMIAL (\mathcal{C})

Input: \mathcal{C} : Circuit
Output: vars : TreeSet
Output: $\text{sgn} \in \{-1, 1\}$ ▷ Algorithm 2 should have been run before this one

```

1:  $\text{vars} \leftarrow \emptyset$ 
2: if  $\mathcal{C}.\text{type} = +$  then ▷ Sample at every + node
3:    $\mathcal{C}_{\text{samp}} \leftarrow$  Sample from left input ( $\mathcal{C}_L$ ) and right input ( $\mathcal{C}_R$ ) w.p.  $\mathcal{C}.\text{Lweight}$  and  $\mathcal{C}.\text{Rweight}$ . ▷ Each call to SAMPLEMONOMIAL uses fresh randomness
4:    $(v, s) \leftarrow$  SAMPLEMONOMIAL( $\mathcal{C}_{\text{samp}}$ )
5:   return  $(v, s)$ 
6: else if  $\mathcal{C}.\text{type} = \times$  then ▷ Multiply the sampled values of all inputs
7:    $\text{sgn} \leftarrow 1$ 
8:   for  $\text{input}$  in  $\mathcal{C}.\text{input}$  do
9:      $(v, s) \leftarrow$  SAMPLEMONOMIAL( $\text{input}$ )
10:     $\text{vars} \leftarrow \text{vars} \cup \{v\}$ 
11:     $\text{sgn} \leftarrow \text{sgn} \times s$ 
12:   end for
13:   return  $(\text{vars}, \text{sgn})$ 
14: else if  $\mathcal{C}.\text{type} = \text{NUM}$  then ▷ The leaf is a coefficient
15:   return  $(\{\}, \text{SGN}(\mathcal{C}.\text{val}))$  ▷  $\text{SGN}(\cdot)$  outputs  $-1$  for  $\mathcal{C}.\text{val} \geq 1$  and  $-1$  for  $\mathcal{C}.\text{val} \leq -1$ 
16: else if  $\mathcal{C}.\text{type} = \text{var}$  then
17:   return  $(\{\mathcal{C}.\text{val}\}, 1)$ 
18: end if

```

1400 $\frac{\mathbf{g}_{k+1_L}.\text{partial}}{\mathbf{g}_{k+1}.\text{partial}}$ and analogously for $\mathbf{g}_{k+1}.\text{Rweight}$. All values needed for each computation have
1401 been correctly computed by the inductive hypothesis.

1402 When $\mathbf{g}_{k+1}.\text{type} = \times$, then line 7 computes $\mathbf{g}_{k+1}.\text{partial} = \mathbf{g}_{k+1_L}.\text{partial} \times \mathbf{g}_{k+1_R}.\text{partial}$,
1403 which indeed by eq. (32) is correct. This concludes the proof of correctness.

1404 Runtime Analysis

1405 It is known that $\text{TOPORD}(G)$ is computable in linear time. There are $\text{SIZE}(\mathcal{C})$ iterations, each
1406 of which takes $O(\overline{\mathcal{M}}(\log(|\mathcal{C}(1 \dots, 1)|), \log(\text{SIZE}(\mathcal{C}))))$ time. This can be seen since each of
1407 all the numbers which the algorithm computes is at most $|\mathcal{C}|(1, \dots, 1)$. Hence, by definition
1408 each such operation takes $\overline{\mathcal{M}}(\log(|\mathcal{C}(1 \dots, 1)|), \log \text{SIZE}(\mathcal{C}))$ time, which proves the claimed
1409 runtime. ◀

1410 D.8 SampleMonomial Remarks

1411 We briefly describe the top-down traversal of SAMPLEMONOMIAL. When $\mathcal{C}.\text{type} = +$, the
1412 input to be visited is sampled from the weighted distribution precomputed by ONEPASS.
1413 When a $\mathcal{C}.\text{type} = \times$ node is visited, both inputs are visited. The algorithm computes two
1414 properties: the set of all variable leaf nodes visited, and the product of the signs of visited
1415 coefficient leaf nodes. We will assume the TreeSet data structure to maintain sets with
1416 logarithmic time insertion and linear time traversal of its elements. While we would like to
1417 take advantage of the space efficiency gained in using a circuit \mathcal{C} instead an expression tree T ,
1418 we do not know that such a method exists when computing a sample of the input polynomial
1419 representation.

1420 The efficiency gains of circuits over trees is found in the capability of circuits to only
 1421 require space for each *distinct* term in the compressed representation. This saves space
 1422 in such polynomials containing non-distinct terms multiplied or added to each other, e.g.,
 1423 x^4 . However, to avoid biased sampling, it is imperative to sample from both inputs of a
 1424 multiplication gate, independently, which is indeed the approach of SAMPLEMONOMIAL.

1425 D.9 Proof of SampleMonomial (Lemma D.7)

1426 **Proof.** We first need to show that SAMPLEMONOMIAL samples a valid monomial \mathbf{v}_m by
 1427 sampling and returning a set of variables \mathbf{v} , such that (\mathbf{v}, \mathbf{c}) is in $\mathbf{E}(\mathbf{C})$ and \mathbf{v}_m is indeed a
 1428 monomial of the $\tilde{\Phi}(\mathbf{X})$ encoded in \mathbf{C} . We show this via induction over the depth of \mathbf{C} .

1429 For the base case, let the depth d of \mathbf{C} be 0. We have that the single gate is either a
 1430 constant c for which by line 15 we return $\{\}$, or we have that $\mathbf{C.type} = \text{VAR}$ and $\mathbf{C.val} = x$,
 1431 and by line 17 we return $\{x\}$. By definition 4.1, both cases return a valid \mathbf{v} for some (\mathbf{v}, \mathbf{c})
 1432 from $\mathbf{E}(\mathbf{C})$, and the base case is proven.

1433 For the inductive hypothesis, assume that for $d \leq k$ for some $k \geq 0$, that it is indeed the
 1434 case that SAMPLEMONOMIAL returns a valid monomial.

1435 For the inductive step, let us take a circuit \mathbf{C} with $d = k + 1$. Note that each input
 1436 has depth $d - 1 \leq k$, and by inductive hypothesis both of them sample a valid monomial.
 1437 Then the sink can be either a $+$ or \times gate. For the case when $\mathbf{C.type} = +$, line 3 of
 1438 SAMPLEMONOMIAL will choose one of the inputs of the source. By inductive hypothesis it is
 1439 the case that some valid monomial is being randomly sampled from each of the inputs. Then
 1440 it follows when $\mathbf{C.type} = +$ that a valid monomial is sampled by SAMPLEMONOMIAL. When
 1441 the $\mathbf{C.type} = \times$, line 10 computes the set union of the monomials returned by the two inputs
 1442 of the sink, and it is trivial to see by definition 4.1 that \mathbf{v}_m is a valid monomial encoded by
 1443 some (\mathbf{v}, \mathbf{c}) of $\mathbf{E}(\mathbf{C})$.

1444 We will next prove by induction on the depth d of \mathbf{C} that for $(\mathbf{v}, \mathbf{c}) \in \mathbf{E}(\mathbf{C})$, \mathbf{v} is sampled
 1445 with a probability $\frac{|\mathbf{c}|}{|\mathbf{C}|(1, \dots, 1)}$.

1446 For the base case $d = 0$, by definition 2.6 we know that the $\text{SIZE}(\mathbf{C}) = 1$ and $\mathbf{C.type} =$
 1447 NUM or VAR. For either case, the probability of the value returned is 1 since there is only
 1448 one value to sample from. When $\mathbf{C.val} = x$, the algorithm always return the variable set
 1449 $\{x\}$. When $\mathbf{C.type} = \text{NUM}$, SAMPLEMONOMIAL will always return \emptyset .

1450 For the inductive hypothesis, assume that for $d \leq k$ and $k \geq 0$ SAMPLEMONOMIAL indeed
 1451 returns \mathbf{v} in (\mathbf{v}, \mathbf{c}) of $\mathbf{E}(\mathbf{C})$ with probability $\frac{|\mathbf{c}|}{|\mathbf{C}|(1, \dots, 1)}$.

1452 We prove now for $d = k + 1$ the inductive step holds. It is the case that the sink of \mathbf{C} has
 1453 two inputs \mathbf{C}_L and \mathbf{C}_R . Since \mathbf{C}_L and \mathbf{C}_R are both depth $d - 1 \leq k$, by inductive hypothesis,
 1454 SAMPLEMONOMIAL will return \mathbf{v}_L in $(\mathbf{v}_L, \mathbf{c}_L)$ of $\mathbf{E}(\mathbf{C}_L)$ and \mathbf{v}_R in $(\mathbf{v}_R, \mathbf{c}_R)$ of $\mathbf{E}(\mathbf{C}_R)$, from \mathbf{C}_L and
 1455 \mathbf{C}_R with probability $\frac{|\mathbf{c}_L|}{|\mathbf{C}_L|(1, \dots, 1)}$ and $\frac{|\mathbf{c}_R|}{|\mathbf{C}_R|(1, \dots, 1)}$.

1456 Consider the case when $\mathbf{C.type} = \times$. For the term (\mathbf{v}, \mathbf{c}) from $\mathbf{E}(\mathbf{C})$ that is being sampled
 1457 it is the case that $\mathbf{v} = \mathbf{v}_L \cup \mathbf{v}_R$, where \mathbf{v}_L is coming from \mathbf{C}_L and \mathbf{v}_R from \mathbf{C}_R . The probability
 1458 that SAMPLEMONOMIAL (\mathbf{C}_L) returns \mathbf{v}_L is $\frac{|\mathbf{c}_{v_L}|}{|\mathbf{C}_L|(1, \dots, 1)}$ and $\frac{|\mathbf{c}_{v_R}|}{|\mathbf{C}_R|(1, \dots, 1)}$ for \mathbf{v}_R . Since both \mathbf{v}_L
 1459 and \mathbf{v}_R are sampled with independent randomness, the final probability for sample \mathbf{v} is
 1460 then $\frac{|\mathbf{c}_{v_L}| \cdot |\mathbf{c}_{v_R}|}{|\mathbf{C}_L|(1, \dots, 1) \cdot |\mathbf{C}_R|(1, \dots, 1)}$. For (\mathbf{v}, \mathbf{c}) in $\mathbf{E}(\mathbf{C})$, by definition 4.1 it is indeed the case that
 1461 $|\mathbf{c}| = |\mathbf{c}_{v_L}| \cdot |\mathbf{c}_{v_R}|$ and that (as shown in eq. (32)) $|\mathbf{C}|(1, \dots, 1) = |\mathbf{C}_L|(1, \dots, 1) \cdot |\mathbf{C}_R|(1, \dots, 1)$,
 1462 and therefore \mathbf{v} is sampled with correct probability $\frac{|\mathbf{c}|}{|\mathbf{C}|(1, \dots, 1)}$.

1463 For the case when $\mathbf{C.type} = +$, SAMPLEMONOMIAL will sample \mathbf{v} from one of its inputs.
 1464 By inductive hypothesis we know that any \mathbf{v}_L in $\mathbf{E}(\mathbf{C}_L)$ and any \mathbf{v}_R in $\mathbf{E}(\mathbf{C}_R)$ will both be
 1465 sampled with correct probability $\frac{|\mathbf{c}_{v_L}|}{|\mathbf{C}_L|(1, \dots, 1)}$ and $\frac{|\mathbf{c}_{v_R}|}{|\mathbf{C}_R|(1, \dots, 1)}$, where either \mathbf{v}_L or \mathbf{v}_R will equal \mathbf{v} ,
 1466 depending on whether \mathbf{C}_L or \mathbf{C}_R is sampled. Assume that \mathbf{v} is sampled from \mathbf{C}_L , and note that

1467 a symmetric argument holds for the case when \mathbf{v} is sampled from \mathbf{C}_R . Notice also that the
 1468 probability of choosing \mathbf{C}_L from \mathbf{C} is $\frac{|\mathbf{C}_L|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)}$ as computed by ONEPASS. Then,
 1469 since SAMPLEMONOMIAL goes top-down, and each sampling choice is independent (which
 1470 follows from the randomness in the root of \mathbf{C} being independent from the randomness used
 1471 in its subtrees), the probability for \mathbf{v} to be sampled from \mathbf{C} is equal to the product of the
 1472 probability that \mathbf{C}_L is sampled from \mathbf{C} and \mathbf{v} is sampled in \mathbf{C}_L , and

$$\begin{aligned}
 1473 \quad & Pr(\text{SAMPLEMONOMIAL}(\mathbf{C}) = \mathbf{v}) = \\
 1474 \quad & Pr(\text{SAMPLEMONOMIAL}(\mathbf{C}_L) = \mathbf{v}) \cdot Pr(\text{SampledChild}(\mathbf{C}) = \mathbf{C}_L) \\
 1475 \quad & = \frac{|\mathbf{c}_v|}{|\mathbf{C}_L|(1, \dots, 1)} \cdot \frac{|\mathbf{C}_L|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)} \\
 1476 \quad & = \frac{|\mathbf{c}_v|}{|\mathbf{C}|(1, \dots, 1)}, \\
 1477
 \end{aligned}$$

1478 and we obtain the desired result.

1479 Lastly, we show by simple induction of the depth d of \mathbf{C} that SAMPLEMONOMIAL indeed
 1480 returns the correct sign value of \mathbf{c} in (\mathbf{v}, \mathbf{c}) .

1481 In the base case, $\mathbf{C.type} = \text{NUM}$ or VAR . For the former, SAMPLEMONOMIAL correctly
 1482 returns the sign value of the gate. For the latter, SAMPLEMONOMIAL returns the correct
 1483 sign of 1, since a variable is a neutral element, and 1 is the multiplicative identity, whose
 1484 product with another sign element will not change that sign element.

1485 For the inductive hypothesis, we assume for a circuit of depth $d \leq k$ and $k \geq 0$ that the
 1486 algorithm correctly returns the sign value of \mathbf{c} .

1487 Similar to before, for a depth $d \leq k + 1$, it is true that \mathbf{C}_L and \mathbf{C}_R both return the correct
 1488 sign of \mathbf{c} . For the case that $\mathbf{C.type} = \times$, the sign value of both inputs are multiplied, which
 1489 is the correct behavior by definition 4.1. When $\mathbf{C.type} = +$, only one input of \mathbf{C} is sampled,
 1490 and the algorithm returns the correct sign value of \mathbf{c} by inductive hypothesis.

1491 Run-time Analysis

1492 It is easy to check that except for lines 3 and 10, all lines take $O(1)$ time. Consider an
 1493 execution of line 10. We note that we will be adding a given set of variables to some set at
 1494 most once: since the sum of the sizes of the sets at a given level is at most $\text{DEG}(\mathbf{C})$, each gate
 1495 visited takes $O(\log \text{DEG}(\mathbf{C}))$. For Line 3, note that we pick \mathbf{C}_L with probability $\frac{a}{a+b}$ where
 1496 $a = \mathbf{C.Lweight}$ and $b = \mathbf{C.Rweight}$. We can implement this step by picking a random number
 1497 $r \in [a + b]$ and then checking if $r \leq a$. It is easy to check that $a + b \leq |\mathbf{C}|(1, \dots, 1)$. This
 1498 means we need to add and compare $\log |\mathbf{C}|(1, \dots, 1)$ -bit numbers, which can certainly be done
 1499 in time $\overline{\mathcal{M}}(\log(|\mathbf{C}|(1, \dots, 1)), \log \text{SIZE}(\mathbf{C}))$ (note that this is an over-estimate). Denote COST
 1500 (\mathbf{C}) (Eq. (35)) to be an upper bound of the number of gates visited by SAMPLEMONOMIAL.
 1501 Then the runtime is $O(\text{COST}(\mathbf{C}) \cdot \log \text{DEG}(\mathbf{C}) \cdot \overline{\mathcal{M}}(\log(|\mathbf{C}|(1, \dots, 1)), \log \text{SIZE}(\mathbf{C})))$.

1502 We now bound the number of recursive calls in SAMPLEMONOMIAL by $O((\text{DEG}(\mathbf{C}) + 1) \cdot$
 1503 $\text{DEPTH}(\mathbf{C}))$, which by the above will prove the claimed runtime.

1504 Let $\text{COST}(\cdot)$ be a function that models an upper bound on the number of gates that can
 1505 be visited in the run of SAMPLEMONOMIAL. We define $\text{COST}(\cdot)$ recursively as follows.

$$1506 \quad \text{COST}(\mathbf{C}) = \begin{cases} 1 + \text{COST}(\mathbf{C}_L) + \text{COST}(\mathbf{C}_R) & \text{if } \mathbf{C.type} = \times \\ 1 + \max(\text{COST}(\mathbf{C}_L), \text{COST}(\mathbf{C}_R)) & \text{if } \mathbf{C.type} = + \\ 1 & \text{otherwise} \end{cases} \quad (35)$$

1507 First note that the number of gates visited in SAMPLEMONOMIAL is $\leq \text{COST}(\mathcal{C})$. To show
 1508 that eq. (35) upper bounds the number of nodes visited by SAMPLEMONOMIAL, note that
 1509 when SAMPLEMONOMIAL visits a gate such that $\mathcal{C}.\text{type} = \times$, line 8 visits each input of \mathcal{C} , as
 1510 defined in (35). For the case when $\mathcal{C}.\text{type} = +$, line 3 visits exactly one of the input gates,
 1511 which may or may not be the subcircuit with the maximum number of gates traversed, which
 1512 makes $\text{COST}(\cdot)$ an upperbound. Finally, it is trivial to see that when $\mathcal{C}.\text{type} \in \{\text{VAR}, \text{NUM}\}$,
 1513 i.e., a source gate, that only one gate is visited.

1514 We prove the following inequality holds.

$$1515 \quad 2(\text{DEG}(\mathcal{C}) + 1) \cdot \text{DEPTH}(\mathcal{C}) + 1 \geq \text{COST}(\mathcal{C}) \quad (36)$$

1516 Note that eq. (36) implies the claimed runtime. We prove eq. (36) for the number of
 1517 gates traversed in SAMPLEMONOMIAL using induction over $\text{DEPTH}(\mathcal{C})$. Recall how degree is
 1518 defined in definition 4.4.

1519 For the base case $\text{DEG}(\mathcal{C}) = \{0, 1\}$, $\text{DEPTH}(\mathcal{C}) = 0$, $\text{COST}(\mathcal{C}) = 1$, and it is trivial to see
 1520 that the inequality $2\text{DEG}(\mathcal{C}) \cdot \text{DEPTH}(\mathcal{C}) + 1 \geq \text{COST}(\mathcal{C})$ holds.

1521 For the inductive hypothesis, we assume the bound holds for any circuit where $\ell \geq$
 1522 $\text{DEPTH}(\mathcal{C}) \geq 0$. Now consider the case when SAMPLEMONOMIAL has an arbitrary circuit \mathcal{C}
 1523 input with $\text{DEPTH}(\mathcal{C}) = \ell + 1$. By definition $\mathcal{C}.\text{type} \in \{+, \times\}$. Note that since $\text{DEPTH}(\mathcal{C}) \geq 1$,
 1524 \mathcal{C} must have input(s). Further we know that by the inductive hypothesis the inputs \mathcal{C}_i for
 1525 $i \in \{\text{L}, \text{R}\}$ of the sink gate \mathcal{C} uphold the bound

$$1526 \quad 2(\text{DEG}(\mathcal{C}_i) + 1) \cdot \text{DEPTH}(\mathcal{C}_i) + 1 \geq \text{COST}(\mathcal{C}_i). \quad (37)$$

1527 In particular, since for any i , eq. (37) holds, then it immediately follows that an inequality
 1528 whose operands consist of a sum of the aforementioned inequalities must also hold. This is
 1529 readily seen in the inequality of eq. (39) and eq. (40), where $2(\text{DEG}(\mathcal{C}_L) + 1) \cdot \text{DEPTH}(\mathcal{C}_L) \geq$
 1530 $\text{COST}(\mathcal{C}_L)$, likewise for \mathcal{C}_R , and $1 \geq 1$. It is also true that $\text{DEPTH}(\mathcal{C}_L) \leq \text{DEPTH}(\mathcal{C}) - 1$ and
 1531 $\text{DEPTH}(\mathcal{C}_R) \leq \text{DEPTH}(\mathcal{C}) - 1$.

1532 If $\mathcal{C}.\text{type} = +$, then $\text{DEG}(\mathcal{C}) = \max(\text{DEG}(\mathcal{C}_L), \text{DEG}(\mathcal{C}_R))$. Otherwise $\mathcal{C}.\text{type} = \times$ and
 1533 $\text{DEG}(\mathcal{C}) = \text{DEG}(\mathcal{C}_L) + \text{DEG}(\mathcal{C}_R) + 1$. In either case it is true that $\text{DEPTH}(\mathcal{C}) = \max(\text{DEPTH}(\mathcal{C}_L), \text{DEPTH}(\mathcal{C}_R)) +$
 1534 1 .

1535 If $\mathcal{C}.\text{type} = \times$, then, by eq. (35), substituting values, the following should hold,

$$1536 \quad 2(\text{DEG}(\mathcal{C}_L) + \text{DEG}(\mathcal{C}_R) + 2) \cdot (\max(\text{DEPTH}(\mathcal{C}_L), \text{DEPTH}(\mathcal{C}_R)) + 1) + 1 \quad (38)$$

$$1537 \quad \geq 2(\text{DEG}(\mathcal{C}_L) + 1) \cdot \text{DEPTH}(\mathcal{C}_L) + 2(\text{DEG}(\mathcal{C}_R) + 1) \cdot \text{DEPTH}(\mathcal{C}_R) + 3 \quad (39)$$

$$1538 \quad \geq 1 + \text{COST}(\mathcal{C}_L) + \text{COST}(\mathcal{C}_R) = \text{COST}(\mathcal{C}). \quad (40)$$

1540 To prove (39), first, eq. (38) expands to,

$$1541 \quad 2\text{DEG}(\mathcal{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathcal{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(\mathcal{C}_L) + 2\text{DEG}(\mathcal{C}_R) + 4 + 1 \quad (41)$$

1542 where DEPTH_{\max} is used to denote the maximum depth of the two input subcircuits. Eq. (39)
 1543 expands to

$$1544 \quad 2\text{DEG}(\mathcal{C}_L) \cdot \text{DEPTH}(\mathcal{C}_L) + 2\text{DEPTH}(\mathcal{C}_L) + 2\text{DEG}(\mathcal{C}_R) \cdot \text{DEPTH}(\mathcal{C}_R) + 2\text{DEPTH}(\mathcal{C}_R) + 3 \quad (42)$$

1545 Putting Eq. (41) and Eq. (42) together we get

$$1546 \quad 2\text{DEG}(\mathcal{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathcal{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(\mathcal{C}_L) + 2\text{DEG}(\mathcal{C}_R) + 5$$

$$1547 \quad \geq 2\text{DEG}(\mathcal{C}_L) \cdot \text{DEPTH}(\mathcal{C}_L) + 2\text{DEG}(\mathcal{C}_R) \cdot \text{DEPTH}(\mathcal{C}_R) + 2\text{DEPTH}(\mathcal{C}_L) + 2\text{DEPTH}(\mathcal{C}_R) + 3 \quad (43)$$

1548

1549 Since the following is always true,

$$\begin{aligned}
1550 \quad & 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 5 \\
1551 \quad & \geq 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_R) + 3,
\end{aligned}$$

1553 then it is the case that Eq. (43) is *always* true.

1554 Now to justify (40) which holds for the following reasons. First, eq. (40) is the result of
1555 Eq. (35) when $\mathbf{C.type} = \times$. Eq. (39) is then produced by substituting the upperbound of
1556 (37) for each $\text{COST}(\mathbf{C}_i)$, trivially establishing the upper bound of (40). This proves eq. (36)
1557 for the \times case.

1558 For the case when $\mathbf{C.type} = +$, substituting values yields

$$1559 \quad 2(\max(\text{DEG}(\mathbf{C}_L), \text{DEG}(\mathbf{C}_R)) + 1) \cdot (\max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) + 1) + 1 \quad (44)$$

$$1560 \quad \geq \max(2(\text{DEG}(\mathbf{C}_L) + 1) \cdot \text{DEPTH}(\mathbf{C}_L) + 1, 2(\text{DEG}(\mathbf{C}_R) + 1) \cdot \text{DEPTH}(\mathbf{C}_R) + 1) + 1 \quad (45)$$

$$1561 \quad \geq 1 + \max(\text{COST}(\mathbf{C}_L), \text{COST}(\mathbf{C}_R)) = \text{COST}(\mathbf{C}) \quad (46)$$

1563 To prove (45), eq. (44) expands to

$$1564 \quad 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEG}_{\max} + 2\text{DEPTH}_{\max} + 2 + 1. \quad (47)$$

1565 Since $\text{DEG}_{\max} \cdot \text{DEPTH}_{\max} \geq \text{DEG}(\mathbf{C}_i) \cdot \text{DEPTH}(\mathbf{C}_i)$, the following upper bound holds for the
1566 expansion of eq. (45):

$$1567 \quad 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEPTH}_{\max} + 2 \quad (48)$$

1568 Putting it together we obtain the following for (45):

$$\begin{aligned}
1569 \quad & 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEG}_{\max} + 2\text{DEPTH}_{\max} + 3 \\
1570 \quad & \geq 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEPTH}_{\max} + 2,
\end{aligned} \quad (49)$$

1572 where it can be readily seen that the inequality stands and (49) follows. This proves (45).

1573 Similar to the case of $\mathbf{C.type} = \times$, (46) follows by equations (35) and (37).1574 This proves (36) as desired. ◀

1575 D.10 Experimental Results

1576 Recall that by definition of BIDB, a query result cannot be derived by a self-join between
1577 non-identical tuples belonging to the same block. Note, that by Theorem 4.7, γ must be
1578 a constant in order for Algorithm 1 to achieve linear time. We would like to determine
1579 experimentally whether queries over BIDB instances in practice generate a constant number
1580 of cancellations or not. Such an experiment would ideally use a database instance with
1581 queries both considered to be typical representations of what is seen in practice.

1582 We ran our experiments using Windows 10 WSL Operating System with an Intel Core i7
1583 2.40GHz processor and 16GB RAM. All experiments used the PostgreSQL 13.0 database
1584 system.

1585 For the data we used the MayBMS data generator [1] tool to randomly generate uncertain
1586 versions of TPC-H tables. The queries computed over the database instance are Q_1 , Q_2 , and
1587 Q_3 from [5], all of which are modified versions of TPC-H queries Q_3 , Q_6 , and Q_7 where all
1588 aggregations have been dropped.

1589 As written, the queries disallow BIDD cross terms. We first ran all queries, noting the
 1590 result size for each. Next the queries were rewritten so as not to filter out the cross terms.
 1591 The comparison of the sizes of both result sets should then suggest in one way or another
 1592 whether or not there exist many cross terms in practice. As seen, the experimental query
 1593 results contain little to no cancelling terms. Fig. 5 shows the result sizes of the queries,
 1594 where column CF is the result size when all cross terms are filtered out, column CI shows
 1595 the number of output tuples when the cancelled tuples are included in the result, and the
 1596 last column is the value of γ . The experiments show γ to be in a range between $[0, 0.1]\%$,
 1597 indicating that only a negligible or constant (compare the result sizes of $Q_1 < Q_2$ and their
 1598 respective γ values) amount of tuples are cancelled in practice when running queries over a
 1599 typical BIDD instance. Interestingly, only one of the three queries had tuples that violated
 1600 the BIDD constraint.

1601 To conclude, the results in Fig. 5 show experimentally that γ is negligible in practice for
 1602 BIDD queries. We also observe that (i) tuple presence is independent across blocks, so the
 1603 corresponding probabilities (and hence p_0) are independent of the number of blocks, and (ii)
 1604 BIDDs model uncertain attributes, so block size (and hence γ) is a function of the “messiness”
 1605 of a dataset, rather than its size. Thus, we expect Theorem 4.7 to hold in general.

| Query | CF | CI | γ |
|-------|---------|---------|----------|
| Q_1 | 46,714 | 46,768 | 0.1% |
| Q_2 | 179,917 | 179,917 | 0% |
| Q_3 | 11,535 | 11,535 | 0% |

■ **Figure 5** Number of Cancellations for Queries Over BIDD.

1606 E Circuits

1607 E.1 Representing Polynomials with Circuits

1608 E.1.1 Circuits for query plans

Atri says: Since this comment is not showing up below, I do not follow why the last sentence of this para is true.

1609
 1610 We now formalize circuits and the construction of circuits for \mathcal{RA}^+ queries. As mentioned
 1611 earlier, we represent lineage polynomials as arithmetic circuits over \mathbb{N} -valued variables with
 1612 $+$, \times . A circuit for query Q and $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ is a directed acyclic graph
 1613 $\langle V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle$ with vertices $V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}$ and directed edges $E_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \subset$
 1614 $V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}^2$. The sink function $\phi_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} : \mathcal{U}^n \rightarrow V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}$ is a partial function that maps the
 1615 tuples of the n -ary

Atri says: In the main paper we have used n to denote the number of input tuples so we need to use some other notation n but since I do not know where all this change will need to be propagated so am not changing it for now.

1616
 1617 relation $Q(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ to vertices. We require that $\phi_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}$'s range be limited to sink vertices
 1618 (i.e., vertices with out-degree 0). A function $\ell_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} : V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rightarrow \{+, \times\} \cup \mathbb{N} \cup \mathbf{X}$ assigns a
 1619 label to each node: Source nodes (i.e., vertices with in-degree 0) are labeled with constants
 1620 or variables (i.e., $\mathbb{N} \cup \mathbf{X}$), while the remaining nodes are labeled with the symbol $+$ or \times . We
 1621 require that vertices have an in-degree of at most two. Note that we can construct circuits for

1622 BIDs in time linear in the time required for deterministic query processing over a possible
 1623 world of the BIDB under the aforementioned assumption that $|\mathcal{D}_{\mathbb{N}[\mathbf{X}]}| \leq c \cdot |D|$.

1624 **Atri says:** I do not follow the last sentence.

1625 E.2 Modeling Circuit Construction

1626 We now connect the size of a circuit (where the size of a circuit is the number of vertices
 1627 in the corresponding DAG) for a given \mathcal{RA}^+ query Q and $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ to the
 1628 runtime $T_{det}^*(Q, D_\Omega)$ of the PDB's deterministic bounding database D_Ω . We do this formally
 1629 by showing that the size of the circuit is asymptotically no worse than the corresponding
 1630 runtime of a large class of deterministic query processing algorithms.

1631 Each vertex $v \in V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}$ in the arithmetic circuit for

$$1632 \langle V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle, E_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle, \phi_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle, \ell_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle$$

1633 encodes a polynomial, realized as

$$1634 \text{lin}(v) = \begin{cases} \sum_{v': (v', v) \in E_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}} \text{lin}(v') & \text{if } \ell(v) = + \\ \prod_{v': (v', v) \in E_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}} \text{lin}(v') & \text{if } \ell(v) = \times \\ \ell(v) & \text{otherwise} \end{cases}$$

1635 We define the circuit for a \mathcal{RA}^+ query Q recursively by cases as follows. In each case, let
 1636 $\langle V_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle, E_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle, \phi_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle, \ell_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle$ denote the circuit for subquery Q_i . We implicitly
 1637 include in all circuits a global zero node v_0 s.t., $\ell_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}(v_0) = 0$ for any $Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}$.

1638 Algorithm 4 defines how the circuit for a query result is constructed. We quickly review
 1639 the number of vertices emitted in each case.

1640 **Base Relation.** This circuit has $|D_\Omega \cdot R|$ vertices.

1641 **Selection.** If we assume dead sinks are iteratively garbage collected, this circuit has at
 1642 most $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}|$ vertices.

1643 **Projection.** This formulation will produce vertices with an in-degree greater than two, a
 1644 problem that we correct by replacing every vertex with an in-degree over two by an equivalent
 1645 fan-in two tree. The resulting structure has at most $|Q_1| - 1$ new vertices. The corrected
 1646 circuit thus has at most $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1|$ vertices.

1647 **Union.** This circuit has $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |V_{Q_2, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1 \cap Q_2|$ vertices.

1648 **k -ary Join.** As in projection, newly created vertices will have an in-degree of k , and a
 1649 fan-in two tree is required. There are $|Q_1 \bowtie \dots \bowtie Q_k|$ such vertices, so the corrected circuit
 1650 has $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + \dots + |V_{Q_k, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + (k-1)|Q_1 \bowtie \dots \bowtie Q_k|$ vertices.

1651 E.2.1 Bounding circuit depth

1652 We first show that the depth of the circuit (DEPTH; Definition 4.3) is bounded by the size
 1653 of the query. Denote by $|Q|$ the number of relational operators in query Q , which recall we
 1654 assume is a constant.

1655 **► Proposition E.1** (Circuit depth is bounded). *Let Q be a relational query and D_Ω be a*
 1656 *deterministic bounding database with n tuples. There exists a (lineage) circuit \mathcal{C}^* encoding*
 1657 *the lineage of all tuples $t \in Q(D_\Omega)$ for which $\text{DEPTH}(\mathcal{C}^*) \leq O(k|Q| \log(n))$.*

Algorithm 4 $LC(Q, D_\Omega, E, V, \ell)$

Input: Q : query**Input:** D_Ω : a deterministic bounding database**Input:** E, V, ℓ : accumulators for the edge list, vertex list, and vertex label list.**Output:** $C = \langle E, V, \phi, \ell \rangle$: a circuit encoding the lineage of each tuple in $Q(D_\Omega)$

```

1: if  $Q$  is  $R$  then                                     ▷ Case 1:  $Q$  is a relation atom
2:   for  $t \in D_\Omega.R$  do
3:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, R(t))\}$            ▷ Allocate a fresh node  $v_t$ 
4:      $\phi(t) \leftarrow v_t$ 
5:   end for
6: else if  $Q$  is  $\sigma_\theta(Q')$  then                       ▷ Case 2:  $Q$  is a Selection
7:    $\langle V, E, \phi', \ell \rangle \leftarrow LC(Q', D_\Omega, V, E, \ell)$ 
8:   for  $t \in \text{DOM}(\phi')$  do
9:     if  $\theta(t)$  then  $\phi(t) \leftarrow \phi'(t)$  else  $\phi(t) \leftarrow v_0$ 
10:  end for
11: else if  $Q$  is  $\pi_{\bar{A}}(Q')$  then                           ▷ Case 3:  $Q$  is a Projection
12:    $\langle V, E, \phi', \ell \rangle \leftarrow LC(Q', D_\Omega, V, E, \ell)$ 
13:   for  $t \in \pi_{\bar{A}}(Q'(D_\Omega))$  do
14:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, +)\}$            ▷ Allocate a fresh node  $v_t$ 
15:      $\phi(t) \leftarrow v_t$ 
16:   end for
17:   for  $t \in Q'(D_\Omega)$  do
18:      $E \leftarrow E \cup \{(\phi'(t), \phi(\pi_{\bar{A}}t))\}$ 
19:   end for
20:   Correct nodes with in-degrees  $> 2$  by appending an equivalent fan-in two tree instead
21: else if  $Q$  is  $Q_1 \cup Q_2$  then                           ▷ Case 4:  $Q$  is a Bag Union
22:    $\langle V, E, \phi_1, \ell \rangle \leftarrow LC(Q_1, D_\Omega, V, E, \ell)$ 
23:    $\langle V, E, \phi_2, \ell \rangle \leftarrow LC(Q_2, D_\Omega, V, E, \ell)$ 
24:    $\phi \leftarrow \phi_1 \cup \phi_2$ 
25:   for  $t \in \text{DOM}(\phi_1) \cap \text{DOM}(\phi_2)$  do
26:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, +)\}$            ▷ Allocate a fresh node  $v_t$ 
27:      $\phi(t) \leftarrow v_t$ 
28:      $E \leftarrow E \cup \{(\phi_1(t), v_t), (\phi_2(t), v_t)\}$ 
29:   end for
30: else if  $Q$  is  $Q_1 \bowtie \dots \bowtie Q_m$  then             ▷ Case 5:  $Q$  is a  $m$ -ary Join
31:   for  $i \in [m]$  do
32:      $\langle V, E, \phi_i, \ell \rangle \leftarrow LC(Q_i, D_\Omega, V, E, \ell)$ 
33:   end for
34:   for  $t \in \text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_m)$  do
35:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, \times)\}$            ▷ Allocate a fresh node  $v_t$ 
36:      $\phi(t) \leftarrow v_t$ 
37:      $E \leftarrow E \cup \{(\phi_i(\pi_{sch(Q_i(D_\Omega))}(t)), v_t) \mid i \in [m]\}$ 
38:   end for
39:   Correct nodes with in-degrees  $> 2$  by appending an equivalent fan-in two tree instead
40: end if

```

1658 **Proof.** We show that the bound of Proposition E.1 holds for the circuit constructed by
 1659 Algorithm 4. First, observe that Algorithm 4 is (recursively) invoked exactly once for every
 1660 relational operator or base relation in Q ; It thus suffices to show that a call to Algorithm 4
 1661 adds at most $O_k(\log(n))$ to the depth of a circuit produced by any recursive invocation.
 1662 Second, observe that modulo the logarithmic fan-in of the projection and join cases, the
 1663 depth of the output is at most one greater than the depth of any input (or at most 1 in the
 1664 base case of relation atoms). For the join case, the number of in-edges can be no greater than
 1665 the join width, which itself is bounded by k . The depth thus increases by at most a constant
 1666 factor of $\lceil \log(k) \rceil = O_k(1)$. For the projection case, observe that the fan-in is bounded by
 1667 $|Q'(D_\Omega)|$, which is in turn bounded by n^k . The depth increase for any projection node is
 1668 thus at most $\lceil \log(n^k) \rceil = O(k \log(n))$, as desired. \blacktriangleleft

1669 E.2.2 Circuit size vs. runtime

1670 **► Lemma E.2.** *Given a $\mathbb{N}[\mathbf{X}]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ with deterministic bounding database D_Ω ,*
 1671 *and an \mathcal{RA}^+ query Q , the runtime of Q over D_Ω has the same or greater complexity as the*
 1672 *size of the lineage of $Q(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$. That is, we have $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq kT_{det}^*(Q, D_\Omega) + 1$, where $k \geq 1$*
 1673 *is the maximal degree of any polynomial in $Q(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$.*

1674 **Proof.** We prove by induction that $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \setminus \{v_0\}| \leq kT_{det}^*(Q, D_\Omega)$. For clarity, we
 1675 implicitly exclude v_0 in the proof below.

1676 The base case is a base relation: $Q = R$ and is trivially true since $|V_{R, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| = |D_\Omega.R| =$
 1677 $T_{det}^*(R, D_\Omega)$ (note that here the degree $k = 1$). For the inductive step, we assume that we
 1678 have circuits for subqueries Q_1, \dots, Q_m such that $|V_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq k_i T_{det}^*(Q_i, D_\Omega)$ where k_i is
 1679 the degree of Q_i .

1680 **Selection.** Assume that $Q = \sigma_\theta(Q_1)$. In the circuit for Q , $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| = |V_{Q_1, D_\Omega}|$ vertices,
 1681 so from the inductive assumption and $T_{det}^*(Q, D_\Omega) = T_{det}^*(Q_1, D_\Omega)$ by definition, we have
 1682 $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq kT_{det}^*(Q, D_\Omega)$.

1683 **Projection.** Assume that $Q = \pi_{\mathbf{A}}(Q_1)$. The circuit for Q has at most $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1|$
 1684 vertices.

$$1685 \quad |V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq |V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1|$$

1687 (From the inductive assumption)

$$1688 \quad \leq kT_{det}^*(Q_1, D_\Omega) + |Q_1|$$

1690 (By definition of $T_{det}^*(Q, D_\Omega)$)

$$1691 \quad \leq kT_{det}^*(Q, D_\Omega).$$

1693 **Union.** Assume that $Q = Q_1 \cup Q_2$. The circuit for Q has $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |V_{Q_2, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1 \cap Q_2|$
 1694 vertices.

$$1695 \quad |V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq |V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |V_{Q_2, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1| + |Q_2|$$

1697 (From the inductive assumption)

$$1698 \quad \leq k(T_{det}^*(Q_1, D_\Omega) + T_{det}^*(Q_2, D_\Omega)) + (|Q_1| + |Q_2|)$$

1700 (By definition of $T_{det}^*(Q, D_\Omega)$)

$$1701 \quad \leq k(T_{det}^*(Q, D_\Omega)).$$

1702

1703 **m -ary Join.** Assume that $Q = Q_1 \bowtie \dots \bowtie Q_m$. Note that $k = \sum_{i=1}^m k_i \geq m$. The circuit
1704 for Q has $|V_{Q_1, \mathcal{D}_{N[\mathbf{x}]}}| + \dots + |V_{Q_k, \mathcal{D}_{N[\mathbf{x}]}}| + (m-1)|Q_1 \bowtie \dots \bowtie Q_k|$ vertices.

$$1705 \quad |V_{Q, \mathcal{D}_{N[\mathbf{x}]}}| = |V_{Q_1, \mathcal{D}_{N[\mathbf{x}]}}| + \dots + |V_{Q_k, \mathcal{D}_{N[\mathbf{x}]}}| + (m-1)|Q_1 \bowtie \dots \bowtie Q_k|$$

1707 From the inductive assumption and noting $\forall i : k_i \leq k$ and $m \leq k$

$$1708 \quad \leq kT_{det}^*(Q_1, D_\Omega) + \dots + kT_{det}^*(Q_k, D_\Omega) +$$

$$1709 \quad (m-1)|Q_1 \bowtie \dots \bowtie Q_m|$$

$$1710 \quad \leq k(T_{det}^*(Q_1, D_\Omega) + \dots + T_{det}^*(Q_m, D_\Omega) +$$

$$1711 \quad |Q_1 \bowtie \dots \bowtie Q_m|)$$

1713 (By definition of $T_{det}^*(Q, D_\Omega)$ and assumption on $T_{join}(\cdot)$)

$$1714 \quad \leq kT_{det}^*(Q, D_\Omega).$$

1716 The property holds for all recursive queries, and the proof holds. \blacktriangleleft

1717 E.2.3 Runtime of LC

1718 We next need to show that we can construct the circuit in time linear in the deterministic
1719 runtime.

1720 **► Lemma E.3.** *Given a query Q over a deterministic bounding database D_Ω and the \mathcal{C}^**
1721 *output by Algorithm 4, the runtime $T_{LC}(Q, D_\Omega, \mathcal{C}^*) \leq O(T_{det}^*(Q, D_\Omega))$.*

1722 **Proof.** By analysis of Algorithm 4, invoked as $\mathcal{C}^* \leftarrow LC(Q, D_\Omega, \{v_0\}, \emptyset, \{(v_0, 0)\})$.

1723 We assume that the vertex list V , edge list E , and vertex label list ℓ are mutable
1724 accumulators with $O(1)$ amortized append. We assume that the tuple to sink mapping ϕ is
1725 a linked hashmap, with $O(1)$ insertions and retrievals, and $O(n)$ iteration over the domain of
1726 keys. We assume that the n -ary join $\text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_n)$ can be computed in time
1727 $T_{join}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_n))$ (Definition 2.10) and that an intersection $\text{DOM}(\phi_1) \cap \text{DOM}(\phi_2)$
1728 can be computed in time $O(|\text{DOM}(\phi_1)| + |\text{DOM}(\phi_2)|)$ (e.g., with a hash table).

1729 Before proving our runtime bound, we first observe that $T_{det}^*(Q, D) \geq \Omega(|Q(D)|)$. This is
1730 true by construction for the relation, projection, and union cases, by Definition 2.10 for joins,
1731 and by the observation that $|\sigma(R)| \leq |R|$.

1732 We show that $T_{det}^*(Q, D_\Omega)$ is an upper-bound for the runtime of Algorithm 4 by recursion.
1733 The base case of a relation atom requires only an $O(|D_\Omega.R|)$ iteration over the source tuples.
1734 For the remaining cases, we make the recursive assumption that for every subquery Q' , it
1735 holds that $O(T_{det}^*(Q', D_\Omega))$ bounds the runtime of Algorithm 4.

1736 **Selection.** Selection requires a recursive call to Algorithm 4, which by the recursive
1737 assumption is bounded by $O(T_{det}^*(Q', D_\Omega))$. Algorithm 4 requires a loop over every element
1738 of $Q'(D_\Omega)$. By the observation above that $T_{det}^*(Q, D) \geq \Omega(|Q(D)|)$, this iteration is also
1739 bounded by $O(T_{det}^*(Q', D_\Omega))$.

1740 **Projection.** Projection requires a recursive call to Algorithm 4, which by the recursive
1741 assumption is bounded by $O(T_{det}^*(Q', D_\Omega))$, which in turn is a term in $T_{det}^*(\pi_{\bar{A}}Q', D_\Omega)$.
1742 What remains is an iteration over $\pi_{\bar{A}}(Q(D_\Omega))$ (lines 13–16), an iteration over $Q'(D_\Omega)$ (lines
1743 17–19), and the construction of a fan-in tree (line 20). The first iteration is $O(|Q(D_\Omega)|) \leq$

1744 $O(T_{det}^*(Q, D_\Omega))$. The second iteration and the construction of the bounded fan-in tree are
 1745 both $O(|Q'(D_\Omega)|) \leq O(T_{det}^*(Q', D_\Omega)) \leq O(T_{det}^*(Q, D_\Omega))$, by the the observation above that
 1746 $T_{det}^*(Q, D) \geq \Omega(|Q(D)|)$.

1747 **Bag Union.** As above, the recursive calls explicitly correspond to terms in the expansion of
 1748 $T_{det}^*(Q_1 \cup Q_2, D_\Omega)$. Initializing ϕ (line 24) can be accomplished in $O(\text{DOM}(\phi_1) + \text{DOM}(\phi_2)) =$
 1749 $O(|Q_1(D_\Omega)| + |Q_2(D_\Omega)|) \leq O(T_{det}^*(Q_1, D_\Omega) + T_{det}^*(Q_2, D_\Omega))$. The remainder requires computing
 1750 $Q_1 \cup Q_2$ (line 25) and iterating over it (lines 25–29), which is $O(|Q_1| + |Q_2|)$ as noted above
 1751 — this directly corresponds to terms in $T_{det}^*(Q_1 \cup Q_2, D_\Omega)$.

1752 **m -ary Join.** As in the prior cases, recursive calls explicitly correspond to terms in our
 1753 target runtime. The remaining logic involves (i) computing $\text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_m)$, (ii)
 1754 iterating over the results, and (iii) creating a fan-in tree. Respectively, these are:

- 1755 (i) $T_{join}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_m))$
 1756 (ii) $O(|Q_1(D_\Omega) \bowtie \dots \bowtie Q_m(D_\Omega)|) \leq O(T_{join}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_m)))$ (Definition 2.10)
 1757 (iii) $O(m|Q_1(D_\Omega) \bowtie \dots \bowtie Q_m(D_\Omega)|)$ (as (ii), noting that $m \leq k = O(1)$) ◀

1758 F Higher Moments

1759 We make a simple observation to conclude the presentation of our results. So far we have only
 1760 focused on the expectation of Φ . In addition, we could e.g. prove bounds of the probability
 1761 of a tuple’s multiplicity being at least 1. Progress can be made on this as follows: For any
 1762 positive integer m we can compute the m -th moment of the multiplicities, allowing us to e.g.
 1763 use the Chebyshev inequality or other high moment based probability bounds on the events
 1764 we might be interested in. We leave further investigations for future work.

1765 G The Karp-Luby Estimator

1766 Computing the marginal probability of a tuple in the output of a set-probabilistic database
 1767 query has been studied extensively. To the best of our knowledge, the current state of
 1768 the art approximation algorithm for this problem is the Karp-Luby estimator [30], which
 1769 first appeared in MayBMS/Sprout [38], and more recently as part of an online “anytime”
 1770 approximation algorithm [20, 15].

The estimator works by observing that for any ℓ random binary (but not necessarily independent) events $\mathbf{W}_1, \dots, \mathbf{W}_\ell$, the probability of at least one event occurring (i.e., $Pr(\mathbf{W}_1 \vee \dots \vee \mathbf{W}_\ell)$) is bounded from above by the sum of the independent event probabilities (i.e., $Pr(\mathbf{W}_1 \vee \dots \vee \mathbf{W}_\ell) \leq Pr(\mathbf{W}_1) + \dots + Pr(\mathbf{W}_\ell)$). Starting from this (‘easily’ computable and large) value, the estimator proceeds to correct the estimate by estimating how much of an over-estimate it is. Specifically, if \mathcal{P} is the joint distribution over \mathbf{W} , the estimator computes an approximation of:

$$\mathcal{O} = \mathbb{E}_{\mathbf{w} \sim \mathcal{P}} \left[|\{i \mid \mathbf{W}_i = 1, i \in [\ell]\}| \right].$$

The accuracy of this estimate is improved by conditioning \mathcal{P} on a W_i chosen uniformly at random (which ensures that the sampled count will be at least 1) and correcting the resulting estimate by $Pr(W_i)$. With an estimate of \mathcal{O} , it can easily be verified that the probability of the disjunction can be computed as:

$$Pr(\mathbf{W}_1 \vee \dots \vee \mathbf{W}_\ell) = Pr(\mathbf{W}_1) + \dots + Pr(\mathbf{W}_\ell) - \mathcal{O}$$

1771 The Karp-Luby estimator is employed on the SMB representation¹⁴ of \mathcal{C} (to solve the
 1772 set-PDB version of Problem 1.6), where each W_i represents the event that one monomial
 1773 is true. By simple inspection, if there are ℓ monomials, this estimator has runtime $\Omega(\ell)$.
 1774 Further, a minimum of $\left\lceil \frac{3 \cdot \ell \cdot \log(\frac{2}{\delta})}{\epsilon^2} \right\rceil$ invocations of the estimator are required to achieve $1 \pm \epsilon$
 1775 approximation with probability at least $1 - \delta$ [38], entailing a runtime at least quadratic in ℓ .
 1776 As an arbitrary lineage circuit \mathcal{C} may encode $\Omega(|\mathcal{C}|^k)$ monomials, the worst case runtime is
 1777 at least $\Omega(|\mathcal{C}|^{2k})$ (where k is the ‘degree’ of lineage polynomial encoded by \mathcal{C}). By contrast
 1778 note that by the discussion after Lemma 4.8 we can solve Problem 1.6 in time $O(|\mathcal{C}|^2)$ for
 1779 all BIDB circuits *independent* of the degree k .

1780 H Parameterized Complexity

1781 In Sec. 3, we utilized common conjectures from fine-grained complexity theory. The notion of
 1782 $\#W[1]$ – *hard* is a standard notion in *parameterized complexity*, which by now is a standard
 1783 complexity tool in providing data complexity bounds on query processing results [22]. E.g.
 1784 the fact that k -matching is $\#W[1]$ – *hard* implies that we cannot have an $n^{\Omega(1)}$ runtime.
 1785 However, these results do not carefully track the exponent in the hardness result. E.g.
 1786 $\#W[1]$ – *hard* for the general k -matching problem does not imply anything specific for the
 1787 3-matching problem. Similar questions have led to intense research into the new sub-field
 1788 of *fine-grained complexity* (see [48]), where we care about the exponent in our hardness
 1789 assumptions as well– e.g. Conjecture 3.3 is based on the popular *Triangle detection hypothesis*
 1790 in this area (cf. [34]).

1791 I Response to first cycle reviewer comments

1792 This paper is a resubmission of our submission to the ICDT first cycle. We thank the
 1793 reviewers for their insightful comments, which we believe has helped improve the presentation
 1794 of the paper tremendously. We use this section to document the changes that have been made
 1795 since our prior submission, and in particular, how we have addressed reviewer comments
 1796 (reviewer comments are shaded followed by our responses).

1797 I.1 Meta Review

1798 Problem definition not stated rigorously nor motivated. Discussion needed on the standard
 PDB approach vs your approach.

1799 We rewrote Sec. 1 to specifically address this concern. The opening paragraph precisely
 1800 and formally states the query evaluation problem in bag-PDBs. We use a series of problem
 1801 statements to clearly define the problem we are addressing as it relates to the query evaluation
 1802 problem. We made the concrete problem statements more precise by more clearly formalizing
 1803 $T_{det}^*(Q, D_\Omega)$ and stating our runtime objectives relative to it (??, 1.5, 1.6).

1804 We have included a discussion of the standard approach, e.g. see the paragraph
 1805 **Relationship to Set-Probabilistic Query Evaluation** on page 4.

1806 Definition 2.6 on reduced BIDB polynomials seem not the right tool for the studied
 problem.

¹⁴Note that since we are in the set semantics, in the lineage polynomial/formula, addition is logical OR and multiplication is logical AND.

1807 We have chosen to stick with a less formal, ad-hoc definition (please see Definition 1.3 and
 1808 ??) of the general problem as suggested by both Reviewer 1 and Reviewer 2. Our earlier
 1809 proof of the current ?? (in the appendix) had a small bug, which also has been fixed.

1810 The paper is very difficult to read. Improvements are needed in particular for the
 presentation of the approximation results and their proofs. Also for the notation. Missing
 definitions for used notions need to be added. Ideally use one instead of three query
 languages (UCQ, RA+, SPJU).

1811 We have chosen one specific query language throughout the paper (\mathcal{RA}^+) and made a
 1812 concerted effort to use clean, defined, non-ambiguous notation. We have also simplified the
 1813 notation by limiting the paper's use of provenance semirings (which are needed solely for
 1814 proofs) to the appendix. To the best of our examination, all notation conflicts have been
 1815 addressed and definitions for used notions are added (see e.g. Definition C.4 appears before
 1816 Lemma C.6 and Lemma C.8).

1817 After the rewrite of Sec. 1, we had even less space for Sec. 4. However, we have modified
 1818 Sec. 4 so that it flows better. In particular, we start off with the algorithm idea first
 1819 (paragraph **Overview of our Techniques** in Sec. 1 also has more details on the intuition
 1820 behind the approximation algorithm) and then state the results (with more details on how
 1821 we argue the claimed runtime). Finally, we clearly state Corollary 4.9 for which queries our
 1822 linear-time approximation result holds.

1823 1.2 Reviewer 1

1824 1.24 "is #W[1]-hard": parameterized by what?

1825 1.103 and 1.105: again, what is the parameter exactly?

1826 While the above references do not exist in the revised Sec. 1 anymore, all theorem statements
 1827 and claims on #W[1] runtime have been stated in a way so as to avoid ambiguity in the
 1828 parameter. Please see e.g. Theorem 3.1 and Theorem 3.6.

1829 You might want to explain your title somewhere (probably in the introduction): in the
 end, what exactly should be considered harmful and why?

1830 We have modified the title to be more descriptive.

1831 1.45 when discussing Dalvi and Suciu's dichotomy, you might want to mention that they
 consider *data complexity*. Currently the second sentence of your introduction ("take a
 query Q and a pdb D") suggests that you are considering combined complexity.

1832 We have made an explicit mention of data complexity when alluding to Dalvi and Suciu's
 1833 dichotomy. We have further rewritten Sec. 1 in such a way as to explicitly note the type(s)
 1834 of complexity we are considering (mostly it's parameterized complexity).

1835 1.51 "Consider ... and tuples are independent random event": so this is actually a set
 PDB... You might want to use an example where the input PDB is actually a bag PDB.
 The last sentence before the example makes the reader *expect* that the example will be
 of a bag PDB that is not a set PDB

1836 Our revision has removed the example referred to above. While the paper considers inputs to
 1837 queries that are equivalent to set-PDB, this is not limiting. Please see ?? on ??. Furthermore,
 1838 we have added a discussion to the appendix that expands on why our results do extend
 1839 beyond set inputs (Appendix A).

- In the case of set semantics, the lineage of a tuple can be defined for *any* query: it is the unique Boolean function that satisfies the if and only if property that you mention on line 70. For bag semantics however, to the best of my knowledge there is no general definition of what is a lineage for an arbitrary query. On line 73, it is not clear at all how the polynomial should be defined, since this will depend on the type of query that you consider

Note that lineage for a set semantics query is as a positive Boolean formula is defined for positive relational algebra. For instance, for aggregate queries a more powerful model ([4]) is needed. The definition of the lineage polynomial (bag PDB) semantics over an arbitrary \mathcal{RA}^+ query Q is now given in Fig. 1. We also note that these semantics are not novel (e.g., similar semantics appear for both provenance [25] and probabilistic database [32, 19] contexts). However, as we were unable to find a formal proof of the equivalence between the expectation of the query multiplicity and of the lineage polynomial in related work, we have included a proof of Proposition 2.5.

1.75 "evaluating the lineage of t over an assignment corresponding to a possible world": here, does the assignment assigns each tuple to true or false? In other words, do the variables X still represent individual tuples? From what I see later in the article it seems that no, so this is confusing if we compare to what is explained in the previous paragraph about set TIDB

The discussion after Problem 1.2 (in particular, the paragraph TIDBs) specifically address these questions. While values for possible worlds assigned are from $\{0, 1\}$, which is analog to Boolean, this is not limiting. Please see ?? (??) and the new appendix section Appendix A.

- 1.135 "polynomial $Q(X)$ ": Q should be reserved for queries... You could use φ or ϕ or... anything else but Q really

We now use $\Phi(\mathbf{X})$ for (lineage) polynomials.

- If we consider data complexity (as did Dalvi and Suciu) and fix an UCQ Q , given as input a bag TIDB PDB we can always compute the lineage in $O(|D|^{|Q|})$ in SOP form and from there compute the expected multiplicity with the same complexity, so in polynomial time. How does this relate to your hardness result? Is it that you are only interested in combined complexity? Why one shouldn't be happy with this running time? Usually queries are much smaller than databases and this motivates studying data complexity.

We have rewritten Sec. 1 in a way to stress that we are primarily interested in data complexity, but we cannot stop there. As the reviewer has noted, the problem we explore requires further analysis, where we require parameterized and fine grained complexity analysis to provide a theoretical foundation for the question we ask in ??. We have discussed this in the prose following Problem 1.2.

A discussion is missing about the difference between the approach usually taken in PDB literature and your approach. In which case would one be more interested in the expected multiplicity or in the marginal probability of a tuple? This should be discussed clearly in the introduction, as currently there is no clear "motivation" to what you do. There is a section about related work at the end but it is mostly a set of facts and there is no insightful comparison to what you do.

We provide more motivating examples in the first paragraph, and include a more detailed discussion of the relationship to sets in paragraph **Relationship to Set-Probabilistic Query Evaluation** after ??. For example, expected multiplicities can model expectation

1866 of a `COUNT(*)` query, while in many contexts computing the probability that this count is
 1867 non-zero is not that useful.

1868 As we now explain in the introduction, another motivation for generalizing marginal
 1869 probability to expected multiplicity is that it is a natural generalization. The marginal
 1870 probability of a tuple t is the expectation of a Boolean random variable that is assigned 1
 1871 in every world where tuple t exists and 0 otherwise. For bag-PDBs the multiplicity of a
 1872 query result tuple can be modeled as a natural-number random variable that for a world D
 1873 is assigned the multiplicity of the tuple in D . Thus, a natural generalization of the marginal
 1874 probability (expectation of a Boolean random variable) to bags is the expectation of this
 1875 variable: the tuple's expected multiplicity.

1876 1.176 "N[X] relations are closed under RA+": is this a *definition* of what it means to
 take an RA+ query and evaluate it over an N[X] database, or does this sentence say
 something more? Also, I think it would be clearer to use UCQs in the whole paper instead
 of constantly changing between UCQs, RA+ and SPJU formalisms

1877 To make the paper more accessible and general, we found it better to not use N[X]-DBs.
 1878 While we wanted to use UCQ, we found the choice of \mathcal{RA}^+ to be more amenable to the
 1879 presentation of the paper, and have, as suggested stuck with one query formalism.

1880 There are too many things undefined in from 1.182 to the end of page. 1.182 and in
 Proposition 2.1 N-PDBs are not defined, the function mod is undefined, etc. The article
 should be self-contained: important definitions should be in the article and the appendix
 should only be used to hide proof details. I think it would not take a lot of space
 to properly define the main concepts that you are using, without hiding things in the
 appendix

1881 All material in Sec. 2 that is proof-related is in the appendix, while Sec. 2 (modulo the
 1882 proofs) is itself now self-contained.

1883 1.622 and 1.632-634: so a N-PDB is a PDB where each possible world is an N-database, but
 an N[X]-PDB is not a PDB where each possible world is an N[X]-database... Confusing
 notation

1884
 1885 The text now refers to latter as an N[X]-encoded PDBs.

1886 If you want to be in the setting of bag PDBs, why not consider that the value of the
 variables are integers rather than Boolean? I.e., consider valuations $\nu : X \rightarrow \mathbb{N}$ (or even
 to \mathbb{R} , why not?) instead of $X \rightarrow \{0, 1\}$; this would seem more natural to me than having
 this ad-hoc "mix" of Boolean and non-Boolean setting. If you consider this then your
 "reduced polynomial" trick does not seem to work anymore.

1887
 1888 Our objective is to establish the feasibility of bag-probabilistic databases as compared to
 1889 existing deterministic query processing systems. Accordingly, we take our input model from
 1890 production database systems like Postgresql, Oracle, DB2, SQLServer, etc. (e.g., see ?? on
 1891 ??), where duplicate tuples are represented as independent entities. As a convenient benefit,
 1892 this leads to a direct translation of TIDBs (which are defined over $\{0, 1\}$ inputs). Finally, as
 1893 we mention earlier, an easy generalization exists to encode a bag-PDB in a set-PDB (which
 1894 then allows for bag inputs). Appendix A.

1895 - 1.656 "Thus, from now on we will solely use such vectors...": this seems to be false.
 Moreover you keep switching notation which makes it very hard to read... Sometimes it is
 φ , sometimes it is small w, sometimes it is big W (1.174 or 1.722), sometimes the database
 is $\varphi(D)$, sometimes it is $\varphi_w(D)$, other times it is $D_{[w]}$ (1.671), and so on.

1896 We have made effort to be consistent with the use of notation, following standard usage
1897 whenever possible.

1898 1.658 "we use $\varphi(D)$ to denote the semiring homomorphism $\mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$ that...": I
1899 don't understand why you need a database to extend an assignment to its semiring
1900 homomorphism from $\mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$

1900 φ [25] lifts the valuation function (with kind $\mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$) to databases (i.e., a mapping
1901 from an $\mathbb{N}[\mathbf{X}]$ -DB to a deterministic \mathbb{N} -DB). We note that the main body of the paper no
1902 longer references $\mathbb{N}[\mathbf{X}]$ -DBs, and thus φ is discussed exclusively in the appendix.

1903 Figure 2, K is undefined

1904 We have updated Fig. 1 (originally figure 2) to not need K .

1905 1.178 " Q_t ", 1.189 "Q will denote a polynomial": this is a very poor choice of notation

1906 1.242 "and query Q": is Q a query or a lineage?

1907 We have reserved Q to mean an \mathcal{RA}^+ query and nothing else.

1908 Section 2.1.1: here you are considering set semantics no? Otherwise, one would think that
1909 for bag semantics the annotation of a tuple could be 0 or something of the form $c \times X$,
1910 where X is a variable and c is a natural number

1909 Please see Appendix A for a discussion on going beyond set inputs.

1910 Proof of Proposition A.3. I seems the proof should end after 1.687, since you already
1911 proved everything from the statement of the proposition. I don't understand what it is
1912 that you do after this line.

1912 This text is an informal proof of Proposition 2.5 originally intended to motivate Proposition B.3.
1913 We agree that this should not be part of the proof of the later, and have removed the text.

1914 1.686 "The closure of ... over K-relations": you should give more details on this part. It is
1915 not obvious to me that the relations from 1.646 hold.

1916 The core of this (otherwise trivial) argument, that semiring homomorphisms commute
1917 through queries, was already proven in [25]. We now make this reference explicit.

1918 We apologize for not explaining this in more detail. In universal algebra [24], it has been
1919 proven (the HSP theorem) that for any variety, the set of all structures (called objects) with
1920 a certain signature that obey a set of equational laws, there exists a "most general" object
1921 called the *free object*. The elements of the free objects are equivalence classes (with respect
1922 to the laws of the variety) of symbolic expressions over a set of variables \mathbf{X} that consist of
1923 the operations of the structure. The operations of the free object are combining symbolic
1924 expression using the operation. It has been shown that for any other object K of a variety,
1925 any assignment $\phi : \mathbf{X} \rightarrow K$ uniquely extends to a homomorphism from the free object to
1926 K by substituting variables for based on ϕ in symbolic expression and then evaluating the
1927 resulting expression in K .

1928 Commutative semirings form a variety where $\mathbb{N}[\mathbf{X}]$ is the free object. Thus, for any
1929 polynomial (element of $\mathbb{N}[\mathbf{X}]$), for any assignment $\phi : \mathbf{X} \rightarrow \mathbb{N}$ (also a semiring) there exists a
1930 unique semiring homomorphism $\text{EVAL}_\phi : \mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$. Homomorphisms by definition commute
1931 with the operations of a semiring. Green et al. [?] did prove that semiring homomorphisms
1932 extend to homomorphisms over K-relations (by applying the homomorphism to each tuple's
1933 annotation) and these homomorphisms over K-relations commute with queries.

1934 1.711 "As already noted...": ah? I don't see where you define which subclass of $\mathbb{N}[\mathbf{X}]$ -PDBs define bag version of TIDBs. If this is supposed to be in Section 2.1.1 this is not clear, since the word "bag" does not even appear there (and as already mentioned everything seems to be set semantics in this section). In fact, nowhere in the article can I see a definition of what are bag TIDBs/BIDBs

1935

1936 The new text precisely defines TIDBs (Sec. 1), and the BIDB generalization (Sec. 2.1.1).
 1937 The specific text referenced in this comment has now been moved to the appendix and
 1938 restructured to reference Definition B.2 (which defines an $\mathbb{N}[\mathbf{X}]$ -encoded PDB defined over
 1939 variables \mathbf{X}) and relate it to the formal structure of BIDBs in Sec. 2.1.1.

1940 - 1.707 "the sum of the probabilities of all the tuples in the same block b is 1": no, traditionally it can be less than 1, which means that there could be no tuple in the block.

1941 The reviewer is correct and we have updated our appendix text accordingly.

1942 it is not clear to me how you can go from 1.733 to 1.736, which is sad because this is actually the whole point of this proof. If I understand correctly, in 1.733, $Q(D)(t)$ is the polynomial annotation of t when you use the semantics of Figure 2 with the semiring K being $\mathbb{N}[\mathbf{X}]$, so I don't see how you go from this to 1.736

1943

1944 This result follows from the inner sum looping only over \mathbf{w} s.t. $\psi_{\mathbf{w}}(D_{\mathbb{N}[\mathbf{X}]}) = D$. As
 1945 a consequence of this constraint, we have $Q(D_{\mathbb{N}[\mathbf{X}]})(t)(\mathbf{w}) = Q(D)(t)$. The latter term is
 1946 independent of the summation, and so can be pulled out by distributivity of addition over
 1947 multiplication.

1948 We agree with the reviewer that this could be presented more clearly, and have now split
 1949 the distributivity argument into a separate step.

1950 1.209-227: so you define what is a polynomial and what is the degree of a polynomial (things that everyone knows), but you don't bother explaining what "taking the mod of $Q(X)$ over all polynomials in S " means? This is a bit weird.

1951 Based on this and other reviewer comments, we removed the earlier definition of $\tilde{\Phi}(\mathbf{X})$
 1952 and have defined it in a more ad-hoc manner, as suggested by the reviewers, including the
 1953 comment immediately following.

1954 Definition 2.6: to me, using polynomial long division to define $\tilde{Q}(\mathbf{X})$ seems like a pedantic way of reformulating something similar to Definition 1.3, which was perfectly fine and understandable already! You could just define $\tilde{Q}(\mathbf{X})$ to set all exponents in the SOP that are >1 to 1 and to remove all monomials with variables from the same block, or using Lemma A.4 as a definition?

1955 As alluded to above, we have incorporated the reviewer's suggestion, c.f. Definition 1.3 and
 1956 ??.

Definition 2.14. It is not clear what is the input exactly. Are the query Q and database D fixed? Moreover, I have the impression that your hardness results have nothing to do with lineages and that you don't need them to express your results. I think the problem you should consider is simply the following: Expected Multiplicity Problem: Input: query Q , $N[X]$ -database D , tuple t . Output: expected multiplicity of t in $Q(D)$. Your main hardness result would then look like this: the Expected Multiplicity problem restricted to conjunctive queries is $\#W[1]$ -hard, parameterized by query size. Indeed if I look at the proof, all you need is the queries Q_G^k . The problem is $\#W[1]$ -hard and it should not matter how one tries to solve it: using an approach with lineages or using anything else. Currently it is confusing because you make it look like the problem is hard only when you consider general arithmetic circuits, but your hardness proof has nothing to do with circuits. Moreover, it is not surprising that computing the expected output of an arithmetic circuit is hard: it is trivial, given a CNF ϕ , to build an arithmetic circuit C such that for any valuation ν of the variables the formula ϕ evaluates to True under ν if C evaluates to 1 and the formula ϕ evaluates to False under ν if C evaluates to 0, so this problem is $\#P$ -hard anyways.

The reviewer is correct. Our hardness results are now stated independently of circuits. We note that the hardness result alluded to at the end of the comment above is not applicable in our case since for fixed queries Q , ?? and Problem 1.2 can be solved in polynomial time.

Further, as we point out in Sec. 1 what is new in our hardness results is that we show a query Q^k such that $T_{det}^*(Q^k, D_\Omega)$ is small (linear in $|D_\Omega|$ but solving ?? and Problem 1.2 is hard. We note that it is well-known that one can reduce the problem of counting k -cliques or k -matchings to a query Q for which computing $Q(D_\Omega)$ is $\#W[1]$ -hard. So our contribution to come up with a different reduction from counting k -matchings so that the hardness manifests itself in the probabilistic computing part of our problem.

Section 3.3. It seems to me the important part of this section is not so much the fact that we have fixed values of p but that the query is now fixed and that you are looking at the fine-grained complexity. If what you really cared about was having fixed value of p , then the result of this section should be exactly like the one in Theorem 3.4, but starting with "fix p ". So something like "Fix p . Computing \tilde{Q}_G^k for arbitrary G is $\#W1$ -hard".

We agree with the reviewer that the result on fixed value of p is mostly of (narrow) theoretical interest. We have added a discussion summarizing the reviewer's point above below Theorem 3.7.

General remark: The story of the paper I think should be this: we can always compute the expected multiplicity for a UCQ Q and $N[X]$ -database D and tuple t by first computing the lineage in SOP form and then using linearity of expectation, which gives an upper bound of (roughly) $O(|D|^{|Q|})$. We show that this exponential dependence in $|Q|$ is unavoidable by proving that this problem is $\#W1$ hard parameterized by $|Q|$ (which implies that we cannot solve it in $f(|Q|)|D|^c$). Furthermore we obtain fine-grained superlinear lower bounds for a fix conjunctive query Q . (Observe how up to here, there is no need to talk about lineages at all). We then obtain an approximation algorithm for this problem for [this class of queries] and [that class of bag PDBs] with [that running time (Q,D)]. The method is to first compute the lineage as an arithmetic circuit C in [this running time (Q,D)], and then from the arithmetic circuit C compute in [running time(C)] an approximation of its expected output. Currently I don't understand to which queries your approximation algorithm can be applied (see later comments).

We have restructured Sec. 1 to more or less follow the reviewer's outline above. The only deviation is that we still introduce lineage polynomials. We do this because the polynomial

1974 view is very helpful in the proofs of our hardness result (in addition to the obvious relevance
1975 for the approximation algorithm). We have also clarified that our approximation result
1976 applied to all \mathcal{RA}^+ queries (see Corollary 4.9).

1977 1.381: Here again, I think it would be simpler to consider that the input of the problem is
the query, the database and a tuple and claim that you can compute an approximation
of the expected multiplicity in linear time. The algo is to first compute the lineage as
an arithmetic circuit, and then to use what you currently use (which could be put in a
lemma or in a proposition).

1978 We have implemented the above overview in Sec. 1 when we move from ?? to Problem 1.6.
1979 For the approximation algorithm we focus on Problem 1.6, which still takes a circuit as an
1980 input.

1981 Definition 4.2: would you mind giving an intuition of what this is? It is not OK to
define something and just tell the reader to refer the appendix to understand what this is
and why this is needed; the article should be understandable without having to look at
the appendix. It is simply something that gives the coefficient of each monomial in the
reduced polynomial?

1982 We have provided an example in directly after Definition 4.1 as well as a sentence pointing
1983 out why this definitions is useful.

1984 - 1.409: how does it matter that the circuit C is the lineage of a UCQ? Doesn't this work
for any arithmetic circuit?

1985 The reviewer is correct that the earlier Theorem 4.9 works for any circuit (this result is now
1986 in the appendix).

1987 1.411: what are $|C|^2(1, \dots, 1)$ and $|C|(1, \dots, 1)$?

1988 We clarify this overloaded notation immediately after Definition 4.2.

1989 Sometimes you consider UCQs, sometimes \mathcal{RA}^+ queries. I think it would be simpler if
you stick to one formalism (probably UCQs is cleaner?)

1990 As alluded to previously, we have followed the reviewer's suggestion and have found \mathcal{RA}^+
1991 queries to be most amenable for this work.

1992 1.432 what is an FAQ query?

1993 We actually no longer need that result since Lemma 4.8 now has a bound on $|C|(1, \dots, 1)$ in
1994 terms of $\text{DEPTH}(C)$ and the latter is used in Corollary 4.9 for all \mathcal{RA}^+ queries. Please see
1995 Lemma 4.8 and the followup discussion for more on this.

1996 Generally speaking, I think I don't understand much about Section 4, and the
convolutedness of the appendix does not help to understand. I don't even see in which
result you get a linear runtime and to which queries the linear runtime applies. Somewhere
there should be a corollary that clearly states a linear time approximation algorithm for
some queries.

1997 We have re-organized sec:algo to address the above comments as follows:
1998

- 1999 ■ We now start off Sec. 4.2 with the algorithm idea.
- 2000 ■ We give a quick overview of how the claimed runtime follows from the algorithm idea
2001 mentioned above.
- 2002 ■ Added Corollary 4.9 that clearly states that we get an $O(T_{det}^*(Q, D_\Omega))$ for all \mathcal{RA}^+ queries
2003 Q .

2004 In section 5, it seems you are arguing that we can compute lineages as arithmetic circuits
at the same time as we would be running an ordinary query evaluation plan. How is that
different from using the relations in Figure 2 for computing the lineage?

2005 There is not a major difference between the two. This observation has persuaded us to
eliminate $\mathbb{N}[\mathbf{X}]$ -DB query evaluation and have only an algorithm for lineage.

2007 We have also re-organized the earlier Section 5 and moved the definition of $T_{det}^*(\cdot)$ (earlier
denoted as $\mathbf{cost}(\cdot)$) to Sec. 2.3 and moved the rest of the material to the appendix.

2009 1.679 where do you use $\max(D_i)$ later in the proof?

2010

2011 Thank you. This reference was unnecessary and has been removed.

2012 1.688 That sentence is hard to parse, consider reformulating it

2013

2014 As the reviewer notes above, this paragraph is unnecessary and we have removed it.

2015 it seems you are defining $\mathbb{N}[\mathbf{X}]$ -PDB at two places in the appendix: once near 1.632, and
another time near 1.652

2016

2017 Thank you. The latter definition has been removed.

2018 1.3 Reviewer 2

2019 First, the paper should state rigorously the problem definition. There are three well-known
definitions in database theory: data complexity, combined complexity, and parameterized
complexity. If I understand correctly, Theorem 3.4 refers to the parameterized complexity,
Theorem 3.6 refers to the data complexity (of a fixed query), while the positive results in
Sec. 4 (e.g. Th. 4.8) introduce yet another notion of complexity, which requires discussion.

2020 We have addressed the concerns in rewriting the entirety of Sec. 1, explicitly mentioning
2021 complexity metrics considered, while forming a series of problem statements that describe
2022 the exact problem we are considering, and the complexity metrics considered. We have
2023 also adjusted the phrasing of the said theorems and definitions to eliminate the ambiguity
2024 described.

2025 The problem definition is supposed to be in Definition 2.14, but this definition is sloppy. It
states that the input to the problem is a circuit C: but then, what is the role of the PDB
and the query Q? Currently Definition 2.14 reads as follows: "Given a circuit C defining
some polynomial $Q(\mathbf{X})$, compute $E[Q(\mathbf{W})]$ ", and, thus, the PDB and the query play no
role at all. All results in Section 4 seem to assume this simplified version of Definition
2.14. On the other hand, if one interprets the definition in the traditional framework
of data complexity (Q is fixed, the inputs are D and C) then the problem is solvable in
PTIME (and there is no need for C), since $E[Q(\mathbf{W})]$ is the sum of expectations of the
monomials in Q (this is mentioned in Example 1.2).

2026 We have rephrased Definition 2.9 to qualify data complexity. The paper (especially in Sec. 1)
2027 builds up the fact that we aren't stopping at polynomial time, but exploring parameterized
2028 complexity and fine grained analysis (as the reviewer aptly noted in the first comment).

2029 Second, Definition 2.6 of Reduced BIDD polynomials is simply wrong. It uses "mod" of two
multivariate polynomials, but "mod" doesn't exist for multivariate polynomials...Either
state Definition 2.6 directly, in an ad-hoc manner (which seems doable), or do a more
thorough job grounding it in the ring of multivariate polynomials and its ideals.

2030 The reviewer is correct in their comment on the "mod" part– we apologize for the error. We
 2031 have implemented the reviewer's ad-hoc suggestion in light of Reviewer 1's similar suggestions.

2032 the paper uses three notations (UCQ, RA+, SPJU) for the same thing, and never defines
 2033 formally any of them.

2033 We have chosen \mathcal{RA}^+ for consistent use throughout the paper. We have included ?? on ??
 2034 for an explicit definition of \mathcal{RA}^+ queries.

2035 G^ℓ is used in Lemma 3.8 but defined only in the Appendix (Def. B.2), without even a
 2036 forward pointer. This is a major omission: Lemma 3.8 is a key step for a key result, but
 2037 it is impossible to read.

2036 We have fixed this mistake. Unfortunately, because of the changes in the paper (especially
 2037 expanding on Sec. 1), the earlier Lemma 3.8 had to be moved to the appendix.

2038 Definition 2.7. "valid worlds η ". This is confusing. A "possible world" is an element of Ω :
 2039 this is not stated explicitly in the paper, but it is implicit on line 163, so I assumed that
 2040 possible worlds refer to elements of Ω . If I assumed correctly, then calling η a "world" in
 2041 Def. 2.7 is misleading, because η is not an element of Ω . More, it is unclear to me why
 2042 this definition is needed: it is used right below, in Lemma 2.8, but that lemma seems to
 2043 continue to hold even if w is not restricted.

2040 We agree with the reviewer that this notation is confusing; η is meant to cope with the
 2041 fact that tuples from the same group in a BIDB can not co-exist, even though our $\{0, 1\}$ -input
 2042 vectors can encode such worlds. We now address this constraint by embedding it directly
 2043 into the reduced polynomial with ??.

2044 line 305: please define what is an "occurrence of H in G". It could mean: a homomorphic
 2045 image, a subgraph of G isomorphic to H, an induced subgraph of G isomorphic to H, or
 2046 maybe something else.

2045 We agree with the reviewer's suggestion and have rephrased the wording to be clear. Please
 2046 see the beginning of Sec. 3.1.

2047 If the proofs are given in the appendix, please say so. Lemmas 3.5 and 3.8 are stated
 2048 without any mention, and one has to guess whether they are obvious, or proven somewhere
 2049 else. On this note: I found Lemma 3.5 quite easy, since the number of k-matching is
 2050 the coefficient of the leading monomial (of degree $2k$) in $Q^k(p, p, \dots, p)$, while Lemma 3.8
 2051 appears much harder. It would help to briefly mention this in the main body of the paper.

2048 We have implemented the reviewer's suggestion. Please see the last sentence of Sec. 1 (as
 2049 well as the expanded discussion on the hardness result in the **Overview of our Techniques**
 2050 part of Sec. 1).

2051 line 177: what is $\Omega_{\mathbb{N}[\mathbf{X}]}$?

2052 We have eliminated the use of $\mathbb{N}[\mathbf{X}]$ -DBs in the paper proper, using them only when necessary
 2053 in the proofs of the appendix.

2054 line 217. The polynomial $X^2 + 2XY + Y^2$ is a poor choice to illustrate the degree. There
 2055 are two standard definitions of the degree of a multivariate polynomial, and one has to
 2056 always clarify which one is meant. One definition is the total degree (which is Def. 2.3 in
 2057 the paper), the other is the maximum degree of any single variable. It is nice that you
 2058 are trying to clarify for the reader which definition you are using, but the polynomial
 2059 $X^2 + 2XY + Y^2$ is worst choice, since here the two coincide.

2055 We have adjusted the example to account for the reviewer's correct observation.

line 220. "we consider only finite degree polynomials". This is a surprise. Polynomials, by definition, are of finite degree; there are extensions (I'm aware of powerseries, maybe you have other extensions in mind), but they are usually not called polynomials, plus, nothing in the paper so far suggests that it might refer to those extensions.

We have removed the redundant terminology the reviewer has pointed out, and refined the discussion surrounding (and including) Eq. (1) to be explicit to the novice reader that polynomials are by definition of finite degree.

"Note that our hardness results even hold for the expression trees". At this point we haven't seen the hardness results, nor their proofs, and we don't know what expression trees are. It's unclear what we can note.

Our hardness results are now stated independently of circuits so the above statement no longer appears in the paper.

paragraph at the top of pp.10 is confusing. My guess is that it is trying to say: "there exists a query Q, such that, for each graph G, there exists a database D s.t. the lineage of Q on D is the polynomial Q_G ."

Our revision has eliminated this statement.

I.4 Reviewer 3

The overall study is then extended to a multiplicative approximation algorithm for the expectation of polynomial circuits in linear time in the size of the polynomial. It was much harder to read this part, and I found the examples and flow in the appendix quite helpful. I suggest to include these examples into the body of the paper.

In our revision we expanded on Sec. 1 to give a better overview of the problems we are considering in this paper. This meant we had to cut out material in later sections, which unfortunately meant we did not have space in Sec. 4 to include any examples that the reviewer suggested above. However, we have tried to make Sec. 4 more readable as a whole.

While ApproximateQ is linear in the size of the circuit, it is quadratic in epsilon and so we need quadratically many samples for the desired accuracy – overall runtime is not linear therefore and it may be better to elaborate this. It may also be helpful to comment on how this relates to Karp, Luby, Madras algorithm [1] for #DNF which is also quadratic in epsilon.

In Problem 1.5 we note explicitly that we care about linear dependence on $T_{det}^*((Q, D_\Omega))$ and do not care about the exact dependence on ϵ . While it would be nice to design an approximation algorithm that is linear in $1/\epsilon$ as well, we believe it is out of scope of this initial work.

The coverage of related work is adequate. Fink et. al seems as the closest related work to me and I would appreciate a more elaborate comparison with this paper. My understanding is that Fink et. al considers exact evaluation only and focuses on knowledge compilation techniques based on decompositions. They also note that "Expected values can lead to unintuitive query answers, for instance when data values and their probabilities follow skewed and non-aligned distributions" attributed to [2]. Does this apply to the current work? Can you please comment on this?

The work is indeed quite close to our own. It targets a broader class of queries (aggregates include COUNT/SUM/MIN/MAX, rather than our more narrow focus on COUNT), but has significantly less general theoretical results. Most notably, their proof of linear runtime in the

2082 size of the input polynomial is based on a tree-based encoding of the polynomial. Tree-based
 2083 representation representation (and hence the Fink et. al. algorithm's runtime) is, as we note
 2084 several times, superlinear in $T_{det}^*(Q, D_\Omega)$. This result is also limited to a specific class of
 2085 (hierarchical) queries, devolving to exponential time (as in [20]) in general. By contrast, our
 2086 results apply to all of \mathcal{RA}^+ . Our revised related work section now addresses both points.

2087 I assume the authors focus on parameterized complexity throughout the paper, and even
 this is not stated unambiguously. The authors should make an extra effort to make the
 paper more accessible by using the explanations and examples from the appendix in the
 body of the paper. It is also important to highlight the differences with the complexity of
 standard query evaluation over PDBs.

2088 Our revision has focused on explicitly mentioning the complexity metrics we are interested
 2089 in. This can be seen in e.g. Sec. 1 and formal statement (theorems, lemmas, etc.), which
 2090 have been rewritten to eliminate ambiguities. We have also taken pains to be promote
 2091 accessibility, keeping the paper self-contained, and using examples for difficult or potentially
 2092 unclear concepts. This can be seen in e.g. eliminating unnecessary machinery (e.g. $\mathbb{N}[\mathbf{X}]$ -
 2093 DB machinery from the paper proper), providing/modifying examples (c.f. Definition 4.1,
 2094 Definition 4.4), and ensuring consistency in notational use, e.g. using one query evaluation
 2095 formalism (\mathcal{RA}^+).

2096 We decided to focus on beefing up Sec. 1 and cleaning up definitions of problems, which
 2097 unfortunately meant we ran out of space to bring back examples from the appendix (especially
 2098 into Sec. 4).

2099 1.5 Reviewer 4

2100 I wonder whether the writing could be revisited to give the reader a better overview of the
 technical challenges, motivation, and the high level ideas of the algorithm and hardness
 results. The current exposition seems slightly too tailored for the expert in probabilistic
 databases rather than the average ICDT attendee. Also the current exposition is structured
 such that the reader needs to get through quite a few definitions and technical lemmas
 until they get to the new ideas in the paper.

2101 We have (as noted throughout this section) revised the writing to provide precision and
 2102 clarity to the problem we explore as well as the results we obtain. Part of this revision
 2103 was a complete rewriting of Sec. 1 where we sought to be extremely precise in language
 2104 and through a series of problem statements to help the reader navigate and understand the
 2105 problem we explore as well as how we have gone about exploring that problem coupled with
 2106 the validity of the exploration strategy. We have simultaneously sought to make the paper
 2107 more accessible by assuming the average ICDT attendee and defining or explaining concepts
 2108 that might not be known to them. Finally, we have expanded on the **Overview of our**
 2109 **Techniques** part of Sec. 1 to provide more intuition on how we prove our lower and upper
 2110 bounds.

2111

2112

2113

2114

2115

2116