

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

Su Feng ✉

Illinois Institute of Technology, Chicago, USA

Boris Glavic ✉

Illinois Institute of Technology, USA

Aaron Huber ✉

University at Buffalo, USA

Oliver Kennedy ✉

University at Buffalo, USA

Atri Rudra ✉

University at Buffalo, USA

## Abstract

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

2012 ACM Subject Classification Information systems → Incomplete data

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

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

## 1 Introduction

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

breaks for BIDBs ← the 1-TIDB reduction, no?  
bag TIDBs understood  
we are defining Bag TIDBs to be this. frame if that way?  
not just a special case

© 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. 23:1–23:63

Leibniz International Proceedings in Informatics

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

Worthy noting?  
 Can't fix the poly bound to e.g.  $n^k$ , since det runtime varies. We are giving a strictly tighter bound!

Might be beneficial to add a note that  $RA^+$  is our  $Q$  lang &  $RA^+$  defines order of execution explicitly

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

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

46  
 47 ▶ **Problem 1.1.** Given a  $c$ -TIDB  $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ ,  $RA^+$  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  $T_{det}(Q, \bar{D}, c) = Q(\bar{D})$  for arbitrary query  $Q$ , deterministic database  $\bar{D}$ , and multiplicity  
 67 bound  $c$ . Let  $T_{det}^*(Q, \bar{D}, c) = \min_{Q': Q' \equiv Q} T_{det}(Q', \bar{D}, c)$  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

69 **Table 1** Our lower bounds for a specific hard query  $Q$  parameterized by  $k$ . For  $\mathcal{D} =$   
 70  $\{\{0, \dots, c\}^D, \mathcal{P}\}$  those with 'Multiple' in the second column need the algorithm to be able to  
 71 handle multiple  $\mathcal{P}$  (for a given  $D$ ). The last column states the hardness assumptions that imply the  
 72 lower bounds in the first column ( $\epsilon_0, C_0, c_0$  are constants that are independent of  $k$ ).

Using a local ref implies that this paper makes it. Suggest adding a cite? eg. [xx] (cor. 3.2)

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

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

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

emph?  
~~no~~  
~~no~~  
 not defined  
 maybe add a clarification above that  $Q^*$  is optimal runtime?





```

112 SELECT 1 FROM OnTime a, Route r, OnTime b
113 WHERE a.city = r.city1 AND b.city = r.city2
114

```

at least an intuition

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  $\Phi^2(A, B, C, E, X, Y, Z)$

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

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

Specifically  
 Introduce the reduced poly at this point before moving on to limitations  
 Also intro  $\tilde{\Phi}(p)$  here

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

129 Considering again our example,

$$\Phi_R^{(AXB)^2}(A, X, B) = \Phi^{(AXB)^2} \left( \sum_{j_1 \in [c]} j_1 A_{j_1}, \sum_{j_2 \in [c]} j_2 X_{j_2}, \sum_{j_3 \in [c]} j_3 B_{j_3} \right)$$

$$= \left( \sum_{j_1 \in [c]} j_1 A_{j_1} \right)^2 \left( \sum_{j_2 \in [c]} j_2 X_{j_2} \right)^2 \left( \sum_{j_3 \in [c]} j_3 B_{j_3} \right)^2$$

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

138 **Definition 1.3.** For any polynomial  $\Phi((X_t)_{t \in D})$  define the reformulated polynomial  
 139  $\Phi_R((X_{t,j})_{t \in D, j \in [c]})$  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  
 140 reduced polynomial  $\tilde{\Phi}((X_{t,j})_{t \in D, j \in [c]})$  to be the polynomial resulting from converting  $\Phi_R$   
 141 into the standard monomial basis (SMB),<sup>1</sup> removing all monomials containing the term  
 142  $X_{t,j} X_{t,j'}$  for  $t \in D, j \neq j' \in [c]$ , and setting all variable exponents  $e > 1$  to 1.

143 Continuing with the example  $\Phi^2(A, B, C, E, X_1, X_2, Y, Z)$  to save clutter we i) do not show  
 144 the full expansion for variables with greatest multiplicity = 1 since e.g. for variable  $A$ ,  
 145 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

<sup>1</sup> 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.

Awk definition  
 $\Phi_R$  defined by construction but missing structure eg  $\Phi_R: X^c \rightarrow \mathbb{N}$   
 or here set  $X_{t,j}$  in what?

It's not the same as this is why

Slightly inconsistent notation  $A_j$  vs  $X_{A,j}$  above

Breaks flow put in footnote?

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  $> 2$ .

149

$$150 \quad \tilde{\Phi}^2(A, B, C, E, X_1, X_2, Y, Z) =$$

$$151 \quad A \left( \sum_{j \in [c]} j^2 X_j \right) B + BYE + BZC + 2A \left( \sum_{j \in [c]} j^2 X_j \right) BYE + 2A \left( \sum_{j \in [c]} j^2 X_j \right) BZC + 2BYE ZC =$$

$$152 \quad ABX_1 + AB(2)^2 X_2 + BYE + BZC + 2AX_1 BYE + 2A(2)^2 X_2 BYE + 2AX_1 BZC + 2A(2)^2 X_2 BZC + 2BYE ZC.$$

153

154 Note that we have argued that for our specific example the expectation that we want is  
 155  $\tilde{\Phi}^2(\Pr(A=1), \Pr(B=1), \Pr(C=1), \Pr(E=1), \Pr(X_1=1), \Pr(X_2=1), \Pr(Y=1), \Pr(Z=1))$ .  
 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  $\mathbf{p} = ((p_{t,j})_{t \in D, j \in [c]})$ .

Here  $\mathbf{p}$  is no longer a product distribution. Triggers BS detection suggest call-out or similar

159 **1.2 Our Techniques**

160 **Lower Bound Proof Techniques.** Our main hardness result shows that computing Problem 1.1  
 161 is  $\#W[1]$  – 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  $\tilde{\Phi}$  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  $\Phi$  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

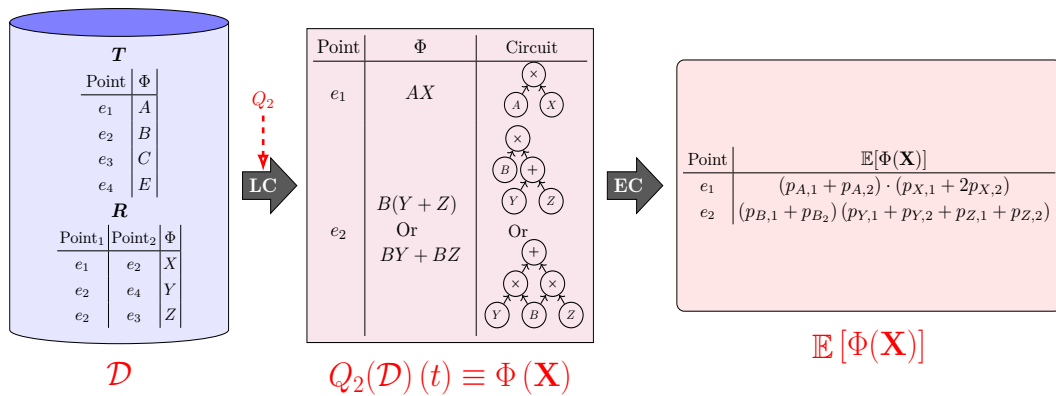


Figure 2 Intensional Query Evaluation Model ( $Q_2 = \pi_{\text{Point}}(T \bowtie_{\text{Point}=\text{Point}_1} R)$  and  $c = 2$ ).

188 **C** with one sink per tuple representing the tuple’s lineage). A key insight of this paper is  
 189 that the representation of **C** matters. For example, if we insist that **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) circuits<sup>2</sup> 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  $\mathbf{C}$ . 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 
$$\Phi^2(\mathbf{p}) = p_A^2 p_X^2 p_B^2 + p_B^2 p_Y^2 p_E^2 + p_B^2 p_Z^2 p_C^2 + 2p_A p_X p_B^2 p_Y p_E + 2p_A p_X p_B^2 p_Z p_C + 2p_B^2 p_Y p_E p_Z p_C$$
  
 210 
$$\leq p_A p_X p_B + p_B p_Y p_E + p_B p_Z p_C + 2p_A p_X p_B p_Y p_E + 2p_A p_X p_B p_Z p_C + 2p_B p_Y p_E p_Z p_C = \tilde{\Phi}(\mathbf{p})$$

212 If we assume that all seven probability values are at least  $p_0 > 0$ , we get that  $\Phi^2(\mathbf{p})$  is  
 213 in the range  $[(p_0)^3 \cdot \tilde{\Phi}(\mathbf{p}), \tilde{\Phi}(\mathbf{p})]$ . 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

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

~~Emph that this is not a good approx?~~

216 equivalent SMB representation of  $\Phi$  (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

### 2.1 Polynomial Definition and Terminology

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

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  $\Phi$ .

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

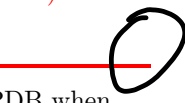
243 As an example, the degree of the polynomial  $X^2 + 2XY^2 + Y^2$  is 3. Product terms in lineage  
244 arise only from join operations (Fig. 1), so intuitively, the degree of a lineage polynomial  
245 is analogous to the largest number of joins needed to produce a result tuple. We call a  
246 polynomial  $\Phi(\mathbf{X})$  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

249 An *incomplete database*  $\Omega$  is a set of deterministic databases  $\omega$  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



Space & params  
break flow

where  $c_{\mathbf{d}} \in \mathbb{N}$

Bas-

emph that here,  $\Omega$  is a set of sets of tuples

Example needed

Discuss relationship to set TIDBs  
e.g. when  $c=1$ , this is analogous to a set-TIDB

concurrently have non zero multiplicities

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  $\Omega$  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  $c$ -BIBD 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** ( $c$ -TIBD reduction). Given  $c$ -TIBD  $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ , let  $\mathcal{D}' =$   
262  $(\Omega, \mathcal{P}')$  be the 1-BIBD 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 BIBD disjoint requirement.

265 For the  $c$ -TIBD  $\mathcal{D}$ , each  $X_t \in [c]$ , while in the reduced 1-BIBD  $\mathcal{D}'$ , each  $X_{t,j} \in \{0, 1\}$ .  
266 Hence, in the setting of 1-BIBD, 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-BIBD,  $\Phi[R, D', t]$  is of the same structure as the reformulated polynomial  $\Phi_R$  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-BIBD  
271 polynomials:  $\mathbb{E}_{\mathbf{w} \sim \mathcal{P}'}[\Phi(\mathbf{W})] = \tilde{\Phi}(\mathbf{p})$ .

**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  $\mathcal{P}' \equiv \mathcal{P}''$ .

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

281 **Corollary 2.4.** If  $\Phi$  is a BIBD-lineage polynomial already in SMB, then the expectation of  
282  $\Phi$ , 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  $\Omega$  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:  $Pr[\omega \in \Omega] = \sum_{\omega' \in \Omega, Q(\omega') = \omega} Pr[\omega']$ .

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

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

awk: shouldn't this be a c-BIBD? new name? r-BIBD

what's the multiplicity of a block

The other way to present this is to emph that the 'c' is the domain of X

footnotes look like exponents put this e.g. on "modeled"

the set of

~~relationship~~

subscript?

has the term "world assignments" been defined yet?



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  $\Omega$ ):  $\mathbb{E}_{\mathbf{D} \sim \mathcal{P}}[Q(\mathbf{D})(t)] = \mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$ .

296 A formal proof of Proposition 2.5 is given in Appendix B.3.<sup>4</sup> 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.

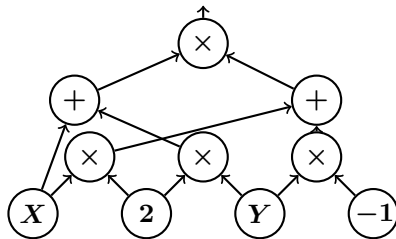
305 ► **Definition 2.6** (Circuit). A circuit  $C$  is a Directed Acyclic Graph (DAG) whose source  
 306 gates (in degree of 0) consist of elements in either  $\mathbb{N}$  or  $\mathbf{X}$ . For each result tuple there exists  
 307 one sink gate. The internal gates have binary input and are either 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.  $C_L$  ( $C_R$ ) denotes the left (right) input of  $C$ .

**Aaron says:** Does the following matter, i.e., does it point anything out special for our 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  $C$ . 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  $\Phi$ . Note that each  
 316 circuit  $C$  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.



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

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

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

<sup>4</sup> 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.

## 23:10 Bag PDB Queries

328  $\mathbf{C}$  need not encode  $\Phi(\mathbf{X})$  in the same, default SMB representation. For instance,  $\mathbf{C}$  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  $\text{POLY}(\mathbf{C}) = \Phi(\mathbf{X})$  is always the equivalent SMB representation.

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

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  $\mathbf{C}$ ,  $\text{POLY}(\mathbf{C}) = \Phi(\mathbf{X})$ .

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  $D_{\overline{\Omega}}$ . Fix an  $\mathcal{RA}^+$  query  $Q$  and a  
338 result tuple  $t$ . The EXPECTED RESULT MULTIPLICITY PROBLEM is defined as follows:

340 **Input:**  $\mathbf{C} \in \text{CSet}(\Phi(\mathbf{X}))$  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,  $T_{\text{join}}(R_1, \dots, R_n)$  is linear in the AGM bound [6].  
349 Our cost model for general query evaluation follows from the join cost:

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

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  $D_{\overline{\Omega}}$ , there  
356 exists a circuit  $\mathbf{C}^*$  such that  $T_{LC}(Q, D_{\overline{\Omega}}, \mathbf{C}^*)$  and  $|\mathbf{C}^*|$  are both  $O(T_{\text{det}}(Q, D_{\overline{\Omega}}))$ . Recall we  
357 assumed these two bounds when we moved from Problem 1.5 to Problem 1.6.

## 358 3 Hardness of Exact Computation

**Aaron says:** If anything need be changed in Sec. 3, it would only be in the following  
(opening) paragraph.

359 In this section, we will prove the hardness results claimed in Table 1 for a specific (family)  
360 of hard instance  $(Q, \mathcal{D})$  for Problem 1.2 where  $\mathcal{D}$  is a TIDB. Note that this implies hardness  
361 for BIDBs and general bag-PDB, showing Problem 1.2 cannot be done in  $O(T_{\text{det}}^*(Q, D_{\overline{\Omega}}))$   
362 runtime.  
363

with framing as a lower-bound on the cost?

364 **3.1 Preliminaries**

365 Our hardness results are based on (exactly) counting the number of (not necessarily induced)  
 366 subgraphs in  $G$  isomorphic to  $H$ . Let  $\#(G, H)$  denote this quantity. We can think of  $H$   
 367 as being of constant size and  $G$  as growing. In particular, we will consider the problems of  
 368 computing the following counts (given  $G$  in its adjacency list representation):  $\#(G, \triangle)$  (the  
 369 number of triangles),  $\#(G, \mathfrak{M}_3)$  (the number of 3-matchings), and the latter's generalization  
 370  $\#(G, \mathfrak{M}_k)$  (the number of  $k$ -matchings). We use  $T_{match}(k, G)$  to denote the optimal  
 371 runtime of computing  $\#(G, \mathfrak{M}_k)$  exactly. Our hardness results in Sec. 3.2 are based on  
 372 the following hardness results/conjectures:

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

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

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

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

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

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

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

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

393 
$$\Phi_G^k(X_1, \dots, X_n) = \left( \sum_{(i,j) \in E} X_i \cdot X_j \right)^k$$

~~Table Names?~~

394 Returning to Fig. 2, it is easy to see that  $\Phi_G^k(\mathbf{X})$  is the lineage polynomial corresponding to  
 395 the query that generalizes our example query from Sec. 1. Let us alias

```
396 SELECT 1 FROM OnTime a, Route r, OnTime b
397 WHERE a.city = r.city1 AND b.city = r.city2
398
```

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

```
401 SELECT COUNT(*) FROM R_1 JOIN R_2 JOIN...JOIN R_k
402
403
```

404 Further, the PDB instance generalizes the one in Fig. 2 as follows. Relation *OnTime* has  
 405  $n$  tuples corresponding to each vertex for  $i$  in  $[n]$ , each with probability  $p_i$  and *Route* has

Why  $R_i$ ? The self-join is needed, no?

406 tuples corresponding to the edges  $E$  (each with probability of 1).<sup>5</sup> In other words, for this  
 407 instance  $D_{\bar{\Omega}}$  contains the set of  $n$  unary tuples in  $OnTime$  (which corresponds to  $V$ ) and  $m$   
 408 binary tuples in  $Route$  (which corresponds to  $E$ ). Note that this implies that  $\Phi_G^k$  is indeed a  
 409 TIDB-lineage polynomial.

410 Next, we note that the runtime for answering  $Q^k$  on deterministic database  $D_{\bar{\Omega}}$ , 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  $D_{\bar{\Omega}}$  be as defined above. Then  $T_{det}^*(Q^k, D_{\bar{\Omega}})$  is  $O(km)$ .*

## 413 3.2 Multiple Distinct $p$ Values

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

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

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

## 428 3.3 Single $p$ value

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

433 ► **Theorem 3.7.** *Fix  $p \in (0, 1)$ . Then assuming Conjecture 3.3 is true, any algorithm that  
 434 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  
 435  $\epsilon_0$  is as defined in Conjecture 3.3.*

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

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

~~Table Names~~

441 **4**  $1 \pm \epsilon$  Approximation Algorithm

442 **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.

443 In Sec. 3, we showed that the answer to Problem 1.6 is no. With this result, we now  
444 design an approximation algorithm for our problem that runs in  $O(|C|)$  for a very broad  
class of circuits (see the discussion after Lemma 4.8 for more). The following approximation  
445 algorithm applies to ~~BIDB~~ lineage polynomials (over  $\mathcal{RA}^+$  queries) though our bounds are  
446 ~~not~~ meaningful for a non-trivial subclass of queries over ~~BIDBs~~ ~~ch~~ contains all queries on  
447 TIDBs as well as the queries of the PDBench benchmark [1]. All proofs and pseudocode can  
448 be found in Appendix D.

449 **Aaron says:** We are going to have to rework  $\gamma$  in this section, as well as the proof for  
our result.

451 **4.1 Preliminaries and some more notation**

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

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

$$457 E(C) = \begin{cases} E(C_L) \circ E(C_R) & \text{if } C.type = + \\ \{(v_L \cup v_R, c_L \cdot c_R) \mid (v_L, c_L) \in E(C_L), (v_R, c_R) \in E(C_R)\} & \text{if } C.type = \times \\ List[(\emptyset, C.val)] & \text{if } C.type = NUM \\ List[(\{C.val\}, 1)] & \text{if } C.type = VAR. \end{cases}$$

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

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

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

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

468 **Definition 4.4** ( $DEG(\cdot)$ ). <sup>6</sup>  $DEG(C)$  is defined recursively as follows:

$$469 DEG(C) = \begin{cases} \max(DEG(C_L), DEG(C_R)) & \text{if } C.type = + \\ DEG(C_L) + DEG(C_R) + 1 & \text{if } C.type = \times \\ 1 & \text{if } C.type = VAR \\ 0 & \text{otherwise.} \end{cases}$$

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

$\frac{DEG(C)+1}{2}$

~~Di + BIDB?~~

~~use ch~~

23:14 Bag PDB Queries

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

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

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

479 ► **Definition 4.6** (Parameter  $\gamma$ ). Given a BIDB circuit  $\mathcal{C}$  define

**Aaron says:** Technically,  $v$  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 ( $v_m$ )

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

482 
$$\gamma(\mathcal{C}) = \frac{\sum_{(v,c) \in E(\mathcal{C})} |c| \cdot \mathbb{1}_{\neg \text{ISIND}(v_m)}}{|\mathcal{C}|(1, \dots, 1)}$$

483 **4.2 Our main result**

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

487 
$$\tilde{\Phi}(p_1, \dots, p_n) = \sum_{(v,c) \in E(\mathcal{C})} \mathbb{1}_{\text{ISIND}(v_m)} \cdot |c| \cdot \prod_{X_i \in v} p_i \tag{2}$$

*Handwritten:* ~~Has ISIND been defined formally by this point?~~

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

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

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

497 
$$\Pr \left( \left| \mathcal{E} - \tilde{\Phi}(p_1, \dots, p_n) \right| > \epsilon' \cdot \tilde{\Phi}(p_1, \dots, p_n) \right) \leq \delta \tag{3}$$

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

*Handwritten:* bin BIDB?

498 can be computed in time

$$499 \quad O \left( \left( \text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta} \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C})}{(\epsilon')^2 \cdot (1-\gamma)^2 \cdot p_0^{2k}} \right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C}))) \right). \quad (4)$$

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

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

504 We briefly connect the runtime in Eq. (4) to the algorithm outline earlier (where we  
 505 ignore the dependence on  $\overline{\mathcal{M}}(\cdot, \cdot)$ , which is needed to handle the cost of arithmetic operations  
 506 over integers). The  $\text{SIZE}(\mathcal{C})$  comes from the time take to run ONEPASS once (ONEPASS  
 507 essentially computes  $|\mathcal{C}|(1, \dots, 1)$  using the natural circuit evaluation algorithm on  $\mathcal{C}$ ). We  
 508 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  $O(k)$   
 509 random sink to source paths in  $\mathcal{C}$  all of which by definition have length at most  $\text{DEPTH}(\mathcal{C})$ ).

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

511 **► Lemma 4.8.** For any <sup>511</sup>BIDB circuit  $\mathcal{C}$  with  $\text{DEG}(\mathcal{C}) = k$ , we have  $|\mathcal{C}|(1, \dots, 1) \leq 2^{2^k \cdot \text{DEPTH}(\mathcal{C})}$ .  
 512 Further, if  $\mathcal{C}$  is a tree, then we have  $|\mathcal{C}|(1, \dots, 1) \leq \text{SIZE}(\mathcal{C})^{O(k)}$ .

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

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

521 **► Corollary 4.9.** Let  $Q$  be an  $\mathcal{RA}^+$  query and  $\mathcal{D}$  be an <sup>521</sup>BIDB with  $p_0 > 0$  and  $\gamma < 1$  (where  
 522  $p_0, \gamma$  as in Theorem 4.7) are absolute constants. Let  $\Phi(\mathbf{X}) = \Phi[Q, D, t]$  for any result tuple  
 523  $t$  with  $\text{deg}(\Phi) = k$ . Then one can compute an approximation satisfying Eq. (3) in time  
 524  $O_{k, |Q|, \epsilon', \delta} (T_{det}^*(Q, D_{\overline{\Omega}}))$  (given  $Q, D_{\overline{\Omega}}$  and  $p_i$  for each  $i \in [n]$  that defines  $\mathcal{P}_{\overline{\Omega}}$ ).

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

## 528 **5 Related Work**

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

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

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

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

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

## 6 Conclusions and Future Work

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

## References

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



- 583 7 R. Iris Bahar, Erica A. Frohm, Charles M. Gaona, Gary D. Hachtel, Enrico Macii, Abelardo  
584 Pardo, and Fabio Somenzi. Algebraic decision diagrams and their applications. In *IEEE CAD*,  
585 1993.
- 586 8 George Beskales, Ihab F. Ilyas, and Lukasz Golab. Sampling the repairs of functional  
587 dependency violations under hard constraints. *Proc. VLDB Endow.*, 3(1):197–207, 2010.
- 588 9 Peter Bürgisser, Michael Clausen, and Mohammad Amin Shokrollahi. *Algebraic complexity*  
589 *theory*, volume 315. Springer, 1997.
- 590 10 Hubie Chen and Martin Grohe. Constraint satisfaction with succinctly specified relations. *J.*  
591 *Comput. Syst. Sci.*, 76(8):847–860, 2010.
- 592 11 Radu Curticapean. Counting matchings of size  $k$  is  $w[1]$ -hard. In *ICALP*, volume 7965, pages  
593 352–363, 2013.
- 594 12 N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. *VLDB*, 16(4):544,  
595 2007.
- 596 13 Nilesh Dalvi and Dan Suciu. The dichotomy of conjunctive queries on probabilistic structures.  
597 In *PODS*, pages 293–302, 2007.
- 598 14 Nilesh Dalvi and Dan Suciu. The dichotomy of probabilistic inference for unions of conjunctive  
599 queries. *JACM*, 59(6):30, 2012.
- 600 15 Maarten Van den Heuvel, Peter Ivanov, Wolfgang Gatterbauer, Floris Geerts, and Martin  
601 Theobald. Anytime approximation in probabilistic databases via scaled dissociations. In  
602 *SIGMOD*, pages 1295–1312, 2019.
- 603 16 Daniel Deutch, Tova Milo, Sudeepa Roy, and Val Tannen. Circuits for datalog provenance. In  
604 *ICDT*, pages 201–212, 2014.
- 605 17 Su Feng, Boris Glavic, Aaron Huber, and Oliver Kennedy. Efficient uncertainty tracking for  
606 complex queries with attribute-level bounds. In *SIGMOD*, 2021.
- 607 18 Su Feng, Aaron Huber, Boris Glavic, and Oliver Kennedy. Uncertainty annotated databases -  
608 a lightweight approach for approximating certain answers. In *SIGMOD*, 2019.
- 609 19 Robert Fink, Larisa Han, and Dan Olteanu. Aggregation in probabilistic databases via  
610 knowledge compilation. *PVLDB*, 5(5):490–501, 2012.
- 611 20 Robert Fink, Jiewen Huang, and Dan Olteanu. Anytime approximation in probabilistic  
612 databases. *VLDBJ*, 22(6):823–848, 2013.
- 613 21 Robert Fink and Dan Olteanu. Dichotomies for queries with negation in probabilistic databases.  
614 *TODS*, 41(1):4:1–4:47, 2016.
- 615 22 Jörg Flum and Martin Grohe. *Parameterized Complexity Theory*. Texts in Theoretical  
616 Computer Science. An EATCS Series. Springer, 2006. doi:10.1007/3-540-29953-X.
- 617 23 Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database systems - the*  
618 *complete book (2. ed.)*. Pearson Education, 2009.
- 619 24 George Grätzer. *Universal algebra*. Springer Science & Business Media, 2008.
- 620 25 Todd J. Green, Gregory Karvounarakis, and Val Tannen. Provenance semirings. In *PODS*,  
621 pages 31–40, 2007.
- 622 26 Todd J Green and Val Tannen. Models for incomplete and probabilistic information. In *EDBT*,  
623 pages 278–296. 2006.
- 624 27 T. Imielinski and W. Lipski. Incomplete information in relational databases. 1989.
- 625 28 Tomasz Imieliński and Witold Lipski Jr. Incomplete information in relational databases.  
626 *JACM*, 31(4):761–791, 1984.
- 627 29 Abhay Kumar Jha and Dan Suciu. Probabilistic databases with markovviews. *PVLDB*,  
628 5(11):1160–1171, 2012.
- 629 30 Richard M. Karp, Michael Luby, and Neal Madras. Monte-carlo approximation algorithms for  
630 enumeration problems. *J. Algorithms*, 10(3):429–448, 1989.
- 631 31 Batya Kenig, Avigdor Gal, and Ofer Strichman. A new class of lineage expressions over  
632 probabilistic databases computable in  $p$ -time. In *SUM*, volume 8078, pages 219–232, 2013.
- 633 32 Oliver Kennedy and Christoph Koch. Pip: A database system for great and small expectations.  
634 In *ICDE*, 2010.

- 635 33 Mahmoud Abo Khamis, Hung Q. Ngo, and Atri Rudra. Faq: Questions asked frequently. In  
636 *PODS*, pages 13–28, 2016.
- 637 34 Tsvi Kopelowitz and Virginia Vassilevska Williams. Towards optimal set-disjointness and  
638 set-intersection data structures. In *ICALP*, volume 168, pages 74:1–74:16, 2020.
- 639 35 Poonam Kumari, Said Achmiz, and Oliver Kennedy. Communicating data quality in on-demand  
640 curation. In *QDB*, 2016.
- 641 36 Hung Q. Ngo. Worst-case optimal join algorithms: Techniques, results, and open problems.  
642 In *PODS*, 2018.
- 643 37 Hung Q. Ngo, Christopher Ré, and Atri Rudra. Skew strikes back: new developments in the  
644 theory of join algorithms. *SIGMOD Rec.*, 42(4):5–16, 2013.
- 645 38 Dan Olteanu, Jiewen Huang, and Christoph Koch. Approximate confidence computation in  
646 probabilistic databases. In *ICDE*, pages 145–156, 2010.
- 647 39 Dan Olteanu and Maximilian Schleich. Factorized databases. *SIGMOD Rec.*, 45(2):5–16, 2016.
- 648 40 Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. Holoclean: Holistic data  
649 repairs with probabilistic inference. *Proc. VLDB Endow.*, 10(11):1190–1201, 2017.
- 650 41 Sudeepa Roy, Vittorio Perduca, and Val Tannen. Faster query answering in probabilistic  
651 databases using read-once functions. In *ICDT*, 2011.
- 652 42 C. Ré and D. Suciu. Materialized views in probabilistic databases: for information exchange  
653 and query optimization. In *VLDB*, pages 51–62, 2007.
- 654 43 Christopher De Sa, Alexander Ratner, Christopher Ré, Jaeho Shin, Feiran Wang, Sen Wu, and  
655 Ce Zhang. Incremental knowledge base construction using deepdive. *VLDB J.*, 26(1):81–105,  
656 2017.
- 657 44 Prithviraj Sen, Amol Deshpande, and Lise Getoor. Read-once functions and query evaluation  
658 in probabilistic databases. *PVLDB*, 3(1):1068–1079, 2010.
- 659 45 Pierre Senellart. Provenance and probabilities in relational databases. *SIGMOD Record*,  
660 46(4):5–15, 2018.
- 661 46 Leslie G. Valiant. The complexity of enumeration and reliability problems. *SIAM J. Comput.*,  
662 8(3):410–421, 1979.
- 663 47 Guy Van den Broeck and Dan Suciu. Query processing on probabilistic data: A survey. 2017.
- 664 48 Virginia Vassilevska Williams. Some open problems in fine-grained complexity. *SIGACT News*,  
665 49(4):29–35, 2018. doi:10.1145/3300150.3300158.
- 666 49 Ying Yang, Niccolò Meneghetti, Ronny Fehling, Zhen Hua Liu, Dieter Gawlick, and Oliver  
667 Kennedy. Lenses: An on-demand approach to etl. *PVLDB*, 8(12):1578–1589, 2015.

## 668 **7 Acknowledgements**

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

## 671 **A** Generalizing Beyond Set Inputs

### 672 **A.1** TIDBs

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

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

$$683 \quad \Phi[R, D_{\bar{\Omega}}, t] = \begin{cases} m_t X_t & \text{if } D_{\bar{\Omega}}.R(t) = m_t \\ 0 & \text{otherwise.} \end{cases}$$

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

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

### 700 **A.2** BIDBs

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

## 705 **B** Missing details from Section 2

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

707 We can use  $\mathcal{K}$ -relations to model bags. A  $\mathcal{K}$ -relation [25] is a relation whose tuples  
 708 are annotated with elements from a commutative semiring  $\mathcal{K} = \{K, \oplus_{\mathcal{K}}, \otimes_{\mathcal{K}}, 0_{\mathcal{K}}, 1_{\mathcal{K}}\}$ . A  
 709 commutative semiring is a structure with a domain  $K$  and associative and commutative  
 710 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  $\oplus_{\mathcal{K}}$ ,

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

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

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

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

Boris says:  
cite

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

$$744 \quad Mod(Q(M)) = Q(Mod(M))$$

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

$$747 \quad Mod(M) = \mathcal{D}$$

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

757 ▶ **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  
 758 equivalent  $\mathbb{N}$ -PDB  $\mathcal{D}_{\mathbb{N}} = \text{Mod}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = (\bar{\Omega}, \mathcal{P}_{\bar{\Omega}}')$  as:

$$759 \quad \bar{\Omega} = \{\psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) \mid \mathbf{w} \in \{0, 1\}^n\}$$

$$760 \quad \forall D \in \bar{\Omega} : Pr(D) = \sum_{\mathbf{w} \in \{0, 1\}^n : \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w})$$

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

$$766 \quad D_{[0,1]}(t_1) = 1 \quad \text{and} \quad D_{[1,1]}(t_1) = 2$$

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

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

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

$$775 \quad D_{\mathbb{N}[\mathbf{X}]}(t) = \sum_{i=1}^{|\bar{\Omega}|} D_i(t) \cdot X_i$$

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

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

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

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

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

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

$$816 \quad \mathbb{E}_{\mathbf{w} \sim \mathcal{P}_{\overline{\Omega}}(\mathbf{p})} [\Phi(\mathbf{W})] = \sum_{\substack{\mathbf{w} \in \{0,1\}^n \\ s.t. w_j, w_{j'} = 1 \rightarrow \exists b_i \supseteq \{t_{i,j}, t_{i,j'}\}}} \Phi(\mathbf{w}) \prod_{\substack{j \in [n] \\ s.t. w_j = 1}} p_j \prod_{\substack{j \in [n] \\ s.t. w_j = 0}} (1 - p_j) \quad (5)$$

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

**Boris says:**

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

### 821 B.3 Proof of Proposition 2.5

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

$$825 \quad \mathbb{E}_{\mathbf{w} \sim \mathcal{P}'_{\overline{\Omega}}} [\Phi(\mathbf{W})] = \sum_{\mathbf{w} \in \{0,1\}^n} Pr(\mathbf{w}) \cdot Q(D_{\mathbb{N}[\mathbf{X}]})(t)(\mathbf{w})$$

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  $\overline{\Omega}$ , so

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

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

$$831 \quad = \sum_{D \in \overline{\Omega}} \sum_{\mathbf{w} \in \{0,1\}^n : \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w}) \cdot Q(D)(t)$$

833 By distributivity of  $+$  over  $\times$

$$834 \quad = \sum_{D \in \overline{\Omega}} Q(D)(t) \sum_{\mathbf{w} \in \{0,1\}^n : \psi_{\mathbf{w}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = D} Pr(\mathbf{w})$$

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

$$837 \quad = \sum_{D \in \bar{\Omega}} Q(D)(t) \cdot Pr(D) = \mathbb{E}_{\mathbf{D} \sim \mathcal{P}_{\bar{\Omega}}} [Q(D)(t)]$$

838

## 840 B.4 Proposition B.4

841 Note the following fact:

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

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

## 849 B.5 Proof for Lemma ??

850 **Proof.** Let  $\Phi$  be a polynomial of  $n$  variables with highest degree =  $B$ , defined as follows:

$$851 \quad \Phi(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \{0, \dots, B\}^n} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n X_i^{d_i}.$$

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

$$854 \quad \mathbb{E}_{\mathbf{W}} [\Phi(\mathbf{W})] = \mathbb{E}_{\mathbf{W}} \left[ \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n W_i^{d_i} + \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \neg \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n W_i^{d_i} \right] \quad (6)$$

$$855 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[ \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n W_i^{d_i} \right] + \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \neg \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[ \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n W_i^{d_i} \right] \quad (7)$$

$$856 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[ \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n W_i^{d_i} \right] \quad (8)$$

$$857 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n \mathbb{E}_{\mathbf{W}} [W_i^{d_i}] \quad (9)$$

$$858 \quad = \sum_{\substack{\mathbf{d} \in \{0, \dots, B\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. \mathbf{d}_i \geq 1}}^n \mathbb{E}_{\mathbf{W}} [W_i] \quad (10)$$

<sup>8</sup> This BIDB notation is used and discussed in sec. 2.1.1

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

$$\begin{aligned}
&= \tilde{\Phi}(p_1, \dots, p_n). \quad (12)
\end{aligned}$$

Eq. (6) is the result of substituting in the definition of  $\Phi$  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 ??.

## B.6 Proof For Corollary 2.4

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

## C Missing details from Section 3

### C.1 Lemma C.1

► **Lemma C.1.** *Assuming that each  $v \in V$  has degree  $\geq 1$ ,<sup>9</sup> the PDB relations encoding the edges for  $\Phi_G^k$  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.

### 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  $Q^1$ ).

$$T_{det}^*(Q^1, D_{\bar{\Omega}}) = |D_{\bar{\Omega}}.V| + |D_{\bar{\Omega}}.E| + |D_{\bar{\Omega}}.V| + T_{join}(D_{\bar{\Omega}}.V, D_{\bar{\Omega}}.E, D_{\bar{\Omega}}.V).$$

We argue that  $T_{join}(D_{\bar{\Omega}}.V, D_{\bar{\Omega}}.E, D_{\bar{\Omega}}.V)$  is at most  $O(m)$  by noting that there exists an algorithm that computes  $D_{\bar{\Omega}}.V \bowtie D_{\bar{\Omega}}.E \bowtie D_{\bar{\Omega}}.V$  in the same runtime<sup>10</sup>. Then by the assumption of Lemma C.1 (each  $v \in V$  has degree  $\geq 1$ ), the sum of the first three terms is  $O(m)$ . We then obtain that  $T_{det}^*(Q^1, D_{\bar{\Omega}}) = O(m) + O(m) = O(m)$ . For  $Q^k = Q_1^1 \times \dots \times Q_k^1$ ,

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

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



891 we have the recurrence  $T_{det}^*(Q^k, D_{\overline{\Omega}}) = T_{det}^*(Q_1^1, D_{\overline{\Omega}}) + \dots + T_{det}^*(Q_k^1, D_{\overline{\Omega}}) + T_{join}(Q_1^1, \dots, Q_k^1)$ .  
 892 Since  $Q^1$  outputs a count, computing the join  $Q_1^1 \bowtie \dots \bowtie Q_k^1$  is just multiplying  $k$  numbers,  
 893 which takes  $O(k)$  time. Thus, we have

$$894 \quad T_{det}^*(Q^k, D_{\overline{\Omega}}) \leq k \cdot O(m) + O(k) \leq O(km),$$

895 as desired. ◀

### 896 C.3 Lemma C.2

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

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

### 901 C.4 Proof of Lemma C.2

902 **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  
 903 follows that  $\Phi_G^k(\mathbf{X})$  has degree  $2k$ . By definition,  $\tilde{\Phi}_G^k(\mathbf{X})$  sets every exponent  $e > 1$  to  $e = 1$ ,  
 904 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  
 905  $\tilde{\Phi}_G^k(p, \dots, p)$  is a univariate polynomial of degree at most  $\text{DEG}(\tilde{\Phi}_G^k) \leq 2k$ . Thus, we can write

$$906 \quad \tilde{\Phi}_G^k(p, \dots, p) = \sum_{i=0}^{2k} c_i p^i$$

907 We note that  $c_i$  is *exactly* the number of monomials in the SMB expansion of  $\Phi_G^k(\mathbf{X})$  composed  
 908 of  $i$  distinct variables.<sup>11</sup>

909 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  
 910 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  $(p_i^0 \dots p_i^{2k})$ ,  
 911  $\mathbf{c}$  is the coefficient vector  $(c_0, \dots, c_{2k})$ , and  $\mathbf{b}$  is the vector such that  $\mathbf{b}[i] = \tilde{\Phi}_G^k(p_i, \dots, p_i)$ .  
 912 In other words, matrix  $\mathbf{M}$  is the Vandermonde matrix, from which it follows that we have  
 913 a matrix with full rank (the  $p_i$ 's are distinct), and we can solve the linear system in  $O(k^3)$   
 914 time (e.g., using Gaussian Elimination) to determine  $\mathbf{c}$  exactly. Thus, after  $O(k^3)$  work, we  
 915 know  $\mathbf{c}$  and in particular,  $c_{2k}$  exactly.

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

$$919 \quad \Phi_G^k(\mathbf{X}) = \sum_{(i_1, j_1), \dots, (i_k, j_k) \in E} X_{i_1} X_{j_1} \dots X_{i_k} X_{j_k},$$

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

<sup>11</sup> 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.

924 the only surviving monomials  $\prod_{\ell=1}^k X_{i_\ell} X_{j_\ell}$  of degree exactly  $2k$  in  $\tilde{\Phi}_G^k(\mathbf{X})$  must have that  
 925 all of  $i_1, j_1, \dots, i_k, j_k$  are distinct in  $\Phi_G^k(\mathbf{X})$ . By the last two statements, only monomials  
 926 composed of  $2k$  distinct variables in  $\Phi_G^k(\mathbf{X})$  (and hence of degree  $2k$  in  $\tilde{\Phi}_G^k(\mathbf{X})$ ) correspond  
 927 to a  $k$ -matching in  $G$ .

928 As noted above, each of the  $k!$  permutations of an arbitrary monomial maps to the same  
 929 distinct  $k$ -matching in  $G$ , and this implies a  $k!$  to 1 mapping between degree  $2k$  monomials  
 930 in  $\tilde{\Phi}_G^k(\mathbf{X})$  and  $k$ -matchings in  $G$ . It then follows that  $c_{2k} = k! \cdot \#(G, \mathfrak{I} \cdots \mathfrak{I}^k)$ . Thus, simply  
 931 dividing  $c_{2k}$  by  $k!$  gives us  $\#(G, \mathfrak{I} \cdots \mathfrak{I}^k)$ , as needed. ◀

### 932 C.5 Proof of Theorem 3.6

933 **Proof.** For the sake of contradiction, assume we can solve our problem in  $o(T_{\text{match}}(k, G))$   
 934 time. Given a graph  $G$  by Lemma C.1 we can compute the PDB encoding in  $O(m)$  time. Then  
 935 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  
 936  $O(k) \cdot o(T_{\text{match}}(k, G))$  time. Lemma C.2 then computes the number of  $k$ -matchings in  $G$  in  
 937  $O(k^3)$  time. Adding the runtime of all of these steps, we have an algorithm for computing  
 938 the number of  $k$ -matchings that runs in time

$$939 \quad O(m) + O(k) \cdot o(T_{\text{match}}(k, G)) + O(k^3) \quad (13)$$

$$940 \quad \leq o(T_{\text{match}}(k, G)). \quad (14)$$

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

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

947 We need all the possible edge patterns in an arbitrary  $G$  with at most three distinct edges.  
 948 We have already seen  $\mathfrak{I}$ ,  $\mathfrak{II}$  and  $\mathfrak{III}$ , so we define the remaining patterns:

- 949 ■ Single Edge ( $\mathfrak{I}$ )
- 950 ■ 2-path ( $\mathfrak{A}$ )
- 951 ■ 2-matching ( $\mathfrak{II}$ )
- 952 ■ 3-star ( $\mathfrak{B}$ )—this is the graph that results when all three edges share exactly one common  
 953 endpoint. The remaining endpoint for each edge is disconnected from any endpoint of  
 954 the remaining two edges.
- 955 ■ Disjoint Two-Path ( $\mathfrak{C}$ )—this subgraph consists of a two-path and a remaining disjoint  
 956 edge.

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

$$959 \quad \#(G, \mathfrak{I}) = m, \quad (15)$$

$$960 \quad \#(G, \mathfrak{A}) = \sum_{i \in V} \binom{d_i}{2} \quad (16)$$

$$961 \quad \#(G, \mathfrak{II}) = \sum_{(i,j) \in E} \frac{m - d_i - d_j + 1}{2} \quad (17)$$

$$962 \quad \#(G, \mathfrak{B}) = \sum_{i \in V} \binom{d_i}{3} \quad (18)$$

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

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

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

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

**Proof of Eq. (17).** For edge  $(i, j)$  connecting arbitrary vertices  $i$  and  $j$ , finding all other edges in  $G$  disjoint to  $(i, j)$  is equivalent to finding all edges that are not connected to either vertex  $i$  or  $j$ . The number of such edges is  $m - d_i - d_j + 1$ , where we add 1 since edge  $(i, j)$  is removed twice when subtracting both  $d_i$  and  $d_j$ . Since the summation is iterating over all edges such that a pair  $((i, j), (k, \ell))$  will also be counted as  $((k, \ell), (i, j))$ , division by 2 then eliminates this double counting. Note that  $m$  and  $d_i$  for all  $i \in V$  can be computed in one pass over the set of edges by simply maintaining counts for each quantity. Finally, the summation is also one traversal through the set of edges where each operation is either a lookup ( $O(1)$  time) or an addition operation (also  $O(1)$ ) time. ◀

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

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

1003 **C.8 Tools to prove Theorem 3.7**

1004 Note that  $\tilde{\Phi}_G^3(p, \dots, p)$  as a polynomial in  $p$  has degree at most six. Next, we figure out the  
 1005 exact coefficients since this would be useful in our arguments:

1006 ► **Lemma C.3.** *For any  $p$ , we have:*

$$1007 \quad \tilde{\Phi}_G^3(p, \dots, p) = \#(G, \mathfrak{I})p^2 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{II})p^4 + 6\#(G, \mathfrak{B})p^3 \\ 1008 \quad + 6\#(G, \mathfrak{AA})p^4 + 6\#(G, \mathfrak{I\I})p^4 + 6\#(G, \mathfrak{I\A})p^5 + 6\#(G, \mathfrak{III})p^6. \quad (21)$$

1010 **C.8.1 Proof for Lemma C.3**

1011 **Proof.** By definition we have that

$$1012 \quad \Phi_G^3(\mathbf{X}) = \sum_{(i_1, j_1), (i_2, j_2), (i_3, j_3) \in E} \prod_{\ell=1}^3 X_{i_\ell} X_{j_\ell}.$$

1013 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  
 1014 coefficient of  $p^\nu$  in  $\tilde{\Phi}_G^3(\mathbf{X})$ , where  $\nu$  is the number of distinct variables in the monomial. Let  
 1015  $e_1 = (i_1, j_1)$ ,  $e_2 = (i_2, j_2)$ , and  $e_3 = (i_3, j_3)$ . We compute  $\tilde{\Phi}_G^3(\mathbf{X})$  by considering each of the  
 1016 three forms that the triple  $(e_1, e_2, e_3)$  can take.

1017 **CASE 1:**  $e_1 = e_2 = e_3$  (all edges are the same). When we have that  $e_1 = e_2 = e_3$ , then  
 1018 the monomial corresponds to  $\#(G, \mathfrak{I})$ . There are exactly  $m$  such triples, each with a  $p^2$   
 1019 factor in  $\tilde{\Phi}_G^3(p, \dots, p)$ .

1020 **CASE 2:** This case occurs when there are two distinct edges of the three, call them  $e$  and  
 1021  $e'$ . When there are two distinct edges, there is then the occurrence when 2 variables in the  
 1022 triple  $(e_1, e_2, e_3)$  are bound to  $e$ . There are three combinations for this occurrence in  $\Phi_G^3(\mathbf{X})$ .  
 1023 Analogously, there are three such occurrences in  $\Phi_G^3(\mathbf{X})$  when there is only one occurrence of  
 1024  $e$ , i.e. 2 of the variables in  $(e_1, e_2, e_3)$  are  $e'$ . This implies that all  $3 + 3 = 6$  combinations of  
 1025 two distinct edges  $e$  and  $e'$  contribute to the same monomial in  $\tilde{\Phi}_G^3$ . Since  $e \neq e'$ , this case  
 1026 produces the following edge patterns:  $\mathfrak{A}, \mathfrak{II}$ , which contribute  $6p^3$  and  $6p^4$  respectively to  
 1027  $\tilde{\Phi}_G^3(p, \dots, p)$ .

1028 **CASE 3:** All  $e_1, e_2$  and  $e_3$  are distinct. For this case, we have  $3! = 6$  permutations  
 1029 of  $(e_1, e_2, e_3)$ , each of which contribute to the same monomial in the SMB representation  
 1030 of  $\Phi_G^3(\mathbf{X})$ . This case consists of the following edge patterns:  $\mathfrak{B}, \mathfrak{AA}, \mathfrak{I\I}, \mathfrak{I\A}, \mathfrak{III}$ , which  
 1031 contribute  $6p^3, 6p^4, 6p^4, 6p^5$  and  $6p^6$  respectively to  $\tilde{\Phi}_G^3(p, \dots, p)$ . ◀

1032 Since  $p$  is fixed, Lemma C.3 gives us one linear equation in  $\#(G, \mathfrak{A})$  and  $\#(G, \mathfrak{III})$  (we  
 1033 can handle the other counts due to equations (15)-(20)). However, we need to generate  
 1034 one more independent linear equation in these two variables. Towards this end we generate  
 1035 another graph related to  $G$ :

1036 ► **Definition C.4.** *For  $\ell \geq 1$ , let graph  $G^{(\ell)}$  be a graph generated from an arbitrary graph  
 1037  $G$ , by replacing every edge  $e$  of  $G$  with an  $\ell$ -path, such that all inner vertexes of an  $\ell$ -path  
 1038 replacement edge are disjoint from all other vertexes.<sup>12</sup>*

1039 We will prove Theorem 3.7 by the following reduction:

1040 ► **Theorem C.5.** *Fix  $p \in (0, 1)$ . Let  $G$  be a graph on  $m$  edges. If we can compute  $\tilde{\Phi}_G^3(p, \dots, p)$   
 1041 exactly in  $T(m)$  time, then we can exactly compute  $\#(G, \mathfrak{A})$  in  $O(T(m) + m)$  time.*

<sup>12</sup>Note that  $G \equiv G^{(1)}$ .

1042 For clarity, we repeat the notion of  $\#(G, H)$  to mean the count of subgraphs in  $G$  isomorphic  
 1043 to  $H$ . The following lemmas relate these counts in  $G^{(2)}$  to  $G^{(1)}$  ( $G$ ). The lemmas are used  
 1044 to prove Lemma C.8.

1045 ► **Lemma C.6.** *The 3-matchings in graph  $G^{(2)}$  satisfy the identity:*

$$1046 \quad \#(G^{(2)}, \mathfrak{III}) = 8 \cdot \#(G^{(1)}, \mathfrak{III}) + 6 \cdot \#(G^{(1)}, \mathfrak{I} \mathfrak{A}) \\ 1047 \quad \quad \quad + 4 \cdot \#(G^{(1)}, \mathfrak{AA}) + 4 \cdot \#(G^{(1)}, \mathfrak{II}) + 2 \cdot \#(G^{(1)}, \mathfrak{A}). \\ 1048$$

1049 ► **Lemma C.7.** *For  $\ell > 1$  and any graph  $G^{(\ell)}$ ,  $\#(G^{(\ell)}, \mathfrak{A}) = 0$ .*

1050 Finally, the following result immediately implies Theorem C.5:

1051 ► **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  $O(m)$   
 1052 time a vector  $\mathbf{b} \in \mathbb{R}^3$  such that*

$$1053 \quad \begin{pmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{pmatrix} \cdot \begin{pmatrix} \#(G, \mathfrak{A}) \\ \#(G, \mathfrak{III}) \end{pmatrix} = \mathbf{b},$$

1054 allowing us to compute  $\#(G, \mathfrak{A})$  and  $\#(G, \mathfrak{III})$  in  $O(1)$  time.

## 1055 C.9 Proofs for Lemma C.6, Lemma C.7, and Lemma C.8

1056 Before proceeding, let us introduce a few more helpful definitions.

1057 ► **Definition C.9** ( $E^{(\ell)}$ ). *For  $\ell > 1$ , we use  $E^{(\ell)}$  to denote the set of edges in  $G^{(\ell)}$ . For  
 1058 any graph  $G^{(\ell)}$ , its edges are denoted by the a pair  $(e, b)$ , such that  $b \in \{0, \dots, \ell - 1\}$  where  
 1059  $(e, 0), \dots, (e, \ell - 1)$  is the  $\ell$ -path that replaces the edge  $e$  for  $e \in E^{(1)}$ .*

1060 ► **Definition C.10** ( $E_S^{(\ell)}$ ). *Given an arbitrary subgraph  $S^{(1)}$  of  $G^{(1)}$ , let  $E_S^{(1)}$  denote the set  
 1061 of edges in  $S^{(1)}$ . Define then  $E_S^{(\ell)}$  for  $\ell > 1$  as the set of edges in the generated subgraph  $S^{(\ell)}$   
 1062 (i.e. when we apply Definition C.4 to  $S$  to generate  $S^{(\ell)}$ ).*

1063 For example, consider  $S^{(1)}$  with edges  $E_S^{(1)} = \{e_1\}$ . Then the edge set of  $S^{(2)}$  is defined  
 1064 as  $E_S^{(2)} = \{(e_1, 0), (e_1, 1)\}$ .

1065 ► **Definition C.11** ( $\binom{E}{t}$  and  $\binom{E}{\leq t}$ ). *Let  $\binom{E}{t}$  denote the set of subsets in  $E$  with exactly  $t$   
 1066 edges. In a similar manner,  $\binom{E}{\leq t}$  is used to mean the subsets of  $E$  with  $t$  or fewer edges.*

1067 The following function  $f_\ell$  is a mapping from every 3-edge shape in  $G^{(\ell)}$  to its ‘projection’  
 1068 in  $G^{(1)}$ .

1069 ► **Definition C.12.** *Let  $f_\ell : \binom{E^{(\ell)}}{3} \rightarrow \binom{E^{(1)}}{\leq 3}$  be defined as follows. For any element  $s \in \binom{E^{(\ell)}}{3}$   
 1070 such that  $s = \{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}$ , define:*

$$1071 \quad f_\ell(\{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}) = \{e_1, e_2, e_3\}.$$

1072 ► **Definition C.13** ( $f_\ell^{-1}$ ). *For an arbitrary subgraph  $S^{(1)}$  of  $G^{(1)}$  with at most  $m \leq 3$  edges,  
 1073 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  
 1074  $s \in \binom{E_S^{(\ell)}}{3}$  such that  $f_\ell(s) = E_S^{(1)}$ .*

1075 Note, importantly, that when we discuss  $f_\ell^{-1}$ , that each *edge* present in  $E_S^{(1)}$  must have an  
 1076 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  
 1077 case that each  $s \in f_\ell^{-1}(E_S^{(1)})$  consists of the following set of edges:  $\{(e_i, b), (e_j, b'), (e_m, b'')\}$ ,  
 1078 where  $i, j$  and  $m$  are distinct.

1079 We are now ready to prove the structural lemmas. To prove the structural lemmas, we  
 1080 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  
 1081 many  $\mathfrak{B}$  and  $\mathfrak{A}$  subgraphs appear in  $f_\ell^{-1}(E_S^{(1)})$ .

### 1082 C.9.1 Proof of Lemma C.6

1083 **Proof.** For each subset  $E_S^{(1)} \in \binom{E_1}{\leq 3}$ , we count the number of 3-matchings in the 3-edge  
 1084 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  $E_S^{(1)}$   
 1085 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  $s \in$   
 1086  $\{(e_1, 0), (e_1, 1), (e_2, 0), (e_2, 1), (e_3, 0), (e_3, 1)\}$  such that  $f_\ell(s) = \{e_1, e_2, e_3\}$ . The size of  
 1087 the output is denoted  $|f_2^{-1}(E_S^{(1)})|$ . For the case where each set of edges of the form  
 1088  $\{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}$  for  $b_i \in [2], i \in [3]$  is present, we have  $|f_2^{-1}(E_S^{(1)})| = 8$ . We count  
 1089 the number of 3-matchings from the set  $f_2^{-1}(E_S^{(1)})$ .

1090 We do a case analysis based on the subgraph  $S^{(1)}$  induced by  $E_S^{(1)}$ .

#### 1091 ■ 3-matching ( $\mathfrak{B}$ )

1092 When  $S^{(1)}$  is isomorphic to  $\mathfrak{B}$ , it is the case that edges in  $E_S^{(2)}$  are *not* disjoint only for the  
 1093 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  
 1094 matching and  $|f_2^{-1}(E_S^{(1)})| = 8$  possible 3-matchings.

#### 1095 ■ Disjoint Two-Path ( $\mathfrak{A}$ )

1096 For  $S^{(1)}$  isomorphic to  $\mathfrak{A}$  edges  $e_2, e_3$  form a 2-path with  $e_1$  being disjoint. This means  
 1097 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  
 1098 disjoint 2-path. We can pick either  $(e_1, 0)$  or  $(e_1, 1)$  for the first edge in the 3-matching, while  
 1099 it is necessary to have a 2-matching from path  $(e_2, 0), \dots, (e_3, 1)$ . Note that the 4-path allows  
 1100 for three possible 2-matchings, specifically,

$$1101 \quad \{(e_2, 0), (e_3, 0)\}, \{(e_2, 0), (e_3, 1)\}, \{(e_2, 1), (e_3, 1)\}.$$

1102 Since these two selections can be made independently,  $|f_2^{-1}(E_S^{(1)})| = 2 \cdot 3 = 6$  *distinct*  
 1103 3-matchings in  $f_2^{-1}(E_S^{(1)})$ .

#### 1104 ■ 3-star ( $\mathfrak{C}$ )

1105 When  $S^{(1)}$  is isomorphic to  $\mathfrak{C}$ , the inner edges  $(e_i, 1)$  of  $S^{(2)}$  are all connected, and the  
 1106 outer edges  $(e_i, 0)$  are all disjoint. Note that for a valid 3-matching it must be the case that  
 1107 at most one inner edge can be part of the set of disjoint edges. For the case of when exactly  
 1108 one inner edge is chosen, there exist 3 possibilities, based on which inner edge is chosen.  
 1109 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  
 1110  $j' \neq i, j' \neq j$ . The remaining possible 3-matching occurs when all 3 outer edges are chosen,  
 1111 and  $|f_2^{-1}(E_S^{(1)})| = 4$ .

#### 1112 ■ 3-path ( $\mathfrak{D}$ )

1113 When  $S^{(1)}$  is isomorphic to  $\mathfrak{I}\mathfrak{I}$  it is the case that all edges beginning with  $e_1$  and ending with  $e_3$   
 1114 are successively connected. This means that the edges of  $E_S^{(2)}$  form a 6-path. For a 3-matching  
 1115 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  
 1116  $j = i + 1$ . There are four such possibilities:  $\{(e_1, 0), (e_2, 0), (e_3, 0)\}$ ,  $\{(e_1, 0), (e_2, 0), (e_3, 1)\}$ ,  
 1117  $\{(e_1, 0), (e_2, 1), (e_3, 1)\}$ ,  $\{(e_1, 1), (e_2, 1), (e_3, 1)\}$  and  $|f_2^{-1}(E_S^{(1)})| = 4$ .

1118 ■ Triangle ( $\mathfrak{A}$ )

1119 For  $S^{(1)}$  isomorphic to  $\mathfrak{A}$ , note that it is the case that the edges in  $E_S^{(2)}$  are connected in a  
 1120 successive manner, but this time in a cycle, such that  $(e_1, 0)$  and  $(e_3, 1)$  are also connected.  
 1121 While this is similar to the discussion of the three path above, the first and last edges are  
 1122 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  
 1123 that  $|f_2^{-1}(E_S^{(1)})| = 2$ .

1124 Let us now consider when  $E_S^{(1)} \in \binom{E_1}{\leq 2}$ , i.e. fixed subgraphs among

1125 ■ 2-matching ( $\mathfrak{I}\mathfrak{I}$ ), 2-path ( $\mathfrak{A}$ ), 1 edge ( $\mathfrak{I}$ )

1126 When  $|E_S^{(1)}| = 2$ , we can only pick one from each of two pairs,  $\{(e_1, 0), (e_1, 1)\}$  and  
 1127  $\{(e_2, 0), (e_2, 1)\}$ . The third edge choice in  $E_S^{(2)}$  will break the disjoint property of a 3-  
 1128 matching. Thus, a 3-matching cannot exist in  $f_2^{-1}(E_S^{(1)})$ . A similar argument holds for  
 1129  $|E_S^{(1)}| = 1$ , where the output of  $f_2^{-1}$  is  $\{\emptyset\}$  since there are not enough edges in the input to  
 1130 produce any other output.

1131 Observe that all of the arguments above focused solely on the property of subgraph  $S^{(1)}$   
 1132 being isomorphic. In other words, all  $E_S^{(1)}$  of a given “shape” yield the same number of  
 1133 3-matchings in  $f_2^{-1}(E_S^{(1)})$ , and this is why we get the required identity using the above case  
 1134 analysis. ◀

### 1135 C.9.2 Proof of Lemma C.7

1136 **Proof.** The number of triangles in  $G^{(\ell)}$  for  $\ell \geq 2$  will always be 0 for the simple fact that all  
 1137 cycles in  $G^{(\ell)}$  will have at least six edges. ◀

### 1138 C.9.3 Proof of Lemma C.8

1139 **Proof.** The proof consists of two parts. First we need to show that a vector  $\mathbf{b}$  satisfying the  
 1140 linear system exists and further can be computed in  $O(m)$  time. Second we need to show  
 1141 that  $\#(G, \mathfrak{A}), \#(G, \mathfrak{I}\mathfrak{I}\mathfrak{I})$  can indeed be computed in time  $O(1)$ .

1142 The lemma claims that for  $\mathbf{M} = \begin{pmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{pmatrix}$ ,  $\mathbf{x} = \begin{pmatrix} \#(G, \mathfrak{A}) \\ \#(G, \mathfrak{I}\mathfrak{I}\mathfrak{I}) \end{pmatrix}$   
 1143 satisfies the linear system  $\mathbf{M} \cdot \mathbf{x} = \mathbf{b}$ .

1144 To prove the first step, we use Lemma C.3 to derive the following equality (dropping the  
 1145 superscript and referring to  $G^{(1)}$  as  $G$ ):

$$1146 \begin{aligned} & \#(G, \mathfrak{I})p^2 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{I}\mathfrak{I})p^4 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{A}\mathfrak{A})p^4 \\ 1147 & + 6\#(G, \mathfrak{I}\mathfrak{I})p^4 + 6\#(G, \mathfrak{I}\mathfrak{A})p^5 + 6\#(G, \mathfrak{I}\mathfrak{I}\mathfrak{I})p^6 = \tilde{\Phi}_G^3(p, \dots, p) \end{aligned} \quad (22)$$

$$1148 \begin{aligned} & \#(G, \mathfrak{A}) + \#(G, \mathfrak{I}\mathfrak{I})p + \#(G, \mathfrak{I}\mathfrak{A})p^2 + \#(G, \mathfrak{I}\mathfrak{I}\mathfrak{I})p^3 \\ 1149 & = \frac{\tilde{\Phi}_G^3(p, \dots, p)}{6p^3} - \frac{\#(G, \mathfrak{I})}{6p} - \#(G, \mathfrak{A}) - \#(G, \mathfrak{I}\mathfrak{I})p - \#(G, \mathfrak{A}\mathfrak{A})p \end{aligned} \quad (23)$$

## 23:32 Bag PDB Queries

$$\begin{aligned}
1150 \quad & \#(G, \mathfrak{A})(1-3p) - \#(G, \mathfrak{III})(3p^2 - p^3) = \\
1151 \quad & \frac{\tilde{\Phi}_G^3(p, \dots, p)}{6p^3} - \frac{\#(G, \mathfrak{I})}{6p} - \#(G, \mathfrak{A}) - \#(G, \mathfrak{II})p - \#(G, \mathfrak{AA})p \\
1152 \quad & - [\#(G, \mathfrak{II})p + 3\#(G, \mathfrak{A})p] - [\#(G, \mathfrak{I} \mathfrak{A})p^2 + 3\#(G, \mathfrak{III})p^2] \\
1153 \quad & \tag{24}
\end{aligned}$$

1154 Eq. (22) is the result of Lemma C.3. We obtain the remaining equations through standard  
1155 algebraic manipulations.

1156 Note that the LHS of Eq. (24) is obtained using eq. (19) and eq. (20) and is indeed the  
1157 product  $\mathbf{M}[1] \cdot \mathbf{x}[1]$ . Further note that this product is equal to the RHS of Eq. (24), where  
1158 every term is computable in  $O(m)$  time (by equations (15)-(20)). We set  $\mathbf{b}[1]$  to the RHS of  
1159 Eq. (24).

1160 We follow the same process in deriving an equality for  $G^{(2)}$ . Replacing occurrences of  
1161  $G$  with  $G^{(2)}$ , we obtain an equation (below) of the form of eq. (24) for  $G^{(2)}$ . Substituting  
1162 identities from lemma C.6 and Lemma C.7 we obtain

$$\begin{aligned}
1163 \quad & 0 - (8\#(G, \mathfrak{III}) + 6\#(G, \mathfrak{I} \mathfrak{A}) + 4\#(G, \mathfrak{AA}) + 4\#(G, \mathfrak{II}) + 2\#(G, \mathfrak{A}))(3p^2 - p^3) = \\
1164 \quad & \frac{\tilde{\Phi}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{I})}{6p} - \#(G^{(2)}, \mathfrak{A}) - \#(G^{(2)}, \mathfrak{II})p - \#(G^{(2)}, \mathfrak{AA})p \\
1165 \quad & - [\#(G^{(2)}, \mathfrak{II})p + 3\#(G^{(2)}, \mathfrak{A})p] - [\#(G^{(2)}, \mathfrak{I} \mathfrak{A})p^2 + 3\#(G^{(2)}, \mathfrak{III})p^2] - [\#(G^{(2)}, \mathfrak{II})p + 3\#(G^{(2)}, \mathfrak{A})p] \\
1166 \quad & \tag{25} \\
1166 \quad & (10\#(G, \mathfrak{A}) + 10G\mathfrak{III})(3p^2 - p^3) = \\
1167 \quad & \frac{\tilde{\Phi}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{I})}{6p} - \#(G^{(2)}, \mathfrak{A}) - \#(G^{(2)}, \mathfrak{II})p - \#(G^{(2)}, \mathfrak{AA})p \\
1168 \quad & - [\#(G^{(2)}, \mathfrak{II})p + 3\#(G^{(2)}, \mathfrak{A})p] - [\#(G^{(2)}, \mathfrak{I} \mathfrak{A})p^2 - 3\#(G^{(2)}, \mathfrak{III})p^2] \\
1169 \quad & + (4\#(G, \mathfrak{AA}) + [6\#(G, \mathfrak{I} \mathfrak{A}) + 18\#(G, \mathfrak{III})] + [4\#(G, \mathfrak{II}) + 12\#(G, \mathfrak{A})]) (3p^2 - p^3) \\
1170 \quad & \tag{26}
\end{aligned}$$

1171 The steps to obtaining eq. (26) are analogous to the derivation immediately preceding. As in  
1172 the previous derivation, note that the LHS of Eq. (26) is the same as  $\mathbf{M}[2] \cdot \mathbf{x}[2]$ . The RHS  
1173 of Eq. (26) has terms all computable (by equations (15)-(20)) in  $O(m)$  time. Setting  $\mathbf{b}[2]$  to  
1174 the RHS then completes the proof of step 1.

1175 Note that if  $\mathbf{M}$  has full rank then one can compute  $\#(G, \mathfrak{A})$  and  $\#(G, \mathfrak{III})$  in  $O(1)$   
1176 using Gaussian elimination.

1177 To show that  $\mathbf{M}$  indeed has full rank, we show in what follows that  $\text{Det}(\mathbf{M}) \neq 0$  for  
1178 every  $p \in (0, 1)$ .  $\text{Det}(\mathbf{M}) =$

$$\begin{aligned}
1179 \quad & \begin{vmatrix} 1-3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{vmatrix} = (1-3p) \cdot 10(3p^2 - p^3) + 10(3p^2 - p^3) \cdot (3p^2 - p^3) \\
1180 \quad & = 10(3p^2 - p^3) \cdot (1-3p + 3p^2 - p^3) = 10(3p^2 - p^3) \cdot (-p^3 + 3p^2 - 3p + 1) \\
1181 \quad & = 10p^2(3-p) \cdot (1-p)^3 \\
1182 \quad & \tag{27}
\end{aligned}$$

1183 From Eq. (27) it can easily be seen that the roots of  $\text{Det}(\mathbf{M})$  are 0, 1, and 3. Hence there  
1184 are no roots in  $(0, 1)$  and Lemma C.8 follows.  $\blacktriangleleft$



1185 **C.10 Proof of Theorem C.5**

1186 **Proof.** We can compute  $G^{(2)}$  from  $G^{(1)}$  in  $O(m)$  time. Additionally, if in time  $O(T(m))$ , we  
 1187 have  $\tilde{\Phi}_{G^{(\ell)}}^3(p, \dots, p)$  for  $\ell \in [2]$ , then the theorem follows by Lemma C.8. ◀

1188 In other words, if Theorem C.5 holds, then so must Theorem 3.7.

1189 **C.11 Proof of Theorem 3.7**

1190 **Proof.** For the sake of contradiction, assume that for any  $G$ , we can compute  $\tilde{\Phi}_G^3(p, \dots, p)$   
 1191 in  $o(m^{1+\epsilon_0})$  time. Let  $G$  be the input graph. It is easy to see that one can compute the  
 1192 expression tree for  $\Phi_G^3(\mathbf{X})$  in  $O(m)$  time. Then by Theorem C.5 we can compute  $\#(G, \mathfrak{A})$   
 1193 in further time  $o(m^{1+\epsilon_0}) + O(m)$ . Thus, the overall, reduction takes  $o(m^{1+\epsilon_0}) + O(m) =$   
 1194  $o(m^{1+\epsilon_0})$  time, which violates Conjecture 3.3. ◀

1195 **D Missing Details from Section 4**

1196 In the following definitions and examples, we use the following polynomial as an example:

$$1197 \quad \Phi(X, Y) = 2X^2 + 3XY - 2Y^2. \quad (28)$$

1198 ▶ **Definition D.1** (Pure Expansion). *The pure expansion of a polynomial  $\Phi$  is formed by*  
 1199 *computing all product of sums occurring in  $\Phi$ , without combining like monomials. The pure*  
 1200 *expansion of  $\Phi$  generalizes Definition 2.1 by allowing monomials  $m_i = m_j$  for  $i \neq j$ .*

1201 Note that similar in spirit to ??,  $\mathbf{E}(\mathbf{C})$  Definition 4.1 reduces all variable exponents  $e > 1$  to  
 1202  $e = 1$ . Further, it is true that  $\mathbf{E}(\mathbf{C})$  is the pure expansion of  $\mathbf{C}$ .

1203 ▶ **Example D.2** (Example of Pure Expansion). *Consider the factorized representation  $(X +$   
 1204  $2Y)(2X - Y)$  of the polynomial in Eq. (28). Its circuit  $\mathbf{C}$  is illustrated in Fig. 3. The*  
 1205 *pure expansion of the product is  $2X^2 - XY + 4XY - 2Y^2$ . As an additional example of*  
 1206 *Definition 4.1,  $\mathbf{E}(\mathbf{C}) = [(X, 2), (XY, -1), (XY, 4), (Y, -2)]$ .*

1207  $\mathbf{E}(\mathbf{C})$  effectively<sup>13</sup> encodes the *reduced* form of  $\text{POLY}(\mathbf{C})$ , decoupling each monomial into a  
 1208 set of variables  $\mathbf{v}$  and a real coefficient  $\mathbf{c}$ . However, unlike the constraint on the input  $\Phi$  to  
 1209 compute  $\tilde{\Phi}$ , the input circuit  $\mathbf{C}$  does not need to be in SMB/SOP form.

1210 ▶ **Example D.3** (Example for Definition 4.2). *Using the same factorization from Example D.2,*  
 1211  *$\text{POLY}(|\mathcal{C}|) = (X + 2Y)(2X + Y) = 2X^2 + XY + 4XY + 2Y^2 = 2X^2 + 5XY + 2Y^2$ . Note that*  
 1212 *this is not the same as the polynomial from Eq. (28). As an example of the slight abuse of*  
 1213 *notation we alluded to,  $\text{POLY}(|\mathcal{C}|(1, \dots, 1)) = 2(1)^2 + 5(1)(1) + 2(1)^2 = 9$ .*

1214 **Aaron says:** Verify whether we need pure expansion or not.

1215 ▶ **Definition D.4** (Subcircuit). *A subcircuit of a circuit  $\mathbf{C}$  is a circuit  $\mathbf{S}$  such that  $\mathbf{S}$  is a DAG*  
 1216 *subgraph of the DAG representing  $\mathbf{C}$ . The sink of  $\mathbf{S}$  has exactly one gate  $g$ .*

1217 The following results assume input circuit  $\mathbf{C}$  computed from an arbitrary  $\mathcal{RA}^+$  query  $Q$   
 1218 and arbitrary BIDB  $\mathcal{D}$ . We refer to  $\mathbf{C}$  as a BIDB circuit.

1219 **Aaron says:** Verify that the proof for Theorem D.5 doesn't rely on properties of  $\mathcal{RA}^+$  or BIDB.

<sup>13</sup>The minor difference here is that  $\mathbf{E}(\mathbf{C})$  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 (ref. 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  $\text{acc}$ 

```

---

1220 ► **Theorem D.5.** Let  $\mathcal{C}$  be an arbitrary BIDB circuit and define  $\Phi(\mathbf{X}) = \text{POLY}(\mathcal{C})$  and let  
 1221  $k = \text{DEG}(\mathcal{C})$ . Then an estimate  $\mathcal{E}$  of  $\tilde{\Phi}(p_1, \dots, p_n)$  can be computed in time

$$1222 \quad O\left(\left(\text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta} \cdot |\mathcal{C}|^2(1, \dots, 1) \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C})}{(\epsilon)^2 \cdot \tilde{\Phi}^2(p_1, \dots, p_n)}\right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right)$$

1223 such that

$$1224 \quad \Pr\left(\left|\mathcal{E} - \tilde{\Phi}(p_1, \dots, p_n)\right| > \epsilon \cdot \tilde{\Phi}(p_1, \dots, p_n)\right) \leq \delta. \quad (29)$$

1225 **Atri says: Aaron:**Just copied over from S4. The above text might need smoothening  
 1226 into the appendix.

1226 The slight abuse of notation seen in  $|\mathcal{C}|(1, \dots, 1)$  is explained after Definition 4.2 and  
 1227 an example is given in Example D.3. The only difference in the use of this notation in  
 1228 Theorem D.5 is that we include an additional exponent to square the quantity.

## 1229 D.1 Proof of Theorem D.5

1230 We prove Theorem D.5 constructively by presenting an algorithm APPROXIMATE $\tilde{\Phi}$  (Algorithm 1)  
 1231 which has the desired runtime and computes an approximation with the desired approximation  
 1232 guarantee. Algorithm APPROXIMATE $\tilde{\Phi}$  uses Algorithm ONEPASS to compute weights on the  
 1233 edges of a circuits. These weights are then used to sample a set of monomials of  $\Phi(\mathcal{C})$  from  
 1234 the circuit  $\mathcal{C}$  by traversing the circuit using the weights to ensure that monomials are sampled  
 1235 with an appropriate probability. The correctness of APPROXIMATE $\tilde{\Phi}$  relies on the correctness  
 1236 (and runtime behavior) of auxiliary algorithms ONEPASS and SAMPLEMONOMIAL that we  
 1237 state in the following lemmas (and prove later in this part of the appendix).

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

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

1238 ONEPASS guarantees two post-conditions: First, for each subcircuit  $\mathcal{S}$  of  $\mathcal{C}$ , we have that  
 1239  $\mathcal{S}.\text{partial}$  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  
 1240 likewise for  $\mathcal{S}.\text{Rweight}$ .

1241 To prove correctness of Algorithm 1, we only use the following fact that follows from the above  
 1242 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*

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

1243 where  $k = \text{DEG}(\mathcal{C})$ . The function returns every  $(v, \text{sign}(c))$  for  $(v, c) \in E(\mathcal{C})$  with probability  
 1244  $\frac{|c|}{|\mathcal{C}|(1, \dots, 1)}$ .

1245 With the above two lemmas, we are ready to argue the following result:

1246 ► **Theorem D.8.** *For any  $\mathcal{C}$  with*

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

1247  $\text{DEG}(\text{poly}(|\mathcal{C}|)) = k$ , algorithm 1 outputs an estimate  $\text{acc}$  of  $\tilde{\Phi}(p_1, \dots, p_n)$  such that

1248 
$$\Pr\left(\left|\text{acc} - \tilde{\Phi}(p_1, \dots, p_n)\right| > \epsilon \cdot |\mathcal{C}|(1, \dots, 1)\right) \leq \delta,$$

1249 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.

1251 Before proving Theorem D.8, we use it to argue the claimed runtime of our main result,  
 1252 Theorem D.5.

1253 **Proof of Theorem D.5.** Set  $\mathcal{E} = \text{APPROXIMATE}\tilde{\Phi}(\mathcal{C}, (p_1, \dots, p_n), \delta, \epsilon')$ , where

1254 
$$\epsilon' = \epsilon \cdot \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)},$$

1255 which achieves the claimed error bound on  $\mathcal{E}$  ( $\text{acc}$ ) trivially due to the assignment to  $\epsilon'$  and  
 1256 theorem D.8, since  $\epsilon' \cdot |\mathcal{C}|(1, \dots, 1) = \epsilon \cdot \frac{\tilde{\Phi}(1, \dots, 1)}{|\mathcal{C}|(1, \dots, 1)} \cdot |\mathcal{C}|(1, \dots, 1) = \epsilon \cdot \tilde{\Phi}(1, \dots, 1)$ .

1257 The claim on the runtime follows from Theorem D.8 since

1258 
$$\frac{1}{(\epsilon')^2} \cdot \log\left(\frac{1}{\delta}\right) = \frac{\log \frac{1}{\delta}}{\epsilon^2 \left(\frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)}\right)^2}$$

1259 
$$= \frac{\log \frac{1}{\delta} \cdot |\mathcal{C}|^2(1, \dots, 1)}{\epsilon^2 \cdot \tilde{\Phi}^2(p_1, \dots, p_n)}.$$

1261

1262 Let us now prove Theorem D.8:

1263 **D.2 Proof of Theorem D.8**

1264 **Proof.** Consider now the random variables  $Y_1, \dots, Y_N$ , where each  $Y_i$  is the value of  $Y_i$  in  
 1265 algorithm 1 after line 8 is executed. Overloading  $\text{ISIND}(\cdot)$  to receive monomial input (recall  
 1266  $\mathbf{v}_m$  is the monomial composed of the variables in the set  $\mathbf{v}$ ), we have

$$1267 \quad Y_i = \mathbb{1}_{(\text{ISIND}(\mathbf{v}_m))} \cdot \prod_{X_i \in \text{VAR}(v)} p_i,$$

1268 where the indicator variable handles the check in Line 6 Then for random variable  $Y_i$ , it is  
 1269 the case that

$$1270 \quad \mathbb{E}[Y_i] = \sum_{(\mathbf{v}, \mathbf{c}) \in \mathbb{E}(\mathcal{C})} \frac{\mathbb{1}_{(\text{ISIND}(\mathbf{v}_m))} \cdot \mathbf{c} \cdot \prod_{X_i \in \text{VAR}(v)} p_i}{|\mathcal{C}|(1, \dots, 1)}$$

$$1271 \quad = \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)},$$

1273 where in the first equality we use the fact that  $\text{sgn}_i \cdot |\mathbf{c}| = \mathbf{c}$  and the second equality follows  
 1274 from Eq. (2) with  $X_i$  substituted by  $p_i$ .

1275 Let  $\bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_i$ . It is also true that

$$1276 \quad \mathbb{E}[\bar{Y}] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[Y_i] = \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)}.$$

1277 Hoeffding's inequality states that if we know that each  $Y_i$  (which are all independent)  
 1278 always lie in the intervals  $[a_i, b_i]$ , then it is true that

$$1279 \quad \Pr(|\bar{Y} - \mathbb{E}[\bar{Y}]| \geq \epsilon) \leq 2 \exp\left(-\frac{2N^2\epsilon^2}{\sum_{i=1}^N (b_i - a_i)^2}\right).$$

1280 Line 5 shows that  $\text{sgn}_i$  has a value in  $\{-1, 1\}$  that is multiplied with  $O(k) p_i \in [0, 1]$ ,  
 1281 which implies the range for each  $Y_i$  is  $[-1, 1]$ . Using Hoeffding's inequality, we then get:

$$1282 \quad \Pr(|\bar{Y} - \mathbb{E}[\bar{Y}]| \geq \epsilon) \leq 2 \exp\left(-\frac{2N^2\epsilon^2}{2^2N}\right) = 2 \exp\left(-\frac{N\epsilon^2}{2}\right) \leq \delta,$$

1283 where the last inequality dictates our choice of  $N$  in Line 2.

1284 For the claimed probability bound of  $\Pr\left(|\text{acc} - \tilde{\Phi}(p_1, \dots, p_n)| > \epsilon \cdot |\mathcal{C}|(1, \dots, 1)\right) \leq \delta$ ,  
 1285 note that in the algorithm,  $\text{acc}$  is exactly  $\bar{Y} \cdot |\mathcal{C}|(1, \dots, 1)$ . Multiplying the rest of the terms  
 1286 by the additional factor  $|\mathcal{C}|(1, \dots, 1)$  yields the said bound.

1287 This concludes the proof for the first claim of theorem D.8. Next, we prove the claim on  
 1288 the runtime.

1289 **Run-time Analysis**

1290 The runtime of the algorithm is dominated first by Line 3 (which by Lemma D.6 takes time  
 1291  $O(\text{SIZE}(\mathcal{C}) \cdot \bar{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C}))))$ ) and then by  $N$  iterations of the loop in  
 1292 Line 4. Each iteration's run time is dominated by the call to  $\text{SAMPLEMONOMIAL}$  in Line 5  
 1293 (which by Lemma D.7 takes  $O(\log k \cdot k \cdot \text{DEPTH}(\mathcal{C}) \cdot \bar{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C}))))$ )  
 1294 and the check Line 6, which by the subsequent argument takes  $O(k \log k)$  time. We sort  
 1295 the  $O(k)$  variables by their block IDs and then check if there is a duplicate block ID or not.  
 1296 Combining all the times discussed here gives us the desired overall runtime. ◀

1297 **D.3 Proof of Theorem 4.7**1298 **Proof.** The result follows by first noting that by definition of  $\gamma$ , we have

1299 
$$\tilde{\Phi}(1, \dots, 1) = (1 - \gamma) \cdot |\mathcal{C}|(1, \dots, 1).$$

1300 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

1301 
$$\tilde{\Phi}(1, \dots, 1) \geq p_0^k \cdot \tilde{\Phi}(1, \dots, 1).$$

1302 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  
1303 bound in the runtime bound in Theorem D.5 gives the first claimed runtime. The final  
1304 runtime of  $O_k\left(\frac{1}{\epsilon^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta} \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right)$  follows by noting  
1305 that  $\text{DEPTH}(\mathcal{C}) \leq \text{SIZE}(\mathcal{C})$  and absorbing all factors that just depend on  $k$ . ◀1306 **D.4 Proof of Lemma 4.8**1307 We will prove Lemma 4.8 by considering the two cases separately. We start by considering  
1308 the case when  $\mathcal{C}$  is a tree:1309 ▶ **Lemma D.9.** *Let  $\mathcal{C}$  be a tree (i.e. the sub-circuits corresponding to two children of a node  
1310 in  $\mathcal{C}$  are completely disjoint). Then we have*

1311 
$$|\mathcal{C}|(1, \dots, 1) \leq (\text{SIZE}(\mathcal{C}))^{\text{DEG}(\mathcal{C})+1}.$$

1312 **Proof of Lemma D.9.** For notational simplicity define  $N = \text{SIZE}(\mathcal{C})$  and  $k = \text{DEG}(\mathcal{C})$ . We  
1313 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  
1314 that  $\text{DEPTH}(\mathcal{C}) = 0$ , and there can only be one node which must contain a coefficient or  
1315 constant. In this case,  $|\mathcal{C}|(1, \dots, 1) = 1$ , and  $\text{SIZE}(\mathcal{C}) = 1$ , and by Definition 4.4 it is the  
1316 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  
1317  $k \in \{0, 1\}$ .1318 Assume for  $\ell > 0$  an arbitrary circuit  $\mathcal{C}$  of  $\text{DEPTH}(\mathcal{C}) \leq \ell$  that it is true that  $|\mathcal{C}|(1, \dots, 1) \leq$   
1319  $N^{k+1}$ .1320 For the inductive step we consider a circuit  $\mathcal{C}$  such that  $\text{DEPTH}(\mathcal{C}) = \ell + 1$ . The sink can  
1321 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  
1322 when sink node is  $\times$ . Then note that

1323 
$$\begin{aligned} |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) \cdot |\mathcal{C}_R|(1, \dots, 1) \\ 1324 &\leq (N - 1)^{k_L+1} \cdot (N - 1)^{k_R+1} \\ 1325 &= (N - 1)^{k+1} \\ 1326 &\leq N^{k+1}. \end{aligned} \tag{30}$$

1328 In the above the first inequality follows from the inductive hypothesis (and the fact that the  
1329 size of either subtree is at most  $N - 1$ ) and Eq. (30) follows by definition 4.4 which states  
1330 that for  $k = \text{DEG}(\mathcal{C})$  we have  $k = k_L + k_R + 1$ .1331 For the case when the sink gate is a  $+$  gate, then for  $N_L = \text{SIZE}(\mathcal{C}_L)$  and  $N_R = \text{SIZE}(\mathcal{C}_R)$   
1332 we have

1333 
$$\begin{aligned} |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) + |\mathcal{C}_R|(1, \dots, 1) \\ 1334 &\leq N_L^{k+1} + N_R^{k+1} \\ 1335 &\leq (N - 1)^{k+1} \end{aligned} \tag{31}$$

$$\leq N^{k+1}.$$

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

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

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

$$|\mathbf{C}|(1, \dots, 1) \leq 2^{2^{\text{DEGC}(\mathbf{C})} \cdot \text{DEPTH}(\mathbf{C})}.$$

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

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

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

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  $k \geq 0$ . ◀

## 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):

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

1376  
1377

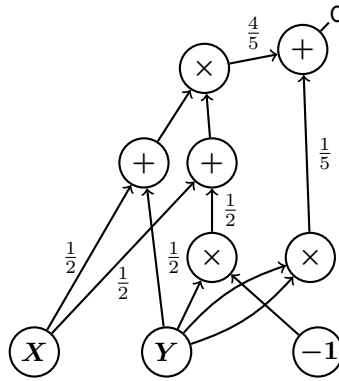
1378 It turns out that for proof of Lemma D.7, we need to argue that when  $C.type = +$ , we  
1379 indeed have

1380 
$$C.Lweight \leftarrow \frac{|C_L|(1, \dots, 1)}{|C_L|(1, \dots, 1) + |C_R|(1, \dots, 1)}; \tag{33}$$

1381 
$$C.Rweight \leftarrow \frac{|C_R|(1, \dots, 1)}{|C_L|(1, \dots, 1) + |C_R|(1, \dots, 1)} \tag{34}$$
  
1382  
1383

1384 **D.6 OnePass Example**

1385 ► **Example D.11.** Let  $T$  encode the expression  $(X + Y)(X - Y) + Y^2$ . After one pass,  
1386 Algorithm 2 would have computed the following weight distribution. For the two inputs of the  
1387 sink gate  $C$ ,  $C.Lweight = \frac{4}{5}$  and  $C.Rweight = \frac{1}{5}$ . Similarly, for  $S$  denoting the left input of  
1388  $C_L$ ,  $S.Lweight = S.Rweight = \frac{1}{2}$ . This is depicted in Fig. 4.



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

1389 **D.7 Proof of OnePass (Lemma D.6)**

1390 **Proof.** We prove the correct computation of **partial**, **Lweight**, **Rweight** values on  $C$  by  
1391 induction over the number of iterations in the topological order **TOPORD** (line 1) of the  
1392 input circuit  $C$ . **TOPORD** follows the standard definition of a topological ordering over the  
1393 DAG structure of  $C$ .

1394 For the base case, we have only one gate, which by definition is a source gate and must be  
1395 either **VAR** or **NUM**. In this case, as per eq. (32), lines 3 and 5 correctly compute **C.partial**  
1396 as 1.

1397 For the inductive hypothesis, assume that ONEPASS correctly computes **S.partial**,  
1398 **S.Lweight**, and **S.Rweight** for all gates  $g$  in  $C$  with  $k \geq 0$  iterations over **TOPORD**.

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

1399

1400 We now prove for  $k+1$  iterations that ONEPASS correctly computes the **partial**, **Lweight**,  
1401 and **Rweight** values for each gate  $g_i$  in  $C$  for  $i \in [k + 1]$ . The  $g_{k + 1}$  must be in the last  
1402 ordering of all gates  $g_i$ . When  $SIZE(C) > 1$ , if  $g_{k+1}$  is a leaf node, we are back to the base  
1403 case. Otherwise  $g_{k+1}$  is an internal node which requires binary input.

---

**Algorithm 2** ONEPASS ( $\mathcal{C}$ )
 

---

**Input:**  $\mathcal{C}$ : Circuit**Output:**  $\mathcal{C}$ : Annotated Circuit**Output:**  $\text{sum} \in \mathbb{N}$ 

```

1: for  $g$  in  $\text{TOPORD}(\mathcal{C})$  do ▷  $\text{TOPORD}(\cdot)$  is the topological order of  $\mathcal{C}$ 
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 ( $\text{sum}, \mathcal{C}$ )

```

---

1404 When  $g_{k+1}.\text{type} = +$ , then by line 9  $g_{k+1}.\text{partial} = g_{k+1_L}.\text{partial} + g_{k+1_R}.\text{partial}$ ,  
 1405 a correct computation, as per eq. (32). Further, lines 10 and 11 compute  $g_{k+1}.L\text{weight} =$   
 1406  $\frac{g_{k+1_L}.\text{partial}}{g_{k+1}.\text{partial}}$  and analogously for  $g_{k+1}.R\text{weight}$ . All values needed for each computation have  
 1407 been correctly computed by the inductive hypothesis.

1408 When  $g_{k+1}.\text{type} = \times$ , then line 7 computes  $g_{k+1}.\text{partial} = g_{k+1_L}.\text{partial} \times g_{k+1_R}.\text{partial}$ ,  
 1409 which indeed by eq. (32) is correct. This concludes the proof of correctness.

### 1410 Runtime Analysis

1411 It is known that  $\text{TOPORD}(G)$  is computable in linear time. There are  $\text{SIZE}(\mathcal{C})$  iterations, each  
 1412 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  
 1413 all the numbers which the algorithm computes is at most  $|\mathcal{C}|(1, \dots, 1)$ . Hence, by definition  
 1414 each such operation takes  $\overline{\mathcal{M}}(\log(|\mathcal{C}(1 \dots, 1)|), \log \text{SIZE}(\mathcal{C}))$  time, which proves the claimed  
 1415 runtime. ◀

## 1416 D.8 SampleMonomial Remarks

1417 We briefly describe the top-down traversal of `SAMPLEMONOMIAL`. When  $\mathcal{C}.\text{type} = +$ , the  
 1418 input to be visited is sampled from the weighted distribution precomputed by `ONEPASS`.  
 1419 When a  $\mathcal{C}.\text{type} = \times$  node is visited, both inputs are visited. The algorithm computes two  
 1420 properties: the set of all variable leaf nodes visited, and the product of the signs of visited  
 1421 coefficient leaf nodes. We will assume the `TreeSet` data structure to maintain sets with  
 1422 logarithmic time insertion and linear time traversal of its elements. While we would like to  
 1423 take advantage of the space efficiency gained in using a circuit  $\mathcal{C}$  instead an expression tree  $T$ ,  
 1424 we do not know that such a method exists when computing a sample of the input polynomial  
 1425 representation.

1426 The efficiency gains of circuits over trees is found in the capability of circuits to only  
 1427 require space for each *distinct* term in the compressed representation. This saves space



---

**Algorithm 3** SAMPLEMONOMIAL ( $\mathcal{C}$ )
 

---

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

```

1:  $\text{vars} \leftarrow \emptyset$ 
2: if  $\mathcal{C}.\text{type} = +$  then  $\triangleright$  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  $\mathcal{C}.\text{Rweight}$ .  $\triangleright$  Each call to SAMPLEMONOMIAL uses fresh randomness
4:    $(\mathbf{v}, \mathbf{s}) \leftarrow$  SAMPLEMONOMIAL( $\mathcal{C}_{\text{samp}}$ )
5:   return  $(\mathbf{v}, \mathbf{s})$ 
6: else if  $\mathcal{C}.\text{type} = \times$  then  $\triangleright$  Multiply the sampled values of all inputs
7:    $\text{sgn} \leftarrow 1$ 
8:   for  $\text{input}$  in  $\mathcal{C}.\text{input}$  do
9:      $(\mathbf{v}, \mathbf{s}) \leftarrow$  SAMPLEMONOMIAL( $\text{input}$ )
10:     $\text{vars} \leftarrow \text{vars} \cup \{\mathbf{v}\}$ 
11:     $\text{sgn} \leftarrow \text{sgn} \times \mathbf{s}$ 
12:   end for
13:   return  $(\text{vars}, \text{sgn})$ 
14: else if  $\mathcal{C}.\text{type} = \text{NUM}$  then  $\triangleright$  The leaf is a coefficient
15:   return  $(\{\}, \text{SGN}(\mathcal{C}.\text{val}))$   $\triangleright$   $\text{SGN}(\cdot)$  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

```

---

1428 in such polynomials containing non-distinct terms multiplied or added to each other, e.g.,  
 1429  $x^4$ . However, to avoid biased sampling, it is imperative to sample from both inputs of a  
 1430 multiplication gate, independently, which is indeed the approach of SAMPLEMONOMIAL.

## 1431 D.9 Proof of SampleMonomial (Lemma D.7)

1432 **Proof.** We first need to show that SAMPLEMONOMIAL samples a valid monomial  $\mathbf{v}_m$  by  
 1433 sampling and returning a set of variables  $\mathbf{v}$ , such that  $(\mathbf{v}, c)$  is in  $\mathbf{E}(\mathcal{C})$  and  $\mathbf{v}_m$  is indeed a  
 1434 monomial of the  $\tilde{\Phi}(\mathbf{X})$  encoded in  $\mathcal{C}$ . We show this via induction over the depth of  $\mathcal{C}$ .

1435 For the base case, let the depth  $d$  of  $\mathcal{C}$  be 0. We have that the single gate is either a  
 1436 constant  $c$  for which by line 15 we return  $\{\}$ , or we have that  $\mathcal{C}.\text{type} = \text{VAR}$  and  $\mathcal{C}.\text{val} = x$ ,  
 1437 and by line 17 we return  $\{x\}$ . By definition 4.1, both cases return a valid  $\mathbf{v}$  for some  $(\mathbf{v}, c)$   
 1438 from  $\mathbf{E}(\mathcal{C})$ , and the base case is proven.

1439 For the inductive hypothesis, assume that for  $d \leq k$  for some  $k \geq 0$ , that it is indeed the  
 1440 case that SAMPLEMONOMIAL returns a valid monomial.

1441 For the inductive step, let us take a circuit  $\mathcal{C}$  with  $d = k + 1$ . Note that each input  
 1442 has depth  $d - 1 \leq k$ , and by inductive hypothesis both of them sample a valid monomial.  
 1443 Then the sink can be either a  $+$  or  $\times$  gate. For the case when  $\mathcal{C}.\text{type} = +$ , line 3 of  
 1444 SAMPLEMONOMIAL will choose one of the inputs of the source. By inductive hypothesis it is  
 1445 the case that some valid monomial is being randomly sampled from each of the inputs. Then  
 1446 it follows when  $\mathcal{C}.\text{type} = +$  that a valid monomial is sampled by SAMPLEMONOMIAL. When  
 1447 the  $\mathcal{C}.\text{type} = \times$ , line 10 computes the set union of the monomials returned by the two inputs  
 1448 of the sink, and it is trivial to see by definition 4.1 that  $\mathbf{v}_m$  is a valid monomial encoded by  
 1449 some  $(\mathbf{v}, c)$  of  $\mathbf{E}(\mathcal{C})$ .

1450 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  
1451 with a probability  $\frac{|\mathbf{c}|}{|\mathbf{C}|(1, \dots, 1)}$ .

1452 For the base case  $d = 0$ , by definition 2.6 we know that the  $\text{SIZE}(\mathbf{C}) = 1$  and  $\mathbf{C.type} =$   
1453  $\text{NUM}$  or  $\text{VAR}$ . For either case, the probability of the value returned is 1 since there is only  
1454 one value to sample from. When  $\mathbf{C.val} = x$ , the algorithm always return the variable set  
1455  $\{x\}$ . When  $\mathbf{C.type} = \text{NUM}$ ,  $\text{SAMPLEMONOMIAL}$  will always return  $\emptyset$ .

1456 For the inductive hypothesis, assume that for  $d \leq k$  and  $k \geq 0$   $\text{SAMPLEMONOMIAL}$  indeed  
1457 returns  $\mathbf{v}$  in  $(\mathbf{v}, \mathbf{c})$  of  $\mathbf{E}(\mathbf{C})$  with probability  $\frac{|\mathbf{c}|}{|\mathbf{C}|(1, \dots, 1)}$ .

1458 We prove now for  $d = k + 1$  the inductive step holds. It is the case that the sink of  $\mathbf{C}$  has  
1459 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,  
1460  $\text{SAMPLEMONOMIAL}$  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  
1461  $\mathbf{C}_R$  with probability  $\frac{|\mathbf{c}_L|}{|\mathbf{C}_L|(1, \dots, 1)}$  and  $\frac{|\mathbf{c}_R|}{|\mathbf{C}_R|(1, \dots, 1)}$ .

1462 Consider the case when  $\mathbf{C.type} = \times$ . For the term  $(\mathbf{v}, \mathbf{c})$  from  $\mathbf{E}(\mathbf{C})$  that is being sampled  
1463 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  
1464 that  $\text{SAMPLEMONOMIAL}(\mathbf{C}_L)$  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  $\mathbf{v}_L$   
1465 and  $\mathbf{v}_R$  are sampled with independent randomness, the final probability for sample  $\mathbf{v}$  is  
1466 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  
1467  $|\mathbf{c}| = |\mathbf{c}_{v_L}| \cdot |\mathbf{c}_{v_R}|$  and that (as shown in eq. (32))  $|\mathbf{C}|(1, \dots, 1) = |\mathbf{C}_L|(1, \dots, 1) \cdot |\mathbf{C}_R|(1, \dots, 1)$ ,  
1468 and therefore  $\mathbf{v}$  is sampled with correct probability  $\frac{|\mathbf{c}|}{|\mathbf{C}|(1, \dots, 1)}$ .

1469 For the case when  $\mathbf{C.type} = +$ ,  $\text{SAMPLEMONOMIAL}$  will sample  $\mathbf{v}$  from one of its inputs.  
1470 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  
1471 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  $\mathbf{v}$ ,  
1472 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  
1473 a symmetric argument holds for the case when  $\mathbf{v}$  is sampled from  $\mathbf{C}_R$ . Notice also that the  
1474 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  $\text{ONEPASS}$ . Then,  
1475 since  $\text{SAMPLEMONOMIAL}$  goes top-down, and each sampling choice is independent (which  
1476 follows from the randomness in the root of  $\mathbf{C}$  being independent from the randomness used  
1477 in its subtrees), the probability for  $\mathbf{v}$  to be sampled from  $\mathbf{C}$  is equal to the product of the  
1478 probability that  $\mathbf{C}_L$  is sampled from  $\mathbf{C}$  and  $\mathbf{v}$  is sampled in  $\mathbf{C}_L$ , and

$$\begin{aligned}
 1479 \quad & Pr(\text{SAMPLEMONOMIAL}(\mathbf{C}) = \mathbf{v}) = \\
 1480 \quad & Pr(\text{SAMPLEMONOMIAL}(\mathbf{C}_L) = \mathbf{v}) \cdot Pr(\text{SampledChild}(\mathbf{C}) = \mathbf{C}_L) \\
 1481 \quad & = \frac{|\mathbf{c}_v|}{|\mathbf{C}_L|(1, \dots, 1)} \cdot \frac{|\mathbf{C}_L|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)} \\
 1482 \quad & = \frac{|\mathbf{c}_v|}{|\mathbf{C}|(1, \dots, 1)}, \\
 1483
 \end{aligned}$$

1484 and we obtain the desired result.

1485 Lastly, we show by simple induction of the depth  $d$  of  $\mathbf{C}$  that  $\text{SAMPLEMONOMIAL}$  indeed  
1486 returns the correct sign value of  $\mathbf{c}$  in  $(\mathbf{v}, \mathbf{c})$ .

1487 In the base case,  $\mathbf{C.type} = \text{NUM}$  or  $\text{VAR}$ . For the former,  $\text{SAMPLEMONOMIAL}$  correctly  
1488 returns the sign value of the gate. For the latter,  $\text{SAMPLEMONOMIAL}$  returns the correct  
1489 sign of 1, since a variable is a neutral element, and 1 is the multiplicative identity, whose  
1490 product with another sign element will not change that sign element.

1491 For the inductive hypothesis, we assume for a circuit of depth  $d \leq k$  and  $k \geq 0$  that the  
1492 algorithm correctly returns the sign value of  $\mathbf{c}$ .

1493 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  
1494 sign of  $\mathbf{c}$ . For the case that  $\mathbf{C.type} = \times$ , the sign value of both inputs are multiplied, which

1495 is the correct behavior by definition 4.1. When  $\mathbf{C.type} = +$ , only one input of  $\mathbf{C}$  is sampled,  
1496 and the algorithm returns the correct sign value of  $c$  by inductive hypothesis.

### 1497 Run-time Analysis

1498 It is easy to check that except for lines 3 and 10, all lines take  $O(1)$  time. Consider an  
1499 execution of line 10. We note that we will be adding a given set of variables to some set at  
1500 most once: since the sum of the sizes of the sets at a given level is at most  $\text{DEG}(\mathbf{C})$ , each gate  
1501 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  
1502  $a = \mathbf{C.Lweight}$  and  $b = \mathbf{C.Rweight}$ . We can implement this step by picking a random number  
1503  $r \in [a + b]$  and then checking if  $r \leq a$ . It is easy to check that  $a + b \leq |\mathbf{C}|(1, \dots, 1)$ . This  
1504 means we need to add and compare  $\log |\mathbf{C}|(1, \dots, 1)$ -bit numbers, which can certainly be done  
1505 in time  $\overline{\mathcal{M}}(\log(|\mathbf{C}|(1, \dots, 1)), \log \text{SIZE}(\mathbf{C}))$  (note that this is an over-estimate). Denote  $\text{COST}$   
1506  $(\mathbf{C})$  (Eq. (35)) to be an upper bound of the number of gates visited by  $\text{SAMPLEMONOMIAL}$ .  
1507 Then the runtime is  $O(\text{COST}(\mathbf{C}) \cdot \log \text{DEG}(\mathbf{C}) \cdot \overline{\mathcal{M}}(\log(|\mathbf{C}|(1, \dots, 1)), \log \text{SIZE}(\mathbf{C})))$ .

1508 We now bound the number of recursive calls in  $\text{SAMPLEMONOMIAL}$  by  $O((\text{DEG}(\mathbf{C}) + 1) \cdot$   
1509  $\text{DEPTH}(\mathbf{C}))$ , which by the above will prove the claimed runtime.

1510 Let  $\text{COST}(\cdot)$  be a function that models an upper bound on the number of gates that can  
1511 be visited in the run of  $\text{SAMPLEMONOMIAL}$ . We define  $\text{COST}(\cdot)$  recursively as follows.

$$1512 \quad \text{COST}(\mathbf{C}) = \begin{cases} 1 + \text{COST}(\mathbf{C}_L) + \text{COST}(\mathbf{C}_R) & \text{if } \mathbf{C.type} = \times \\ 1 + \max(\text{COST}(\mathbf{C}_L), \text{COST}(\mathbf{C}_R)) & \text{if } \mathbf{C.type} = + \\ 1 & \text{otherwise} \end{cases} \quad (35)$$

1513 First note that the number of gates visited in  $\text{SAMPLEMONOMIAL}$  is  $\leq \text{COST}(\mathbf{C})$ . To show  
1514 that eq. (35) upper bounds the number of nodes visited by  $\text{SAMPLEMONOMIAL}$ , note that  
1515 when  $\text{SAMPLEMONOMIAL}$  visits a gate such that  $\mathbf{C.type} = \times$ , line 8 visits each input of  $\mathbf{C}$ , as  
1516 defined in (35). For the case when  $\mathbf{C.type} = +$ , line 3 visits exactly one of the input gates,  
1517 which may or may not be the subcircuit with the maximum number of gates traversed, which  
1518 makes  $\text{COST}(\cdot)$  an upperbound. Finally, it is trivial to see that when  $\mathbf{C.type} \in \{\text{VAR}, \text{NUM}\}$ ,  
1519 i.e., a source gate, that only one gate is visited.

1520 We prove the following inequality holds.

$$1521 \quad 2(\text{DEG}(\mathbf{C}) + 1) \cdot \text{DEPTH}(\mathbf{C}) + 1 \geq \text{COST}(\mathbf{C}) \quad (36)$$

1522 Note that eq. (36) implies the claimed runtime. We prove eq. (36) for the number of  
1523 gates traversed in  $\text{SAMPLEMONOMIAL}$  using induction over  $\text{DEPTH}(\mathbf{C})$ . Recall how degree is  
1524 defined in definition 4.4.

1525 For the base case  $\text{DEG}(\mathbf{C}) = \{0, 1\}$ ,  $\text{DEPTH}(\mathbf{C}) = 0$ ,  $\text{COST}(\mathbf{C}) = 1$ , and it is trivial to see  
1526 that the inequality  $2\text{DEG}(\mathbf{C}) \cdot \text{DEPTH}(\mathbf{C}) + 1 \geq \text{COST}(\mathbf{C})$  holds.

1527 For the inductive hypothesis, we assume the bound holds for any circuit where  $\ell \geq$   
1528  $\text{DEPTH}(\mathbf{C}) \geq 0$ . Now consider the case when  $\text{SAMPLEMONOMIAL}$  has an arbitrary circuit  $\mathbf{C}$   
1529 input with  $\text{DEPTH}(\mathbf{C}) = \ell + 1$ . By definition  $\mathbf{C.type} \in \{+, \times\}$ . Note that since  $\text{DEPTH}(\mathbf{C}) \geq 1$ ,  
1530  $\mathbf{C}$  must have input(s). Further we know that by the inductive hypothesis the inputs  $\mathbf{C}_i$  for  
1531  $i \in \{L, R\}$  of the sink gate  $\mathbf{C}$  uphold the bound

$$1532 \quad 2(\text{DEG}(\mathbf{C}_i) + 1) \cdot \text{DEPTH}(\mathbf{C}_i) + 1 \geq \text{COST}(\mathbf{C}_i). \quad (37)$$

1533 In particular, since for any  $i$ , eq. (37) holds, then it immediately follows that an inequality  
1534 whose operands consist of a sum of the aforementioned inequalities must also hold. This is

## 23:44 Bag PDB Queries

1535 readily seen in the inequality of eq. (39) and eq. (40), where  $2(\text{DEG}(\mathbf{C}_L) + 1) \cdot \text{DEPTH}(\mathbf{C}_L) \geq$   
 1536  $\text{COST}(\mathbf{C}_L)$ , likewise for  $\mathbf{C}_R$ , and  $1 \geq 1$ . It is also true that  $\text{DEPTH}(\mathbf{C}_L) \leq \text{DEPTH}(\mathbf{C}) - 1$  and  
 1537  $\text{DEPTH}(\mathbf{C}_R) \leq \text{DEPTH}(\mathbf{C}) - 1$ .

1538 If  $\mathbf{C.type} = +$ , then  $\text{DEG}(\mathbf{C}) = \max(\text{DEG}(\mathbf{C}_L), \text{DEG}(\mathbf{C}_R))$ . Otherwise  $\mathbf{C.type} = \times$  and  
 1539  $\text{DEG}(\mathbf{C}) = \text{DEG}(\mathbf{C}_L) + \text{DEG}(\mathbf{C}_R) + 1$ . In either case it is true that  $\text{DEPTH}(\mathbf{C}) = \max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) +$   
 1540  $1$ .

1541 If  $\mathbf{C.type} = \times$ , then, by eq. (35), substituting values, the following should hold,

$$1542 \quad 2(\text{DEG}(\mathbf{C}_L) + \text{DEG}(\mathbf{C}_R) + 2) \cdot (\max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) + 1) + 1 \quad (38)$$

$$1543 \quad \geq 2(\text{DEG}(\mathbf{C}_L) + 1) \cdot \text{DEPTH}(\mathbf{C}_L) + 2(\text{DEG}(\mathbf{C}_R) + 1) \cdot \text{DEPTH}(\mathbf{C}_R) + 3 \quad (39)$$

$$1544 \quad \geq 1 + \text{COST}(\mathbf{C}_L) + \text{COST}(\mathbf{C}_R) = \text{COST}(\mathbf{C}). \quad (40)$$

1546 To prove (39), first, eq. (38) expands to,

$$1547 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) + 4 + 1 \quad (41)$$

1548 where  $\text{DEPTH}_{\max}$  is used to denote the maximum depth of the two input subcircuits. Eq. (39)  
 1549 expands to

$$1550 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_R) + 3 \quad (42)$$

1551 Putting Eq. (41) and Eq. (42) together we get

$$1552 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) + 5$$

$$1553 \quad \geq 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_R) + 3 \quad (43)$$

1554 Since the following is always true,

$$1556 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 5$$

$$1557 \quad \geq 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_R) + 3,$$

1559 then it is the case that Eq. (43) is *always* true.

1560 Now to justify (40) which holds for the following reasons. First, eq. (40) is the result of  
 1561 Eq. (35) when  $\mathbf{C.type} = \times$ . Eq. (39) is then produced by substituting the upperbound of  
 1562 (37) for each  $\text{COST}(\mathbf{C}_i)$ , trivially establishing the upper bound of (40). This proves eq. (36)  
 1563 for the  $\times$  case.

1564 For the case when  $\mathbf{C.type} = +$ , substituting values yields

$$1565 \quad 2(\max(\text{DEG}(\mathbf{C}_L), \text{DEG}(\mathbf{C}_R)) + 1) \cdot (\max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) + 1) + 1 \quad (44)$$

$$1566 \quad \geq \max(2(\text{DEG}(\mathbf{C}_L) + 1) \cdot \text{DEPTH}(\mathbf{C}_L) + 1, 2(\text{DEG}(\mathbf{C}_R) + 1) \cdot \text{DEPTH}(\mathbf{C}_R) + 1) + 1 \quad (45)$$

$$1567 \quad \geq 1 + \max(\text{COST}(\mathbf{C}_L), \text{COST}(\mathbf{C}_R)) = \text{COST}(\mathbf{C}) \quad (46)$$

1569 To prove (45), eq. (44) expands to

$$1570 \quad 2\text{DEG}_{\max} \text{DEPTH}_{\max} + 2\text{DEG}_{\max} + 2\text{DEPTH}_{\max} + 2 + 1. \quad (47)$$

1571 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  
 1572 expansion of eq. (45):

$$1573 \quad 2\text{DEG}_{\max} \text{DEPTH}_{\max} + 2\text{DEPTH}_{\max} + 2 \quad (48)$$

1574 Putting it together we obtain the following for (45):

$$\begin{aligned}
 1575 \quad & 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEG}_{\max} + 2\text{DEPTH}_{\max} + 3 \\
 1576 \quad & \geq 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEPTH}_{\max} + 2, \tag{49} \\
 1577
 \end{aligned}$$

1578 where it can be readily seen that the inequality stands and (49) follows. This proves (45).

1579 Similar to the case of `C.type = ×`, (46) follows by equations (35) and (37).

1580 This proves (36) as desired. ◀

## 1581 D.10 Experimental Results

1582 Recall that by definition of BIDB, a query result cannot be derived by a self-join between  
 1583 non-identical tuples belonging to the same block. Note, that by Theorem 4.7,  $\gamma$  must be  
 1584 a constant in order for Algorithm 1 to achieve linear time. We would like to determine  
 1585 experimentally whether queries over BIDB instances in practice generate a constant number  
 1586 of cancellations or not. Such an experiment would ideally use a database instance with  
 1587 queries both considered to be typical representations of what is seen in practice.

1588 We ran our experiments using Windows 10 WSL Operating System with an Intel Core i7  
 1589 2.40GHz processor and 16GB RAM. All experiments used the PostgreSQL 13.0 database  
 1590 system.

1591 For the data we used the MayBMS data generator [1] tool to randomly generate uncertain  
 1592 versions of TPC-H tables. The queries computed over the database instance are  $Q_1$ ,  $Q_2$ , and  
 1593  $Q_3$  from [5], all of which are modified versions of TPC-H queries  $Q_3$ ,  $Q_6$ , and  $Q_7$  where all  
 1594 aggregations have been dropped.

1595 As written, the queries disallow BIDB cross terms. We first ran all queries, noting the  
 1596 result size for each. Next the queries were rewritten so as not to filter out the cross terms.  
 1597 The comparison of the sizes of both result sets should then suggest in one way or another  
 1598 whether or not there exist many cross terms in practice. As seen, the experimental query  
 1599 results contain little to no cancelling terms. Fig. 5 shows the result sizes of the queries,  
 1600 where column CF is the result size when all cross terms are filtered out, column CI shows  
 1601 the number of output tuples when the cancelled tuples are included in the result, and the  
 1602 last column is the value of  $\gamma$ . The experiments show  $\gamma$  to be in a range between  $[0, 0.1]\%$ ,  
 1603 indicating that only a negligible or constant (compare the result sizes of  $Q_1 < Q_2$  and their  
 1604 respective  $\gamma$  values) amount of tuples are cancelled in practice when running queries over a  
 1605 typical BIDB instance. Interestingly, only one of the three queries had tuples that violated  
 1606 the BIDB constraint.

1607 To conclude, the results in Fig. 5 show experimentally that  $\gamma$  is negligible in practice for  
 1608 BIDB queries. We also observe that (i) tuple presence is independent across blocks, so the  
 1609 corresponding probabilities (and hence  $p_0$ ) are independent of the number of blocks, and (ii)  
 1610 BIDBs model uncertain attributes, so block size (and hence  $\gamma$ ) is a function of the “messiness”  
 1611 of a dataset, rather than its size. Thus, we expect Theorem 4.7 to hold in general.

Query	CF	CI	$\gamma$
$Q_1$	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 BIDB.

## E Circuits

### E.1 Representing Polynomials with Circuits

#### 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.

We now formalize circuits and the construction of circuits for  $\mathcal{RA}^+$  queries. As mentioned earlier, we represent lineage polynomials as arithmetic circuits over  $\mathbb{N}$ -valued variables with  $+$ ,  $\times$ . A circuit for query  $Q$  and  $\mathbb{N}[\mathbf{X}]$ -encoded PDB  $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$  is a directed acyclic graph  $\langle V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle$  with vertices  $V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}$  and directed edges  $E_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \subset V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}^2$ . 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 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.

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 (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 label to each node: Source nodes (i.e., vertices with in-degree 0) are labeled with constants or variables (i.e.,  $\mathbb{N} \cup \mathbf{X}$ ), while the remaining nodes are labeled with the symbol  $+$  or  $\times$ . We require that vertices have an in-degree of at most two. Note that we can construct circuits for BIBDs in time linear in the time required for deterministic query processing over a possible world of the BIBD under the aforementioned assumption that  $|\mathcal{D}_{\mathbb{N}[\mathbf{X}]}| \leq c \cdot |D|$ .

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

### E.2 Modeling Circuit Construction

We now connect the size of a circuit (where the size of a circuit is the number of vertices 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 runtime  $T_{det}^*(Q, D_{\overline{\Omega}})$  of the PDB's deterministic bounding database  $D_{\overline{\Omega}}$ . We do this formally by showing that the size of the circuit is asymptotically no worse than the corresponding runtime of a large class of deterministic query processing algorithms.

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

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

encodes a polynomial, realized as

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

We define the circuit for a  $\mathcal{RA}^+$  query  $Q$  recursively by cases as follows. In each case, let  $\langle V_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \rangle$  denote the circuit for subquery  $Q_i$ . We implicitly 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  $Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ .

---

**Algorithm 4**  $LC(Q, D_{\bar{\Omega}}, E, V, \ell)$ 


---

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

```

1: if  $Q$  is  $R$  then                                     ▷ Case 1:  $Q$  is a relation atom
2:   for  $t \in D_{\bar{\Omega}}.R$  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:    $\langle V, E, \phi', \ell \rangle \leftarrow LC(Q', D_{\bar{\Omega}}, V, E, \ell)$ 
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:    $\langle V, E, \phi', \ell \rangle \leftarrow LC(Q', D_{\bar{\Omega}}, V, E, \ell)$ 
13:   for  $t \in \pi_{\bar{A}}(Q'(D_{\bar{\Omega}}))$  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  $t \in Q'(D_{\bar{\Omega}})$  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:    $\langle V, E, \phi_1, \ell \rangle \leftarrow LC(Q_1, D_{\bar{\Omega}}, V, E, \ell)$ 
23:    $\langle V, E, \phi_2, \ell \rangle \leftarrow LC(Q_2, D_{\bar{\Omega}}, V, E, \ell)$ 
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:      $\langle V, E, \phi_i, \ell \rangle \leftarrow LC(Q_i, D_{\bar{\Omega}}, V, E, \ell)$ 
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:      $E \leftarrow E \cup \left\{ (\phi_i(\pi_{sch(Q_i(D_{\bar{\Omega}})}(t))), v_t) \mid i \in [n] \right\}$ 
38:   end for
39:   Correct nodes with in-degrees  $> 2$  by appending an equivalent fan-in two tree instead
40: end if

```

---

1644 Algorithm 4 defines how the circuit for a query result is constructed. We quickly review  
1645 the number of vertices emitted in each case.

1646 **Base Relation.** This circuit has  $|D_{\bar{\Omega}}.R|$  vertices.

1647 **Selection.** If we assume dead sinks are iteratively garbage collected, this circuit has at  
1648 most  $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}|$  vertices.

1649 **Projection.** This formulation will produce vertices with an in-degree greater than two, a  
1650 problem that we correct by replacing every vertex with an in-degree over two by an equivalent  
1651 fan-in two tree. The resulting structure has at most  $|Q_1| - 1$  new vertices. The corrected  
1652 circuit thus has at most  $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| + |Q_1|$  vertices.

1653 **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.

1654  **$k$ -ary Join.** As in projection, newly created vertices will have an in-degree of  $k$ , and a  
1655 fan-in two tree is required. There are  $|Q_1 \bowtie \dots \bowtie Q_k|$  such vertices, so the corrected circuit  
1656 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.

## 1657 E.2.1 Bounding circuit depth

1658 We first show that the depth of the circuit (DEPTH; Definition 4.3) is bounded by the size  
1659 of the query. Denote by  $|Q|$  the number of relational operators in query  $Q$ , which recall we  
1660 assume is a constant.

1661 **► Proposition E.1** (Circuit depth is bounded). *Let  $Q$  be a relational query and  $D_{\bar{\Omega}}$  be a*  
1662 *deterministic bounding database with  $n$  tuples. There exists a (lineage) circuit  $\mathcal{C}^*$  encoding*  
1663 *the lineage of all tuples  $t \in Q(D_{\bar{\Omega}})$  for which  $\text{DEPTH}(\mathcal{C}^*) \leq O(k|Q| \log(n))$ .*

1664 **Proof.** We show that the bound of Proposition E.1 holds for the circuit constructed by  
1665 Algorithm 4. First, observe that Algorithm 4 is (recursively) invoked exactly once for every  
1666 relational operator or base relation in  $Q$ ; It thus suffices to show that a call to Algorithm 4  
1667 adds at most  $O_k(\log(n))$  to the depth of a circuit produced by any recursive invocation.  
1668 Second, observe that modulo the logarithmic fan-in of the projection and join cases, the  
1669 depth of the output is at most one greater than the depth of any input (or at most 1 in the  
1670 base case of relation atoms). For the join case, the number of in-edges can be no greater than  
1671 the join width, which itself is bounded by  $k$ . The depth thus increases by at most a constant  
1672 factor of  $\lceil \log(k) \rceil = O_k(1)$ . For the projection case, observe that the fan-in is bounded by  
1673  $|Q'(D_{\bar{\Omega}})|$ , which is in turn bounded by  $n^k$ . The depth increase for any projection node is  
1674 thus at most  $\lceil \log(n^k) \rceil = O(k \log(n))$ , as desired. ◀

## 1675 E.2.2 Circuit size vs. runtime

1676 **► Lemma E.2.** *Given a  $\mathbb{N}[\mathbf{X}]$ -encoded PDB  $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$  with deterministic bounding database  $D_{\bar{\Omega}}$ ,*  
1677 *and an  $\mathcal{RA}^+$  query  $Q$ , the runtime of  $Q$  over  $D_{\bar{\Omega}}$  has the same or greater complexity as the*  
1678 *size of the lineage of  $Q(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ . That is, we have  $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq kT_{det}^*(Q, D_{\bar{\Omega}}) + 1$ , where  $k \geq 1$*   
1679 *is the maximal degree of any polynomial in  $Q(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ .*

1680 **Proof.** We prove by induction that  $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}} \setminus \{v_0\}| \leq kT_{det}^*(Q, D_{\bar{\Omega}})$ . For clarity, we  
1681 implicitly exclude  $v_0$  in the proof below.

1682 The base case is a base relation:  $Q = R$  and is trivially true since  $|V_{R, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| = |D_{\bar{\Omega}}.R| =$   
1683  $T_{det}^*(R, D_{\bar{\Omega}})$  (note that here the degree  $k = 1$ ). For the inductive step, we assume that we  
1684 have circuits for subqueries  $Q_1, \dots, Q_m$  such that  $|V_{Q_i, \mathcal{D}_{\mathbb{N}[\mathbf{X}]}}| \leq k_i T_{det}^*(Q_i, D_{\bar{\Omega}})$  where  $k_i$  is  
1685 the degree of  $Q_i$ .



1686 **Selection.** Assume that  $Q = \sigma_\theta(Q_1)$ . In the circuit for  $Q$ ,  $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| = |V_{Q_1, D_{\overline{\Omega}}}|$  vertices,  
 1687 so from the inductive assumption and  $T_{det}^*(Q, D_{\overline{\Omega}}) = T_{det}^*(Q_1, D_{\overline{\Omega}})$  by definition, we have  
 1688  $|V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| \leq kT_{det}^*(Q, D_{\overline{\Omega}})$ .

1689 **Projection.** Assume that  $Q = \pi_{\mathbf{A}}(Q_1)$ . The circuit for  $Q$  has at most  $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + |Q_1|$   
 1690 vertices.

$$1691 \quad |V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| \leq |V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + |Q_1|$$

1693 (From the inductive assumption)

$$1694 \quad \leq kT_{det}^*(Q_1, D_{\overline{\Omega}}) + |Q_1|$$

1696 (By definition of  $T_{det}^*(Q, D_{\overline{\Omega}})$ )

$$1697 \quad \leq kT_{det}^*(Q, D_{\overline{\Omega}}).$$

1699 **Union.** Assume that  $Q = Q_1 \cup Q_2$ . The circuit for  $Q$  has  $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + |V_{Q_2, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + |Q_1 \cap Q_2|$   
 1700 vertices.

$$1701 \quad |V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| \leq |V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + |V_{Q_2, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + |Q_1| + |Q_2|$$

1703 (From the inductive assumption)

$$1704 \quad \leq k(T_{det}^*(Q_1, D_{\overline{\Omega}}) + T_{det}^*(Q_2, D_{\overline{\Omega}})) + (|Q_1| + |Q_2|)$$

1706 (By definition of  $T_{det}^*(Q, D_{\overline{\Omega}})$ )

$$1707 \quad \leq k(T_{det}^*(Q, D_{\overline{\Omega}})).$$

1709  **$m$ -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  
 1710 for  $Q$  has  $|V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + \dots + |V_{Q_k, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + (m-1)|Q_1 \bowtie \dots \bowtie Q_k|$  vertices.

$$1711 \quad |V_{Q, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| = |V_{Q_1, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + \dots + |V_{Q_k, \mathcal{D}_{\mathbb{N}[\mathbf{x}]}}| + (m-1)|Q_1 \bowtie \dots \bowtie Q_k|$$

1713 From the inductive assumption and noting  $\forall i : k_i \leq k$  and  $m \leq k$

$$1714 \quad \leq kT_{det}^*(Q_1, D_{\overline{\Omega}}) + \dots + kT_{det}^*(Q_k, D_{\overline{\Omega}}) +$$

$$1715 \quad (m-1)|Q_1 \bowtie \dots \bowtie Q_m|$$

$$1716 \quad \leq k(T_{det}^*(Q_1, D_{\overline{\Omega}}) + \dots + T_{det}^*(Q_m, D_{\overline{\Omega}}) +$$

$$1717 \quad |Q_1 \bowtie \dots \bowtie Q_m|)$$

1719 (By definition of  $T_{det}^*(Q, D_{\overline{\Omega}})$  and assumption on  $T_{join}(\cdot)$ )

$$1720 \quad \leq kT_{det}^*(Q, D_{\overline{\Omega}}).$$

1722 The property holds for all recursive queries, and the proof holds.  $\blacktriangleleft$

### 1723 E.2.3 Runtime of LC

1724 We next need to show that we can construct the circuit in time linear in the deterministic  
 1725 runtime.

1726  $\blacktriangleright$  **Lemma E.3.** *Given a query  $Q$  over a deterministic bounding database  $D_{\overline{\Omega}}$  and the  $\mathcal{C}^*$*   
 1727 *output by Algorithm 4, the runtime  $T_{LC}(Q, D_{\overline{\Omega}}, \mathcal{C}^*) \leq O(T_{det}^*(Q, D_{\overline{\Omega}}))$ .*

1728 **Proof.** By analysis of Algorithm 4, invoked as  $\mathcal{C}^* \leftarrow LC(Q, D_{\bar{\Omega}}, \{v_0\}, \emptyset, \{(v_0, 0)\})$ .

1729 We assume that the vertex list  $V$ , edge list  $E$ , and vertex label list  $\ell$  are mutable  
1730 accumulators with  $O(1)$  amortized append. We assume that the tuple to sink mapping  $\phi$  is  
1731 a linked hashmap, with  $O(1)$  insertions and retrievals, and  $O(n)$  iteration over the domain of  
1732 keys. We assume that the  $n$ -ary join  $\text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_n)$  can be computed in time  
1733  $T_{\text{join}}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_n))$  (Definition 2.10) and that an intersection  $\text{DOM}(\phi_1) \cap \text{DOM}(\phi_2)$   
1734 can be computed in time  $O(|\text{DOM}(\phi_1)| + |\text{DOM}(\phi_2)|)$  (e.g., with a hash table).

1735 Before proving our runtime bound, we first observe that  $T_{\text{det}}^*(Q, D) \geq \Omega(|Q(D)|)$ . This is  
1736 true by construction for the relation, projection, and union cases, by Definition 2.10 for joins,  
1737 and by the observation that  $|\sigma(R)| \leq |R|$ .

1738 We show that  $T_{\text{det}}^*(Q, D_{\bar{\Omega}})$  is an upper-bound for the runtime of Algorithm 4 by recursion.  
1739 The base case of a relation atom requires only an  $O(|D_{\bar{\Omega}}.R|)$  iteration over the source tuples.  
1740 For the remaining cases, we make the recursive assumption that for every subquery  $Q'$ , it  
1741 holds that  $O(T_{\text{det}}^*(Q', D_{\bar{\Omega}}))$  bounds the runtime of Algorithm 4.

1742 **Selection.** Selection requires a recursive call to Algorithm 4, which by the recursive  
1743 assumption is bounded by  $O(T_{\text{det}}^*(Q', D_{\bar{\Omega}}))$ . Algorithm 4 requires a loop over every element  
1744 of  $Q'(D_{\bar{\Omega}})$ . By the observation above that  $T_{\text{det}}^*(Q, D) \geq \Omega(|Q(D)|)$ , this iteration is also  
1745 bounded by  $O(T_{\text{det}}^*(Q', D_{\bar{\Omega}}))$ .

1746 **Projection.** Projection requires a recursive call to Algorithm 4, which by the recursive  
1747 assumption is bounded by  $O(T_{\text{det}}^*(Q', D_{\bar{\Omega}}))$ , which in turn is a term in  $T_{\text{det}}^*(\pi_{\bar{A}}Q', D_{\bar{\Omega}})$ .  
1748 What remains is an iteration over  $\pi_{\bar{A}}(Q(D_{\bar{\Omega}}))$  (lines 13–16), an iteration over  $Q'(D_{\bar{\Omega}})$  (lines  
1749 17–19), and the construction of a fan-in tree (line 20). The first iteration is  $O(|Q(D_{\bar{\Omega}})|) \leq$   
1750  $O(T_{\text{det}}^*(Q, D_{\bar{\Omega}}))$ . The second iteration and the construction of the bounded fan-in tree are  
1751 both  $O(|Q'(D_{\bar{\Omega}})|) \leq O(T_{\text{det}}^*(Q', D_{\bar{\Omega}})) \leq O(T_{\text{det}}^*(Q, D_{\bar{\Omega}}))$ , by the the observation above that  
1752  $T_{\text{det}}^*(Q, D) \geq \Omega(|Q(D)|)$ .

1753 **Bag Union.** As above, the recursive calls explicitly correspond to terms in the expansion of  
1754  $T_{\text{det}}^*(Q_1 \cup Q_2, D_{\bar{\Omega}})$ . Initializing  $\phi$  (line 24) can be accomplished in  $O(\text{DOM}(\phi_1) + \text{DOM}(\phi_2)) =$   
1755  $O(|Q_1(D_{\bar{\Omega}})| + |Q_2(D_{\bar{\Omega}})|) \leq O(T_{\text{det}}^*(Q_1, D_{\bar{\Omega}}) + T_{\text{det}}^*(Q_2, D_{\bar{\Omega}}))$ . The remainder requires computing  
1756  $Q_1 \cup Q_2$  (line 25) and iterating over it (lines 25–29), which is  $O(|Q_1| + |Q_2|)$  as noted above  
1757 — this directly corresponds to terms in  $T_{\text{det}}^*(Q_1 \cup Q_2, D_{\bar{\Omega}})$ .

1758  **$m$ -ary Join.** As in the prior cases, recursive calls explicitly correspond to terms in our  
1759 target runtime. The remaining logic involves (i) computing  $\text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_m)$ , (ii)  
1760 iterating over the results, and (iii) creating a fan-in tree. Respectively, these are:

- 1761 (i)  $T_{\text{join}}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_m))$
- 1762 (ii)  $O(|Q_1(D_{\bar{\Omega}}) \bowtie \dots \bowtie Q_m(D_{\bar{\Omega}})|) \leq O(T_{\text{join}}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_m)))$  (Definition 2.10)
- 1763 (iii)  $O(m|Q_1(D_{\bar{\Omega}}) \bowtie \dots \bowtie Q_m(D_{\bar{\Omega}})|)$  (as (ii), noting that  $m \leq k = O(1)$ ) ◀

## 1764 F Higher Moments

1765 We make a simple observation to conclude the presentation of our results. So far we have only  
1766 focused on the expectation of  $\Phi$ . In addition, we could e.g. prove bounds of the probability  
1767 of a tuple's multiplicity being at least 1. Progress can be made on this as follows: For any  
1768 positive integer  $m$  we can compute the  $m$ -th moment of the multiplicities, allowing us to e.g.  
1769 use the Chebyshev inequality or other high moment based probability bounds on the events  
1770 we might be interested in. We leave further investigations for future work.

## 1771 **G** The Karp-Luby Estimator

1772 Computing the marginal probability of a tuple in the output of a set-probabilistic database  
 1773 query has been studied extensively. To the best of our knowledge, the current state of  
 1774 the art approximation algorithm for this problem is the Karp-Luby estimator [30], which  
 1775 first appeared in MayBMS/Sprout [38], and more recently as part of an online “anytime”  
 1776 approximation algorithm [20, 15].

The estimator works by observing that for any  $\ell$  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:

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

1777 The Karp-Luby estimator is employed on the SMB representation<sup>14</sup> of  $\mathcal{C}$  (to solve the  
 1778 set-PDB version of Problem 1.6), where each  $W_i$  represents the event that one monomial  
 1779 is true. By simple inspection, if there are  $\ell$  monomials, this estimator has runtime  $\Omega(\ell)$ .  
 1780 Further, a minimum of  $\left\lceil \frac{3 \cdot \ell \cdot \log(\frac{3}{\delta})}{\epsilon^2} \right\rceil$  invocations of the estimator are required to achieve  $1 \pm \epsilon$   
 1781 approximation with probability at least  $1 - \delta$  [38], entailing a runtime at least quadratic in  $\ell$ .  
 1782 As an arbitrary lineage circuit  $\mathcal{C}$  may encode  $\Omega(|\mathcal{C}|^k)$  monomials, the worst case runtime is  
 1783 at least  $\Omega(|\mathcal{C}|^{2k})$  (where  $k$  is the ‘degree’ of lineage polynomial encoded by  $\mathcal{C}$ ). By contrast  
 1784 note that by the discussion after Lemma 4.8 we can solve Problem 1.6 in time  $O(|\mathcal{C}|^2)$  for  
 1785 all BIDB circuits *independent* of the degree  $k$ .

## 1786 **H** Parameterized Complexity

1787 In Sec. 3, we utilized common conjectures from fine-grained complexity theory. The notion of  
 1788  $\#W[1]$  – *hard* is a standard notion in *parameterized complexity*, which by now is a standard  
 1789 complexity tool in providing data complexity bounds on query processing results [22]. E.g.  
 1790 the fact that  $k$ -matching is  $\#W[1]$  – *hard* implies that we cannot have an  $n^{\Omega(1)}$  runtime.  
 1791 However, these results do not carefully track the exponent in the hardness result. E.g.  
 1792  $\#W[1]$  – *hard* for the general  $k$ -matching problem does not imply anything specific for the  
 1793 3-matching problem. Similar questions have led to intense research into the new sub-field  
 1794 of *fine-grained complexity* (see [48]), where we care about the exponent in our hardness  
 1795 assumptions as well— e.g. Conjecture 3.3 is based on the popular *Triangle detection hypothesis*  
 1796 in this area (cf. [34]).

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

## 1797 **I** Response to first cycle reviewer comments

1798 This paper is a resubmission of our submission to the ICDT first cycle. We thank the  
 1799 reviewers for their insightful comments, which we believe has helped improve the presentation  
 1800 of the paper tremendously. We use this section to document the changes that have been made  
 1801 since our prior submission, and in particular, how we have addressed reviewer comments  
 1802 (reviewer comments are shaded followed by our responses).

### 1803 **I.1** Meta Review

1804 Problem definition not stated rigorously nor motivated. Discussion needed on the standard  
 1805 PDB approach vs your approach.

1806 We rewrote Sec. 1 to specifically address this concern. The opening paragraph precisely  
 1807 and formally states the query evaluation problem in bag-PDBs. We use a series of problem  
 1808 statements to clearly define the problem we are addressing as it relates to the query evaluation  
 1809 problem. We made the concrete problem statements more precise by more clearly formalizing  
 1810  $T_{det}^*(Q, D_{\overline{\Omega}})$  and stating our runtime objectives relative to it (??, 1.5, 1.6).

1811 We have included a discussion of the standard approach, e.g. see the paragraph  
 1812 **Relationship to Set-Probabilistic Query Evaluation** on page 4.

1813 Definition 2.6 on reduced BIDD polynomials seem not the right tool for the studied  
 1814 problem.

1815 We have chosen to stick with a less formal, ad-hoc definition (please see Definition 1.3 and  
 1816 ??) of the general problem as suggested by both Reviewer 1 and Reviewer 2. Our earlier  
 1817 proof of the current ?? (in the appendix) had a small bug, which also has been fixed.

1818 The paper is very difficult to read. Improvements are needed in particular for the  
 1819 presentation of the approximation results and their proofs. Also for the notation. Missing  
 1820 definitions for used notions need to be added. Ideally use one instead of three query  
 1821 languages (UCQ, RA+, SPJU).

1822 We have chosen one specific query language throughout the paper ( $\mathcal{RA}^+$ ) and made a  
 1823 concerted effort to use clean, defined, non-ambiguous notation. We have also simplified the  
 1824 notation by limiting the paper's use of provenance semirings (which are needed solely for  
 1825 proofs) to the appendix. To the best of our examination, all notation conflicts have been  
 1826 addressed and definitions for used notions are added (see e.g. Definition C.4 appears before  
 1827 Lemma C.6 and Lemma C.8).

1828 After the rewrite of Sec. 1, we had even less space for Sec. 4. However, we have modified  
 1829 Sec. 4 so that it flows better. In particular, we start off with the algorithm idea first  
 1830 (paragraph **Overview of our Techniques** in Sec. 1 also has more details on the intuition  
 1831 behind the approximation algorithm) and then state the results (with more details on how  
 1832 we argue the claimed runtime). Finally, we clearly state Corollary 4.9 for which queries our  
 1833 linear-time approximation result holds.

### 1829 **I.2** Reviewer 1

1830 1.24 "is #W[1]-hard": parameterized by what?

1831 1.103 and 1.105: again, what is the parameter exactly?

1832 While the above references do not exist in the revised Sec. 1 anymore, all theorem statements

1833 and claims on  $\#W[1]$  runtime have been stated in a way so as to avoid ambiguity in the  
 1834 parameter. Please see e.g. Theorem 3.1 and Theorem 3.6.

1835 You might want to explain your title somewhere (probably in the introduction): in the  
 1836 end, what exactly should be considered harmful and why?

1836 We have modified the title to be more descriptive.

1837 1.45 when discussing Dalvi and Suciu's dichotomy, you might want to mention that they  
 1838 consider \*data complexity\*. Currently the second sentence of your introduction ("take a  
 1839 query  $Q$  and a  $\text{pdb } D$ ") suggests that you are considering combined complexity.

1838 We have made an explicit mention of data complexity when alluding to Dalvi and Suciu's  
 1839 dichotomy. We have further rewritten Sec. 1 in such a way as to explicitly note the type(s)  
 1840 of complexity we are considering (mostly it's parameterized complexity).

1841 1.51 "Consider ... and tuples are independent random event": so this is actually a set  
 1842 PDB... You might want to use an example where the input PDB is actually a bag PDB.  
 1843 The last sentence before the example makes the reader \*expect\* that the example will be  
 1844 of a bag PDB that is not a set PDB

1842 Our revision has removed the example referred to above. While the paper considers inputs to  
 1843 queries that are equivalent to set-PDB, this is not limiting. Please see ?? on ?. Furthermore,  
 1844 we have added a discussion to the appendix that expands on why our results do extend  
 1845 beyond set inputs (Appendix A).

1846 - In the case of set semantics, the lineage of a tuple can be defined for \*any\* query: it is  
 1847 the unique Boolean function that satisfies the if and only if property that you mention  
 1848 on line 70. For bag semantics however, to the best of my knowledge there is no general  
 1849 definition of what is a lineage for an arbitrary query. On line 73, it is not clear at all how  
 1850 the polynomial should be defined, since this will depend on the type of query that you  
 1851 consider

1848 Note that lineage for a set semantics query is as a positive Boolean formula is defined for  
 1849 positive relational algebra. For instance, for aggregate queries a more powerful model ([4])  
 1850 is needed. The definition of the lineage polynomial (bag PDB) semantics over an arbitrary  
 1851  $\mathcal{RA}^+$  query  $Q$  is now given in Fig. 1. We also note that these semantics are not novel  
 1852 (e.g., similar semantics appear for both provenance [25] and probabilistic database [32, 19]  
 1853 contexts). However, as we were unable to find a formal proof of the equivalence between the  
 1854 expectation of the query multiplicity and of the lineage polynomial in related work, we have  
 1855 included a proof of Proposition 2.5.

1856 1.75 "evaluating the lineage of  $t$  over an assignment corresponding to a possible world":  
 1857 here, does the assignment assigns each tuple to true or false? In other words, do the  
 1858 variables  $X$  still represent individual tuples? From what I see later in the article it seems  
 1859 that no, so this is confusing if we compare to what is explained in the previous paragraph  
 1860 about set TIDB

1857 The discussion after Problem 1.2 (in particular, the paragraph TIDBs) specifically address  
 1858 these questions. While values for possible worlds assigned are from  $\{0, 1\}$ , which is analog to  
 1859 Boolean, this is not limiting. Please see ?? (??) and the new appendix section Appendix A.

1860 - 1.135 "polynomial  $Q(X)$ ":  $Q$  should be reserved for queries... You could use  $\varphi$  or  $\phi$  or...  
 1861 anything else but  $Q$  really

1861 We now use  $\Phi(\mathbf{X})$  for (lineage) polynomials.

1862 - 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.

1863 We have rewritten Sec. 1 in a way to stress that we are primarily interested in data  
1864 complexity, but we cannot stop there. As the reviewer has noted, the problem we explore  
1865 requires further analysis, where we require parameterized and fine grained complexity analysis  
1866 to provide a theoretical foundation for the question we ask in ???. We have discussed this in  
1867 the prose following Problem 1.2.

1868 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.

1869 We provide more motivating examples in the first paragraph, and include a more detailed  
1870 discussion of the relationship to sets in paragraph **Relationship to Set-Probabilistic**  
1871 **Query Evaluation** after ???. For example, expected multiplicities can model expectation  
1872 of a `COUNT(*)` query, while in many contexts computing the probability that this count is  
1873 non-zero is not that useful.

1874 As we now explain in the introduction, another motivation for generalizing marginal  
1875 probability to expected multiplicity is that it is a natural generalization. The marginal  
1876 probability of a tuple  $t$  is the expectation of a Boolean random variable that is assigned 1  
1877 in every world where tuple  $t$  exists and 0 otherwise. For bag-PDBs the multiplicity of a  
1878 query result tuple can be modeled as a natural-number random variable that for a world  $D$   
1879 is assigned the multiplicity of the tuple in  $D$ . Thus, a natural generalization of the marginal  
1880 probability (expectation of a Boolean random variable) to bags is the expectation of this  
1881 variable: the tuple's expected multiplicity.

1882 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

1883 To make the paper more accessible and general, we found it better to not use  $N[\mathbf{X}]$ -DBs.  
1884 While we wanted to use UCQ, we found the choice of  $\mathcal{RA}^+$  to be more amenable to the  
1885 presentation of the paper, and have, as suggested stuck with one query formalism.

1886 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

1887 All material in Sec. 2 that is proof-related is in the appendix, while Sec. 2 (modulo the  
1888 proofs) is itself now self-contained.

1889 1.622 and 1.632-634: so a N-PDB is a PDB where each possible world is an N-database, but  
 1890 an  $N[X]$ -PDB is not a PDB where each possible world is an  $N[X]$ -database... Confusing  
 1891 notation

1891 The text now refers to latter as an  $N[X]$ -encoded PDBs.

1892 If you want to be in the setting of bag PDBs, why not consider that the value of the  
 1893 variables are integers rather than Boolean? I.e., consider valuations  $\nu : X \rightarrow \mathbb{N}$  (or even  
 1894 to  $\mathbb{R}$ , why not?) instead of  $X \rightarrow \{0, 1\}$ ; this would seem more natural to me than having  
 1895 this ad-hoc "mix" of Boolean and non-Boolean setting. If you consider this then your  
 1896 "reduced polynomial" trick does not seem to work anymore.

1893 Our objective is to establish the feasibility of bag-probabilistic databases as compared to  
 1894 existing deterministic query processing systems. Accordingly, we take our input model from  
 1895 production database systems like Postgresql, Oracle, DB2, SQLServer, etc. (e.g., see ?? on  
 1896 ??), where duplicate tuples are represented as independent entities. As a convenient benefit,  
 1897 this leads to a direct translation of TIDBs (which are defined over  $\{0, 1\}$  inputs). Finally, as  
 1898 we mention earlier, an easy generalization exists to encode a bag-PDB in a set-PDB (which  
 1899 then allows for bag inputs). Appendix A.

1901 - 1.656 "Thus, from now on we will solely use such vectors...": this seems to be false.  
 1902 Moreover you keep switching notation which makes it very hard to read... Sometimes it is  
 1903  $\varphi$ , sometimes it is small  $w$ , sometimes it is big  $W$  (1.174 or 1.722), sometimes the database  
 1904 is  $\varphi(D)$ , sometimes it is  $\varphi_w(D)$ , other times it is  $D_{[w]}$  (1.671), and so on.

1902 We have made effort to be consistent with the use of notation, following standard usage  
 1903 whenever possible.

1904 1.658 "we use  $\varphi(D)$  to denote the semiring homomorphism  $N[X] \rightarrow \mathbb{N}$  that...": I  
 1905 don't understand why you need a database to extend an assignment to its semiring  
 1906 homomorphism from  $N[X] \rightarrow \mathbb{N}$

1906  $\varphi$  [25] lifts the valuation function (with kind  $N[X] \rightarrow \mathbb{N}$ ) to databases (i.e., a mapping  
 1907 from an  $N[X]$ -DB to a deterministic  $\mathbb{N}$ -DB). We note that the main body of the paper no  
 1908 longer references  $N[X]$ -DBs, and thus  $\varphi$  is discussed exclusively in the appendix.

1909 Figure 2,  $K$  is undefined

1910 We have updated Fig. 1 (originally figure 2) to not need  $K$ .

1911 1.178 " $Q_t$ ", 1.189 "Q will denote a polynomial": this is a very poor choice of notation

1912 1.242 "and query Q": is Q a query or a lineage?

1913 We have reserved  $Q$  to mean an  $\mathcal{RA}^+$  query and nothing else.

1914 Section 2.1.1: here you are considering set semantics no? Otherwise, one would think that  
 1915 for bag semantics the annotation of a tuple could be 0 or something of the form  $c \times X$ ,  
 1916 where  $X$  is a variable and  $c$  is a natural number

1915 Please see Appendix A for a discussion on going beyond set inputs.

1916 Proof of Proposition A.3. I seems the proof should end after 1.687, since you already  
 1917 proved everything from the statement of the proposition. I don't understand what it is  
 1918 that you do after this line.

1918 This text is an informal proof of Proposition 2.5 originally intended to motivate Proposition B.3.

1919 We agree that this should not be part of the proof of the later, and have removed the text.

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

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

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

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

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

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

- 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.

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

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

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

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

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.



1957 Based on this and other reviewer comments, we removed the earlier definition of  $\tilde{\Phi}(\mathbf{X})$   
 1958 and have defined it in a more ad-hoc manner, as suggested by the reviewers, including the  
 1959 comment immediately following.

1960 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?

1961 As alluded to above, we have incorporated the reviewer's suggestion, c.f. Definition 1.3 and  
 1962 ??.

1963 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  $\phi$ , to build an arithmetic circuit  $C$   
 such that for any valuation  $\nu$  of the variables the formula  $\phi$  evaluates to True under  $\nu$  if  
 $C$  evaluates to 1 and the formula  $\phi$  evaluates to False under  $\nu$  if  $C$  evaluates to 0, so this  
 problem is  $\#P$ -hard anyways.

1964 The reviewer is correct. Our hardness results are now stated independently of circuits. We  
 1965 note that the hardness result alluded to at the end of the comment above is not applicable  
 1966 in our case since for fixed queries  $Q$ , ?? and Problem 1.2 can be solved in polynomial time.

1967 Further, as we point out in Sec. 1 what is new in our hardness results is that we show a  
 1968 query  $Q^k$  such that  $T_{det}^*(Q^k, D_{\overline{\Omega}})$  is small (linear in  $|D_{\overline{\Omega}}|$  but solving ?? and Problem 1.2 is  
 1969 hard. We note that it is well-known that one can reduce the problem of counting  $k$ -cliques or  
 1970  $k$ -matchings to a query  $Q$  for which computing  $Q(D_{\overline{\Omega}})$  is  $\#W[1]$ -hard. So our contribution to  
 1971 come up with a different reduction from counting  $k$ -matchings so that the hardness manifests  
 1972 itself in the probabilistic computing part of our problem.

1973 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  $Q_G^k$  for arbitrary  $G$  is  $\#W1$ -hard".

1974 We agree with the reviewer that the result on fixed value of  $p$  is mostly of (narrow)  
 1975 theoretical interest. We have added a discussion summarizing the reviewer's point above  
 1976 below Theorem 3.7.

1977 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).

1978 We have restructured Sec. 1 to more or less follow the reviewer's outline above. The only  
1979 deviation is that we still introduce lineage polynomials. We do this because the polynomial  
1980 view is very helpful in the proofs of our hardness result (in addition to the obvious relevance  
1981 for the approximation algorithm). We have also clarified that our approximation result  
1982 applied to all  $\mathcal{RA}^+$  queries (see Corollary 4.9).

1983 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).

1984 We have implemented the above overview in Sec. 1 when we move from ?? to Problem 1.6.  
1985 For the approximation algorithm we focus on Problem 1.6, which still takes a circuit as an  
1986 input.

1987 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?

1988 We have provided an example in directly after Definition 4.1 as well as a sentence pointing  
1989 out why this definitions is useful.

1990 - 1.409: how does it matter that the circuit  $C$  is the lineage of a UCQ? Doesn't this work for any arithmetic circuit?

1991 The reviewer is correct that the earlier Theorem 4.9 works for any circuit (this result is now  
1992 in the appendix).

1993 1.411: what are  $|C|^2(1, \dots, 1)$  and  $|C|(1, \dots, 1)$ ?

1994 We clarify this overloaded notation immediately after Definition 4.2.

1995 Sometimes you consider UCQs, sometimes  $\mathcal{RA}^+$  queries. I think it would be simpler if you stick to one formalism (probably UCQs is cleaner?)

1996 As alluded to previously, we have followed the reviewer's suggestion and have found  $\mathcal{RA}^+$   
1997 queries to be most amenable for this work.

1998 1.432 what is an FAQ query?

1999 We actually no longer need that result since Lemma 4.8 now has a bound on  $|C|(1, \dots, 1)$  in

2000 terms of  $\text{DEPTH}(\mathcal{C})$  and the latter is used in Corollary 4.9 for all  $\mathcal{RA}^+$  queries. Please see  
 2001 Lemma 4.8 and the followup discussion for more on this.

2002 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.

2003

2004 We have re-organized sec:algo to address the above comments as follows:

- 2005 ■ We now start off Sec. 4.2 with the algorithm idea.
- 2006 ■ We give a quick overview of how the claimed runtime follows from the algorithm idea  
 2007 mentioned above.
- 2008 ■ Added Corollary 4.9 that clearly states that we get an  $O(T_{det}^*(Q, D_{\overline{\Omega}}))$  for *all*  $\mathcal{RA}^+$  queries  
 2009  $Q$ .

2010 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?

2011 There is not a major difference between the two. This observation has persuaded us to  
 2012 eliminate  $\mathbb{N}[\mathbf{X}]$ -DB query evaluation and have only an algorithm for lineage.

2013 We have also re-organized the earlier Section 5 and moved the definition of  $T_{det}^*(\cdot)$  (earlier  
 2014 denoted as  $\mathbf{cost}(\cdot)$ ) to Sec. 2.3 and moved the rest of the material to the appendix.

2015 1.679 where do you use  $\max(D_i)$  later in the proof?

2016

2017 Thank you. This reference was unnecessary and has been removed.

2018 1.688 That sentence is hard to parse, consider reformulating it

2019

2020 As the reviewer notes above, this paragraph is unnecessary and we have removed it.

2021 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

2022

2023 Thank you. The latter definition has been removed.

### 2024 1.3 Reviewer 2

2025 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.

2026 We have addressed the concerns in rewriting the entirety of Sec. 1, explicitly mentioning  
 2027 complexity metrics considered, while forming a series of problem statements that describe  
 2028 the exact problem we are considering, and the complexity metrics considered. We have  
 2029 also adjusted the phrasing of the said theorems and definitions to eliminate the ambiguity  
 2030 described.

2031 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  $Q(X)$ , compute  $E[Q(W)]$ ", 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  $E[Q(W)]$  is the sum of expectations of the monomials in  $Q$  (this is mentioned in Example 1.2).

2032 We have rephrased Definition 2.9 to qualify data complexity. The paper (especially in Sec. 1)  
2033 builds up the fact that we aren't stopping at polynomial time, but exploring parameterized  
2034 complexity and fine grained analysis (as the reviewer aptly noted in the first comment).

2035 Second, Definition 2.6 of Reduced BIDB 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.

2036 The reviewer is correct in their comment on the "mod" part— we apologize for the error. We  
2037 have implemented the reviewer's ad-hoc suggestion in light of Reviewer 1's similar suggestions.

2038 the paper uses three notations (UCQ, RA+, SPJU) for the same thing, and never defines formally any of them.

2039 We have chosen  $\mathcal{RA}^+$  for consistent use throughout the paper. We have included ?? on ??  
2040 for an explicit definition of  $\mathcal{RA}^+$  queries.

2041  $G^\ell$  is used in Lemma 3.8 but defined only in the Appendix (Def. B.2), without even a forward pointer. This is a major omission: Lemma 3.8 is a key step for a key result, but it is impossible to read.

2042 We have fixed this mistake. Unfortunately, because of the changes in the paper (especially  
2043 expanding on Sec. 1), the earlier Lemma 3.8 had to be moved to the appendix.

2044 Definition 2.7. "valid worlds  $\eta$ ". This is confusing. A "possible world" is an element of  $\overline{\Omega}$ : this is not stated explicitly in the paper, but it is implicit on line 163, so I assumed that possible worlds refer to elements of  $\overline{\Omega}$ . If I assumed correctly, then calling  $\eta$  a "world" in Def. 2.7 is misleading, because  $\eta$  is not an element of  $\overline{\Omega}$ . More, it is unclear to me why this definition is needed: it is used right below, in Lemma 2.8, but that lemma seems to continue to hold even if  $w$  is not restricted.

2045  
2046 We agree with the reviewer that this notation is confusing;  $\eta$  is meant to cope with the  
2047 fact that tuples from the same group in a BIDB can not co-exist, even though our  $\{0, 1\}$ -input  
2048 vectors can encode such worlds. We now address this constraint by embedding it directly  
2049 into the reduced polynomial with ??.

2050 line 305: please define what is an "occurrence of  $H$  in  $G$ ". It could mean: a homomorphic image, a subgraph of  $G$  isomorphic to  $H$ , an induced subgraph of  $G$  isomorphic to  $H$ , or maybe something else.

2051 We agree with the reviewer's suggestion and have rephrased the wording to be clear. Please  
2052 see the beginning of Sec. 3.1.

2053 If the proofs are given in the appendix, please say so. Lemmas 3.5 and 3.8 are stated  
without any mention, and one has to guess whether they are obvious, or proven somewhere  
else. On this note: I found Lemma 3.5 quite easy, since the number of  $k$ -matching is  
the coefficient of the leading monomial (of degree  $2k$ ) in  $Q^k(p, p, \dots, p)$ , while Lemma 3.8  
appears much harder. It would help to briefly mention this in the main body of the paper.

2054 We have implemented the reviewer's suggestion. Please see the last sentence of Sec. 1 (as  
2055 well as the expanded discussion on the hardness result in the **Overview of our Techniques**  
2056 part of Sec. 1).

2057 line 177: what is  $\Omega_{\mathbb{N}[\mathbf{X}]}$ ?

2058 We have eliminated the use of  $\mathbb{N}[\mathbf{X}]$ -DBs in the paper proper, using them only when necessary  
2059 in the proofs of the appendix.

2060 line 217. The polynomial  $X^2 + 2XY + Y^2$  is a poor choice to illustrate the degree. There  
are two standard definitions of the degree of a multivariate polynomial, and one has to  
always clarify which one is meant. One definition is the total degree (which is Def. 2.3 in  
the paper), the other is the maximum degree of any single variable. It is nice that you  
are trying to clarify for the reader which definition you are using, but the polynomial  
 $X^2 + 2XY + Y^2$  is worst choice, since here the two coincide.

2061 We have adjusted the example to account for the reviewer's correct observation.

2062 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.

2063 We have removed the redundant terminology the reviewer has pointed out, and refined  
2064 the discussion surrounding (and including) Eq. (1) to be explicit to the novice reader that  
2065 polynomials are by definition of finite degree.

2066 "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.

2067 Our hardness results are now stated independently of circuits so the above statement no  
2068 longer appears in the paper.

2069 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  $Q_G$ ."

2070 Our revision has eliminated this statement.

## 2071 I.4 Reviewer 3

2072 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.

2073 In our revision we expanded on Sec. 1 to give a better overview of the problems we are  
2074 considering in this paper. This meant we had to cut out material in later sections, which  
2075 unfortunately meant we did not have space in Sec. 4 to include any examples that the  
2076 reviewer suggested above. However, we have tried to make Sec. 4 more readable as a whole.

2077 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.

2078

2079 In Problem 1.5 we note explicitly that we care about linear dependence on  $T_{det}^*(Q, D_{\bar{\Omega}})$  and do not care about the exact dependence on  $\epsilon$ . 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.

2082

2083 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?

2084

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

2092

2093 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.

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

2102 We decided to focus on beefing up Sec. 1 and cleaning up definitions of problems, which unfortunately meant we ran out of space to bring back examples from the appendix (especially into Sec. 4).

2104

2105 **I.5 Reviewer 4**

2106 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 ICDD 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.

2107 We have (as noted throughout this section) revised the writing to provide precision and  
2108 clarity to the problem we explore as well as the results we obtain. Part of this revision  
2109 was a complete rewriting of Sec. 1 where we sought to be extremely precise in language  
2110 and through a series of problem statements to help the reader navigate and understand the  
2111 problem we explore as well as how we have gone about exploring that problem coupled with  
2112 the validity of the exploration strategy. We have simultaneously sought to make the paper  
2113 more accessible by assuming the average ICDD attendee and defining or explaining concepts  
2114 that might not be known to them. Finally, we have expanded on the **Overview of our**  
2115 **Techniques** part of Sec. 1 to provide more intuition on how we prove our lower and upper  
2116 bounds.

2117

2118

2119

2120

2121

2122