

# Standard Operating Procedure in Bag PDB Queries Considered Harmful

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) can be reduced to calculating the probability of the *lineage formula* of the result, a Boolean formula over random variables representing the existence of tuples in the database's possible worlds. The analog for bag semantics is a natural number-valued polynomial over random variables that evaluates to the multiplicity of the tuple in each world. In this work, we study the problem of calculating the expectation of such polynomials (a tuple's expected multiplicity) exactly and approximately. For tuple-independent databases (TIDBs), the expected multiplicity of a query result tuple can trivially be computed in linear time in the size of the tuple's lineage, if this polynomial is encoded as a sum of products (the standard operating procedure for Set-PDBs). However, using a reduction from the problem of counting  $k$ -matchings, we demonstrate that calculating the expectation is  $\#W[1]$ -hard when the polynomial is compressed, for example through factorization. Such factorized representations are exploited by modern join algorithms (e.g., worst-case optimal joins), and so our results imply that computing probabilities for Bag-PDB based on the results produced by such algorithms introduces super-linear overhead. The problem stays hard even for polynomials generated by conjunctive queries (CQs) if all input tuples have a fixed probability  $p$  (s.t.  $p \in (0, 1)$ ). We proceed to study polynomials of result tuples of union of conjunctive queries (UCQs) over 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 expectation of polynomial circuits in linear time in the size of the polynomial. By removing Bag-PDB's reliance on the sum-of-products representation of polynomials, this result paves the way for future work on PDBs that are competitive with deterministic databases.

**2012 ACM Subject Classification** Information systems  $\rightarrow$  Incomplete data

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

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

## 1 Introduction

A *probabilistic database*  $\mathcal{D} = (\Omega, \mathbf{P})$  is set of deterministic databases  $\Omega = \{D_1, \dots, D_n\}$  called possible worlds, paired with a probability distribution  $\mathbf{P}$  over these worlds. A well-studied problem in probabilistic databases is to take a query  $Q$  and a probabilistic database  $\mathcal{D}$ , and compute the *marginal probability* of a tuple  $t$  (i.e., its probability of appearing in the result of query  $Q$  over  $\mathcal{D}$ ). This problem is  $\#P$ -hard for set semantics, even for *tuple-independent probabilistic databases* [35] (TIDBs), which are a subclass of probabilistic databases where



© 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).

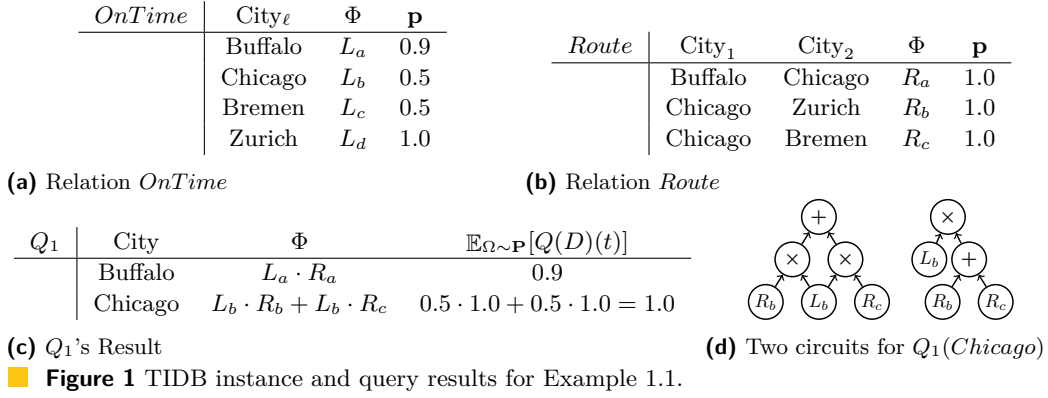
Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:44

Leibniz International Proceedings in Informatics



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

23:2 Bag PDB Queries



45 tuples are independent events. The dichotomy of Dalvi and Suciu [10] separates the hard  
 46 cases, from cases that are in PTIME for unions of conjunctive queries (UCQs). In this work  
 47 we consider bag semantics, where each tuple is associated with a multiplicity  $D_i(t)$  in each  
 48 possible world  $D_i$  and study the analogous problem of computing the expectation of the  
 49 multiplicity of a query result tuple  $t$  (denoted  $Q(D)(t)$ ):

$$50 \quad \mathbb{E}_{\mathbf{D} \sim \mathbf{P}} [Q(\mathbf{D})(t)] = \sum_{D \in \Omega} Q(D)(t) \cdot \mathbf{P}(D) \quad (\text{Expected Result Multiplicity}) \quad (1)$$

51 ► **Example 1.1.** Consider the bag-TIDB relations shown in Fig. 1. We define a TIDB under  
 52 bag semantics analogously to the set case: each input tuple is associated with a probability of  
 53 having a multiplicity of one (and otherwise multiplicity zero), and tuples are independent  
 54 random events. Ignore column  $\Phi$  for now. In this example, we have shipping routes that  
 55 are certain (probability 1.0) and information about whether shipping at locations is on time  
 56 (with a certain probability). Query  $Q_1$ , shown below returns starting points of shipping routes  
 57 where shipment processing is on time.

$$Q_1(City) :- OnTime(City), Route(City, \_)$$

58 Fig. 1c shows the possible results of this query. For example, there is a 90% probability  
 59 there is a single route starting in Buffalo that is on time, and the expected multiplicity of  
 60 this result tuple is 0.9. There are two shipping routes starting in Chicago. Since the Chicago  
 61 location has a 50% probability of being on schedule (we assume that delays are linked), the  
 62 expected multiplicity of this result tuple is  $0.5 + 0.5 = 1.0$ .

63 A well-known result in probabilistic databases is that under set semantics, the marginal  
 64 probability of a query result  $t$  can be computed based on the tuple's lineage. The lineage of  
 65 a tuple is a Boolean formula (an element of the semiring  $\text{PosBool}[\mathbf{X}]$  [19] of positive Boolean  
 66 expressions) over random variables ( $\mathbf{X} = (X_1, \dots, X_n)$ ) that encode the existence of input  
 67 tuples. Each possible world  $D$  corresponds to an assignment  $\{0, 1\}^n$  of the variables in  $\mathbf{X}$   
 68 to either true (the tuple exists in this world) or false (the tuple does not exist in this world).  
 69 Importantly, the following holds: if the lineage formula for  $t$  evaluates to true under the  
 70 assignment for a world  $D$ , then  $t \in Q(D)$ . Thus, the marginal probability of tuple  $t$  is equal  
 71 to the probability that its lineage evaluates to true (with respect to the obvious analog of  
 72 probability distribution  $\mathbf{P}$  defined over  $\mathbf{X}$ ).

73 For bag semantics, the lineage of a tuple is a polynomial over variables  $\mathbf{X} = (X_1, \dots, X_n)$   
 74 with coefficients in the set of natural numbers  $\mathbb{N}$  (an element of semiring  $\mathbb{N}[\mathbf{X}]$ ). Analogously  
 75 to sets, evaluating the lineage for  $t$  over an assignment corresponding to a possible world

76 yields the multiplicity of the result tuple  $t$  in this world. Thus, instead of using Eq. (1)  
 77 to compute the expected result multiplicity of a tuple  $t$ , we can equivalently compute the  
 78 expectation of the lineage polynomial of  $t$ , which for this example we denote as  $\Phi_{Q,\mathcal{D}}^t$  or  
 79  $\Phi$  if the parameters are clear from the context<sup>1</sup>. In this work, we study the complexity of  
 80 computing the expectation of such polynomials encoded as arithmetic circuits.

81 ► **Example 1.2.** *Associating a lineage variable with every input tuple as shown in Fig. 1, we*  
 82 *can compute the lineage of every result tuple as shown in Fig. 1b. For example, the tuple*  
 83 *Chicago is in the result, because  $L_b$  joins with both  $R_b$  and  $R_c$ . Its lineage is  $\Phi = L_b \cdot R_b + L_b \cdot R_c$ .*  
 84 *The expected multiplicity of this result tuple is calculated by summing the multiplicity of the*  
 85 *result tuple, weighted by its probability, over all possible worlds. In this example,  $\Phi$  is a sum of*  
 86 *products (SOP), and so we can use linearity of expectation to solve the problem in linear time*  
 87 *(in the size of  $\Phi$ ). The expectation of the sum is the sum of the expectations of monomials.*  
 88 *The expectation of each monomial is then computed by multiplying the probabilities of the*  
 89 *variables (tuples) occurring in the monomial. The expected multiplicity for Chicago is 1.0.*

90 The expected multiplicity of a query result can be computed in linear time (in the size  
 91 of the result’s lineage) if the lineage is in SOP form. However, this need not be true for  
 92 compressed representations of polynomials, including factorized polynomials or arithmetic  
 93 circuits. For instance, Fig. 1d shows two circuits encoding the lineage of the result tuple  
 94 (*Chicago*) from Example 1.2. The left circuit encodes the lineage as a SOP while the right  
 95 circuit uses distributivity to push the addition gate below the multiplication, resulting in a  
 96 smaller circuit. Given that there is a large body of work (on, e.g., deterministic bag-relational  
 97 query processing) that can output such compressed representations [24, 30], an interesting  
 98 question is whether computing expectations is still in linear time for such compressed  
 99 representations. If the answer is in the affirmative, then probabilities over bag-PDBs can  
 100 be computed with linear overhead (in the size of the compressed representation) using any  
 101 algorithm that computes compressed lineage polynomials. Unfortunately, we prove that this  
 102 is not the case: computing the expected count of a query result tuple is super-linear under  
 103 standard complexity assumptions ( $\#W[1]$ -hard) in the size of a lineage circuit.

104 Concretely, we make the following contributions: (i) We show that the expected result  
 105 multiplicity problem (Definition 2.14) for conjunctive queries for bag-TIDBs is  $\#W[1]$ -hard  
 106 in the size of a lineage circuit by reduction from counting the number of  $k$ -matchings over  
 107 an arbitrary graph; (ii) We present an  $(1 \pm \epsilon)$ -multiplicative approximation algorithm for  
 108 bag-TIDBs and show that for typical database usage patterns (e.g. when the circuit is a  
 109 tree or is generated by recent worst-case optimal join algorithms or their FAQ followups [24])  
 110 its complexity is linear in the size of the compressed lineage encoding; (iii) We generalize the  
 111 approximation algorithm to bag-BIDBs, a more general model of probabilistic data; (iv)  
 112 We further prove that for  $\mathcal{RA}^+$  queries (an equivalently expressive, but factorizable form  
 113 of UCQs), we can approximate the expected output tuple multiplicities with only  $O(\log Z)$   
 114 overhead (where  $Z$  is the number of output tuples) over the runtime of a broad class of query  
 115 processing algorithms. We also observe that our results trivially extend to higher moments  
 116 of the tuple multiplicity (instead of just the expectation).

117 **Overview of our Techniques.** All of our results rely on working with a *reduced* form of the  
 118 lineage polynomial  $\Phi$ . In fact, it turns out that for the TIDB (and BIDB) case, computing  
 119 the expected multiplicity is *exactly* the same as evaluating this reduced polynomial over the

<sup>1</sup> In later sections, where we focus on a single lineage polynomial, we will simply refer to  $\Phi_{Q,\mathcal{D}}^t$  as  $Q$ .

120 probabilities that define the TIDB/BIDB. Next, we motivate this reduced polynomial by  
121 continuing Example 1.1.

122 Consider the query  $Q() :- OnTime(City), Route(City, City'), OnTime(City')$  over the  
123 bag relations of Fig. 1. It can be verified that  $\Phi$  for  $Q$  is  $L_a L_b + L_b L_d + L_b L_c$ . Now consider  
124 the product query  $Q^2() :- Q(), Q()$ . The lineage polynomial for  $Q^2$  is given by  $\Phi^2$ :

$$125 \quad (L_a L_b + L_b L_d + L_b L_c)^2 = L_a^2 L_b^2 + L_b^2 L_d^2 + L_b^2 L_c^2 + 2L_a L_b^2 L_d + 2L_a L_b^2 L_c + 2L_b^2 L_d L_c.$$

126 The expectation  $\mathbb{E}[\Phi^2]$  then is:

$$127 \quad \mathbb{E}[L_a] \mathbb{E}[L_b^2] + \mathbb{E}[L_b^2] \mathbb{E}[L_d^2] + \mathbb{E}[L_b^2] \mathbb{E}[L_c^2] + 2\mathbb{E}[L_a] \mathbb{E}[L_b^2] \mathbb{E}[L_d]$$

$$128 \quad + 2\mathbb{E}[L_a] \mathbb{E}[L_b^2] \mathbb{E}[L_c] + 2\mathbb{E}[L_b^2] \mathbb{E}[L_d] \mathbb{E}[L_c]$$

131 If the domain of a random variable  $W$  is  $\{0, 1\}$ , then for any  $k > 0$ ,  $\mathbb{E}[W^k] = \mathbb{E}[W]$ , which  
132 means that  $\mathbb{E}[\Phi^2]$  simplifies to:

$$133 \quad \mathbb{E}[L_a^2] \mathbb{E}[L_b] + \mathbb{E}[L_b] \mathbb{E}[L_d] + \mathbb{E}[L_b] \mathbb{E}[L_c] + 2\mathbb{E}[L_a] \mathbb{E}[L_b] \mathbb{E}[L_d] + 2\mathbb{E}[L_a] \mathbb{E}[L_b] \mathbb{E}[L_c] + 2\mathbb{E}[L_b] \mathbb{E}[L_d] \mathbb{E}[L_c]$$

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

135 **► Definition 1.3.** For any polynomial  $Q(\mathbf{X})$ , define the reduced polynomial  $\tilde{Q}(\mathbf{X})$  to be the  
136 polynomial obtained by setting all exponents  $e > 1$  in the SOP form of  $Q(\mathbf{X})$  to 1.

137 With  $\Phi^2$  as an example, we have:

$$138 \quad \tilde{\Phi}^2(L_a, L_b, L_c, L_d) = L_a L_b + L_b L_d + L_b L_c + 2L_a L_b L_d + 2L_a L_b L_c + 2L_b L_d L_c$$

140 It can be verified that the reduced polynomial is a closed form of the expected count (i.e.,  
141  $\mathbb{E}[\Phi^2] = \tilde{\Phi}^2(P[L_a = 1], P[L_b = 1], P[L_c = 1], P[L_d = 1])$ ). In fact, we show in Lemma 2.8  
142 that this equivalence holds for *all* UCQs over TIDB/BIDB.

143 To prove our hardness result we show that for the same  $Q$  considered in the running  
144 example, the query  $Q^k$  is able to encode various hard graph-counting problems. We do so by  
145 analyzing how the coefficients in the (univariate) polynomial  $\tilde{\Phi}(p, \dots, p)$  relate to counts of  
146 various sub-graphs on  $k$  edges in an arbitrary graph  $G$  (which is used to define the relations  
147 in  $Q$ ). For the upper bound it is easy to check that if all the probabilities are constant  
148 then  $\tilde{\Phi}(P[X_1 = 1], \dots, P[X_n = 1])$  (i.e. evaluating the original lineage polynomial over  
149 the probability values) is a constant factor approximation. To get an  $(1 \pm \epsilon)$ -multiplicative  
150 approximation we sample monomials from  $\Phi$  and ‘adjust’ their contribution to  $\tilde{\Phi}(\cdot)$ .

151 **Paper Organization.** We present relevant background and notation in Sec. 2. We then  
152 prove our main hardness results in Sec. 3 and present our approximation algorithm in Sec. 4.  
153 We present some (easy) generalizations of our results in Sec. 5 and also discuss extensions  
154 from computing expectations of polynomials to the expected result multiplicity problem  
155 (Definition 2.14). Finally, we discuss related work in Sec. 6 and conclude in Sec. 7.

## 156 **2 Background and Notation**

### 157 **2.1 Probabilistic Databases (PDBs)**

158 An *incomplete database*  $\Omega$  is a set of deterministic databases  $D$  called possible worlds. Denote  
159 the schema of  $D$  as  $sch(D)$ . A *probabilistic database*  $\mathcal{D}$  is a pair  $(\Omega, \mathbf{P})$  where  $\Omega$  is an  
160 incomplete database and  $\mathbf{P}$  is a probability distribution over  $\Omega$ . Queries over probabilistic  
161 databases are evaluated using the so-called possible world semantics. Under the possible

$$\begin{aligned}
\llbracket \pi_A(R) \rrbracket_D(t) &= \sum_{t': \pi_A(t')=t} \llbracket R \rrbracket_D(t') & \llbracket (R_1 \cup R_2) \rrbracket_D(t) &= \llbracket R_1 \rrbracket_D(t) + \llbracket R_2 \rrbracket_D(t) \\
\llbracket \sigma_\theta(R) \rrbracket_D(t) &= \begin{cases} \llbracket R \rrbracket_D(t) & \text{if } \theta(t) \\ 0_{\mathcal{K}} & \text{otherwise.} \end{cases} & \llbracket (R_1 \bowtie R_2) \rrbracket_D(t) &= \llbracket R_1 \rrbracket_D(\pi_{sch(R_1)}(t)) \\
& & & \cdot \llbracket R_2 \rrbracket_D(\pi_{sch(R_2)}(t)) \\
& & \llbracket R \rrbracket_D(t) &= R(t)
\end{aligned}$$

■ **Figure 2** Evaluation semantics  $\llbracket \cdot \rrbracket_D$  for  $\mathbb{N}[\mathbf{X}]$ -DBs [19].

162 world semantics, the result of a query  $Q$  over an incomplete database  $\Omega$  is the set of query  
 163 answers produced by evaluating  $Q$  over each possible world:  $Q(\Omega) = \{ Q(D) \mid D \in \Omega \}$ .

164 For a probabilistic database  $\mathcal{D} = (\Omega, \mathbf{P})$ , the result of a query is the pair  $(Q(\Omega), \mathbf{P}')$  where  
 165  $\mathbf{P}'$  is a probability distribution over  $Q(\Omega)$  that assigns to each possible query result the sum  
 166 of the probabilities of the worlds that produce this answer:

$$167 \quad \forall D \in Q(\Omega) : \mathbf{P}'(D) = \sum_{D' \in \Omega: Q(D')=D} \mathbf{P}(D')$$

168 Let  $\mathbb{N}[\mathbf{X}]$  denote the set of polynomials over variables  $\mathbf{X} = (X_1, \dots, X_n)$  with natural  
 169 number coefficients and exponents. We model incomplete relations using Green et. al.'s  
 170  $\mathbb{N}[\mathbf{X}]$ -databases [19], discussed in detail in Appendix A.1 and summarized here. In an  $\mathbb{N}[\mathbf{X}]$ -  
 171 database, relations are defined as functions from tuples to elements of  $\mathbb{N}[\mathbf{X}]$ , typically called  
 172 annotations. We write  $R(t)$  to denote the polynomial annotating tuple  $t$  in relation  $R$ . Note  
 173 that  $R(t)$  is the lineage polynomial for  $t$ . Each possible world is defined by an assignment of  
 174  $N$  binary values  $\mathbf{W} \in \{0, 1\}^{|\mathbf{X}|}$  to  $\mathbf{X}$ . The multiplicity of  $t \in R$  in this possible world, denoted  
 175  $R(t)(\mathbf{W})$ , is obtained by evaluating the polynomial annotating  $t$  on  $\mathbf{W}$ .  $\mathbb{N}[\mathbf{X}]$ -relations are  
 176 closed under  $\mathcal{RA}^+$  (Fig. 2).

177 We will use  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$ , defined as the tuple  $(\Omega_{\mathbb{N}[\mathbf{X}]}, \mathbf{P})$ , where  $\mathbb{N}[\mathbf{X}]$ -database  $\Omega_{\mathbb{N}[\mathbf{X}]}$  is  
 178 paired with probability distribution  $\mathbf{P}$ . We denote by  $Q_t$  the annotation of tuple  $t$  in the  
 179 result of  $Q$  on an implicit  $\mathbb{N}[\mathbf{X}]$ -PDB (i.e.,  $Q_t = Q(\mathbf{D})(t)$  for some  $\mathbf{D}$ ) and as before, interpret  
 180 it as a function  $Q_t : \{0, 1\}^{|\mathbf{X}|} \rightarrow \mathbb{N}$  from vectors of variable assignments to the corresponding  
 181 value of the annotating polynomial.  $\mathbb{N}[\mathbf{X}]$ -PDBs and a function  $Mod$  (which transforms an  
 182  $\mathbb{N}[\mathbf{X}]$ -PDB to classical, or  $\mathbb{N}$ -PDB [19, 14]) are both formalized in Appendix A.1.

183 ► **Proposition 2.1** (Expectation of polynomials). *Given an  $\mathbb{N}$ -PDB  $\mathcal{D} = (\Omega, \mathbf{P})$  and  $\mathbb{N}[\mathbf{X}]$ -PDB*  
 184  *$\mathbf{D} = (\Omega'_{\mathbb{N}[\mathbf{X}]}, \mathbf{P}')$  where  $Mod(\mathbf{D}) = \mathcal{D}$ , we have:  $\mathbb{E}_{\Omega \sim \mathbf{P}}[Q(\Omega)(t)] = \mathbb{E}_{\mathbf{W} \sim \mathbf{P}'}[Q_t(\mathbf{W})]$ .<sup>2</sup>*

185 A formal proof of Proposition 2.1 is given in Appendix A.3. This proposition shows that  
 186 computing expected tuple multiplicities is equivalent to computing the expectation of a  
 187 polynomial (for that tuple) from a probability distribution over all possible assignments of  
 188 variables in the polynomial to  $\{0, 1\}$ . We focus on this problem from now on, assume an  
 189 implicit result tuple, and so drop the subscript from  $Q_t$  (i.e.,  $Q$  will denote a polynomial).

<sup>2</sup> Although assumed by most prior work on set-probabilistic databases, e.g., as an obvious consequence of [21]'s Theorem 7.1, we are unaware of any formal proof for bag-probabilistic databases.

### 190 2.1.1 TIDBs and BIDBs

191 In this paper, we focus on two popular forms of PDBs: Block-Independent (BIDB) and  
 192 Tuple-Independent (TIDB) PDBs. A BIDB  $\mathbf{D} = (\Omega_{\mathbb{N}[\mathbf{X}]}, \mathbf{P})$  is an  $\mathbb{N}[\mathbf{X}]$ -PDB such that (i)  
 193 every tuple is annotated with either 0 (i.e., the tuple does not exist) or a unique variable  $X_i$   
 194 and (ii) that the tuples  $t$  of  $\mathbf{D}$  for which  $\mathbf{D}(t) \neq 0$  can be partitioned into a set of blocks  
 195 such that variables from separate blocks are independent of each other and variables from  
 196 the same block are disjoint events. In other words, each random variable corresponds to the  
 197 event of a single tuple's presence. A *TIDB* is a BIDB where each block contains exactly  
 198 one tuple. Appendix A.2 explains TIDBs and BIDBs in greater detail. In a BIDB (and by  
 199 extension a TIDB)  $\mathbf{D}$ , tuples are partitioned into  $\ell$  blocks  $b_1, \dots, b_\ell$  where tuple  $t_{i,j} \in b_i$  is  
 200 associated with a probability  $p_{t_{i,j}} = \mathbf{P}[X_{i,j} = 1]$ , and is annotated with a unique variable  
 201  $X_{i,j}$ .<sup>3</sup> Because blocks are independent and tuples from the same block are disjoint, the  
 202 probabilities  $p_{t_{i,j}}$  and the blocks induce the probability distribution  $\mathbf{P}$  of  $\mathbf{D}$ . We will write  
 203 a BIDB-lineage polynomial  $Q(\mathbf{X})$  for a BIDB with  $\ell$  blocks as  $Q(\mathbf{X}) = Q(X_{1,1}, \dots, X_{1,|b_1|},$   
 204  $\dots, X_{\ell,|b_\ell|})$ , where  $|b_i|$  denotes the size of  $b_i$ .<sup>4</sup>

## 205 2.2 Reduced Polynomials and Equivalences

206 We now introduce some terminology and develop a reduced form (a closed form of the  
 207 polynomial's expectation) for polynomials over probability distributions derived from a BIDB  
 208 or TIDB. Note that a polynomial over  $\mathbf{X} = (X_1, \dots, X_n)$  is formally defined as:

$$209 \quad Q(X_1, \dots, X_n) = \sum_{\mathbf{d}=(d_1, \dots, d_n) \in \mathbb{N}^n} c_{\mathbf{d}} \cdot \prod_{i=1}^n X_i^{d_i}. \quad (2)$$

210 **► Definition 2.2** (Standard Monomial Basis). *From above, the term  $\prod_{i=1}^n X_i^{d_i}$  is a monomial.*  
 211 *A polynomial  $Q(\mathbf{X})$  is in standard monomial basis (SMB) when we keep only the terms with*  
 212  *$c_{\mathbf{d}} \neq 0$  from Eq. (2).*

213 We consider SMB as the default representation of a polynomial. We use  $\text{SMB}(Q)$  to denote  
 214 the SMB form of a polynomial  $Q$ .

215 **► Definition 2.3** (Degree). *The degree of polynomial  $Q(\mathbf{X})$  is the largest  $\sum_{i=1}^n d_i$  such that*  
 216  *$c_{(d_1, \dots, d_n)} \neq 0$ .*

217 The degree of the polynomial  $X^2 + 2XY + Y^2$  is 2. Product terms in lineage arise only  
 218 from join operations (Fig. 2), so intuitively, the degree of a lineage polynomial is analogous  
 219 to the largest number of joins in any clause of the UCQ query that created it. In this paper  
 220 we consider only finite degree polynomials. We call a polynomial  $Q(\mathbf{X})$  a *BIDB-lineage*  
 221 *polynomial* (resp., *TIDB-lineage polynomial*, or simply lineage polynomial), if there exists a  
 222  $\mathcal{RA}^+$  query  $Q$ , BIDB  $\mathbf{D}$  (TIDB  $\mathbf{D}$ , or  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$ ), and tuple  $t$  such that  $Q(\mathbf{X}) = Q(\mathbf{D})(t)$ .

223 **► Definition 2.4** (Modding with a set). *Let  $S$  be a set of polynomials over  $\mathbf{X}$ . Then  $Q(\mathbf{X})$*   
 224 *mod  $S$  is the polynomial obtained by taking the mod of  $Q(\mathbf{X})$  over all polynomials in  $S$  (order*  
 225 *does not matter).*

<sup>3</sup> Although only a single independent,  $[|b_i| + 1]$ -valued variable is customarily used per block, we decompose it into  $|b_i|$  correlated  $\{0, 1\}$ -valued variables per block that can be used directly in polynomials (without an indicator function). For  $t_j \in b_i$ , the event  $(X_{i,j} = 1)$  corresponds to the event  $(X_i = j)$  in the customary annotation scheme.

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

226 For example for a set of polynomials  $S = \{X^2 - X, Y^2 - Y\}$ , taking the polynomial  $2X^2 +$   
 227  $3XY - 2Y^2 \pmod S$  yields  $2X + 3XY - 2Y$ .

228 ► **Definition 2.5** ( $\mathcal{B}, \mathcal{T}$ ). Given the set of BIDB variables  $\{X_{i,j}\}$ , define

$$229 \quad \mathcal{B} = \{ X_{i,j} \cdot X_{i,j'} \mid i \in [\ell], j, j' \in [ |b_i| ] \} \quad \mathcal{T} = \{ X_{i,j}^2 - X_{i,j} \mid i \in [\ell], j \in [ |b_i| ] \}$$

230 ► **Definition 2.6** (Reduced BIDB Polynomials). Let  $Q(\mathbf{X})$  be a BIDB-lineage polynomial. The  
 231 reduced form  $\tilde{Q}(\mathbf{X})$  of  $Q(\mathbf{X})$  is:  $\tilde{Q}(\mathbf{X}) = Q(\mathbf{X}) \pmod{(\mathcal{T} \cup \mathcal{B})}$

232 All exponents  $e > 1$  in  $\text{SMB}(Q(\mathbf{X}))$  are reduced to  $e = 1$  via  $\pmod{\mathcal{T}}$ . Performing the  
 233 modulus of  $\tilde{Q}(\mathbf{X})$  with  $\mathcal{B}$  ensures the disjoint condition of BIDB, removing monomials with  
 234 lineage variables from the same block. For the special case of TIDBs, the second step is not  
 235 necessary since every block contains a single tuple.

236 ► **Definition 2.7** (Valid Worlds). For probability distribution  $\mathbf{P}$ , the set of valid worlds  $\eta$   
 237 consists of all the worlds with probability value greater than 0; i.e., for variable vector  $\mathbf{W}$

$$238 \quad \eta = \{ \mathbf{w} \mid P[\mathbf{W} = \mathbf{w}] > 0 \}$$

239 Next, we show why the reduced form is useful for our purposes:

240 ► **Lemma 2.8.** Let  $\mathbf{D}$  be a BIDB over variables  $\mathbf{X} = \{X_1, \dots, X_n\}$  and with probability  
 241 distribution  $\mathbf{P}$  produced by the tuple probability vector  $\mathbf{p} = (p_1, \dots, p_n)$  over all  $\mathbf{w}$  in  $\eta$ . For  
 242 any BIDB-lineage polynomial  $Q(\mathbf{X})$  based on  $\mathbf{D}$  and query  $Q$  we have:

$$243 \quad \mathbb{E}_{\mathbf{W} \sim \mathbf{P}} [Q(\mathbf{W})] = \tilde{Q}(\mathbf{p}).$$

244 Note that in the preceding lemma, we have assigned  $\mathbf{p}$  to the variables  $\mathbf{X}$ . Intuitively,  
 245 Lemma 2.8 states that when we replace each variable  $X_i$  with its probability  $p_i$  in the reduced  
 246 form of a BIDB-lineage polynomial and evaluate the resulting expression in  $\mathbb{R}$ , then the  
 247 result is the expectation of the polynomial.

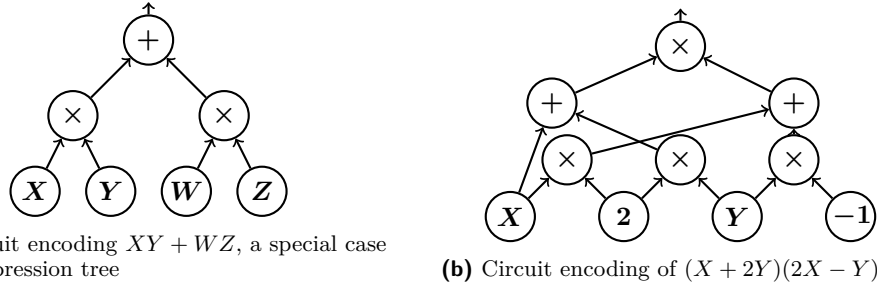
248 ► **Corollary 2.9.** If  $Q$  is a BIDB-lineage polynomial, then the expectation of  $Q$ , i.e.,  $\mathbb{E}[Q] =$   
 249  $\tilde{Q}(p_1, \dots, p_n)$  can be computed in  $O(\text{SIZE}(\text{SMB}(Q)))$ , where  $\text{SIZE}(Q)$  (Definition 4.4) is  
 250 proportional to the total number of multiplication/addition operators in  $Q$ .

## 251 2.3 Problem Definition

252 We first formally define circuits, an encoding of polynomials that we use throughout the paper.  
 253 Since we are particularly using *lineage* circuits, we drop the term lineage and only refer to  
 254 them as circuits. For illustrative purposes consider the polynomial  $Q(\mathbf{X}) = 2X^2 + 3XY - 2Y^2$   
 255 over  $\mathbf{X} = [X, Y]$ .

256 We represent query polynomials via *arithmetic circuits* [6], a standard way to represent  
 257 polynomials over fields (particularly in the field of algebraic complexity) that we use for  
 258 polynomials over  $\mathbb{N}$  in the obvious way.

259 ► **Definition 2.10** (Circuit). A circuit  $\mathcal{C}$  is a Directed Acyclic Graph (DAG) whose source  
 260 nodes (in degree of 0) consist of elements in either  $\mathbb{R}$  or  $\mathbf{X}$ . The internal nodes and (the  
 261 single) sink node of  $\mathcal{C}$  (corresponding to the result tuple  $t$ ) have binary input and are either  
 262 sum (+) or product ( $\times$ ) gates. Each node in a circuit  $\mathcal{C}$  has the following members: **type**,  
 263 **val**, **partial**, **input**, **degree** and **Lweight**, **Rweight**, where **type** is the type of value stored  
 264 in the node (one of  $\{+, \times, \text{VAR}, \text{NUM}\}$ ), **val** is the value stored (a constant or variable), and  
 265 **input** is the list of the nodes inputs. We use  $C_L$  to denote the left input and  $C_R$  the right input



(a) Circuit encoding  $XY + WZ$ , a special case of an expression tree

(b) Circuit encoding of  $(X + 2Y)(2X - Y)$

■ **Figure 3** Example circuit encodings

266 or the sink of circuit  $\mathcal{C}$ . When the underlying DAG is a tree (with edges pointing towards the  
 267 root), we will refer to the structure as an expression tree  $T$ . Note that in such a case, the root  
 268 of  $T$  is analogous to the sink of  $\mathcal{C}$ .

269 As stated in Definition 2.10, every internal node has at most two in-edges, is labeled as  
 270 an addition or a multiplication node, and has no limit on its outdegree. Note that if we limit  
 271 the outdegree to one, then we get expression trees. We ignore the fields `partial`, `Lweight`,  
 272 and `Rweight` until Sec. 4.

273 ► **Example 2.11.** The circuit  $\mathcal{C}$  in Fig. 3a encodes the polynomial  $XY + WZ$ . Note that  
 274 circuit  $\mathcal{C}$  encodes a tree, with edges pointing towards the root.

275 The semantics of circuits follows the obvious interpretation. We next define its relationship  
 276 with polynomials formally:

277 ► **Definition 2.12** ( $\text{POLY}(\cdot)$ ). Denote  $\text{POLY}(\mathcal{C})$  to be the function from circuit  $\mathcal{C}$  to its  
 278 corresponding polynomial.  $\text{POLY}(\cdot)$  is recursively defined on  $\mathcal{C}$  as follows, with addition and  
 279 multiplication following the standard interpretation for polynomials:

$$280 \quad \text{POLY}(\mathcal{C}) = \begin{cases} \text{POLY}(\mathcal{C}_L) + \text{POLY}(\mathcal{C}_R) & \text{if } \mathcal{C}.\text{type} = + \\ \text{POLY}(\mathcal{C}_L) \cdot \text{POLY}(\mathcal{C}_R) & \text{if } \mathcal{C}.\text{type} = \times \\ \mathcal{C}.\text{val} & \text{if } \mathcal{C}.\text{type} = \text{VAR OR NUM}. \end{cases}$$

281 Note that  $\mathcal{C}$  need not encode an expression in SMB. For instance,  $\mathcal{C}$  could represent a  
 282 compressed form of the running example, such as  $(X + 2Y)(2X - Y)$ , as shown in Fig. 3b,  
 283 while  $\text{POLY}(\mathcal{C}) = 2X^2 + 3XY - 2Y^2$ .

284 ► **Definition 2.13** (Circuit Set).  $\text{CSet}(Q(\mathbf{X}))$  is the set of all possible circuits  $\mathcal{C}$  such that  
 285  $\text{POLY}(\mathcal{C}) = Q(\mathbf{X})$ .

286 The circuit of Fig. 3b is an element of  $\text{CSet}(2X^2 + 3XY - 2Y^2)$ . One can think of  
 287  $\text{CSet}(Q(\mathbf{X}))$  as the infinite set of circuits each of which model an encoding (factorization)  
 288 equal to  $\text{POLY}(\mathcal{C})$ . Note that Definition 2.13 implies that  $\mathcal{C} \in \text{CSet}(\text{POLY}(\mathcal{C}))$ .

289 We are now ready to formally state our **main problem**.

290 ► **Definition 2.14** (The Expected Result Multiplicity Problem). Let  $\mathbf{X} = (X_1, \dots, X_n)$ , and  
 291  $\mathcal{D}$  be an  $\mathbb{N}[\mathbf{X}]$ -PDB over  $\mathbf{X}$  with probability distribution  $\mathbf{P}$  over assignments  $\mathbf{X} \rightarrow \{0, 1\}$ ,  $Q$   
 292 an  $n$ -ary query, and  $t$  an  $n$ -ary tuple. The **EXPECTED RESULT MULTIPLICITY PROBLEM** is  
 293 defined as follows:

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



### 3 Hardness of exact computation

In this section, we will prove that computing  $\mathbb{E}_{\mathbf{W} \sim \mathbf{P}}[Q(\mathbf{W})]$  exactly for a TIDB-lineage polynomial  $Q(\mathbf{X})$  generated from a project-join query (even an expression tree representation) is  $\#W[1]$ -hard. Note that this implies hardness for BIDs and general  $\mathbb{N}[\mathbf{X}]$ -PDBs under bag semantics. Furthermore, we demonstrate in Sec. 3.3 that the problem remains hard, even if  $P[X_i = 1] = p$  for all  $X_i$  and any fixed valued  $p \in (0, 1)$  as long as certain popular hardness conjectures in fine-grained complexity hold.

#### 3.1 Preliminaries

Our hardness results are based on (exactly) counting the number of occurrences of a subgraph  $H$  in  $G$ . Let  $\#(G, H)$  denote the number of occurrences of  $H$  in graph  $G$ . We can think of  $H$  as being of constant size and  $G$  as growing. In particular, we will consider the problems of computing the following counts (given  $G$  as an input and its adjacency list representation):  $\#(G, \triangle)$  (the number of triangles),  $\#(G, \text{3-matchings})$  (the number of 3-matchings), and the latter's generalization  $\#(G, \text{3} \cdots \text{3}^k)$  (the number of  $k$ -matchings). Our hardness result in Sec. 3.2 is based on the following result:

► **Theorem 3.1** ([8]). *Given positive integer  $k$  and undirected graph  $G$  with no self-loops or parallel edges, computing  $\#(G, \text{3} \cdots \text{3}^k)$  exactly is  $\#W[1]$ -hard (parameterization is in  $k$ ).*

The above result means that we cannot hope to count the number of  $k$ -matchings in  $G = (V, E)$  in time  $f(k) \cdot |V|^c$  for any function  $f$  and constant  $c$  independent of  $k$ . In fact, all known algorithms to solve this problem take time  $|V|^{\Omega(k)}$ . Our hardness result in Section 3.3 is based on the following conjectured hardness result:

► **Conjecture 3.2.** *There exists a constant  $\epsilon_0 > 0$  such that given an undirected graph  $G = (V, E)$ , computing exactly  $\#(G, \triangle)$  cannot be done in time  $o(|E|^{1+\epsilon_0})$ .*

Based on the so called *Triangle detection hypothesis* (cf. [25]), which states that detection of whether  $G$  has a triangle or not takes time  $\Omega(|E|^{4/3})$ , implies that in Conjecture 3.2 we can take  $\epsilon_0 \geq \frac{1}{3}$ .

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

► **Definition 3.3.** *For any graph  $G = ([n], E)$  and  $k \geq 1$ , define*

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

Our hardness results only need a TIDB instance; We also consider the special case when all the tuple probabilities (probabilities assigned to  $X_i$  by  $\mathbf{p}$ ) are the same value. Note that our hardness results even hold for the expression trees.

Returning to Fig. 1, it is easy to see that  $Q_G^k(\mathbf{X})$  generalizes our running example query:

$$Q_G^k := \text{OnTime}(C_1), \text{Route}(C_1, C'_1), \text{OnTime}(C'_1), \dots, \text{OnTime}(C_k), \text{Route}(C_k, C'_k), \text{OnTime}(C'_k)$$

334 where adapting the PDB instance in Fig. 1, relation *OnTime* has  $n$  tuples corresponding  
 335 to each vertex in  $V = [n]$  each with probability  $p$  and *Route*(City<sub>1</sub>, City<sub>2</sub>) has tuples  
 336 corresponding to the edges  $E$  (each with probability of 1).<sup>5</sup> Note that this implies that our  
 337 hard query polynomial can be represented as an expression tree produced by a project-join  
 338 query with same probability value for each input tuple  $p_i$ .

### 339 3.2 Multiple Distinct $p$ Values

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

341 ► **Theorem 3.4.** *Computing  $\tilde{Q}_G^k(p_i, \dots, p_i)$  for arbitrary  $G$  and any  $(2k + 1)$  distinct values*  
 342  *$p_i$  ( $0 \leq i \leq 2k$ ) is  $\#W[1]$ -hard (parameterization is in  $k$ ).*

343 We will prove the above result by reducing from the problem of computing the number of  
 344  $k$ -matchings in  $G$ . Given the current best-known algorithm for this counting problem, our  
 345 results imply that unless the state-of-the-art  $k$ -matching algorithms are improved, we cannot  
 346 hope to solve our problem in time better than  $\Omega_k(m^{k/2})$  where  $m = |E|$ , which is only  
 347 quadratically faster than expanding  $Q_G^k(\mathbf{X})$  into its SMB form and then using Corollary 2.9.  
 348 By contrast the approximation algorithm we present in Sec. 4 has runtime  $O_k(m)$  for this  
 349 query.

350 The following lemma reduces the problem of counting  $k$ -matchings in a graph to our problem  
 351 (and proves Theorem 3.4):

352 ► **Lemma 3.5.** *Let  $p_0, \dots, p_{2k}$  be distinct values in  $(0, 1]$ . Then given the values  $\tilde{Q}_G^k(p_i, \dots, p_i)$*   
 353 *for  $0 \leq i \leq 2k$ , the number of  $k$ -matchings in  $G$  can be computed in  $O(k^3)$  time.*

### 354 3.3 Single $p$ value

355 While Theorem 3.4 shows that computing  $\tilde{Q}(p, \dots, p)$  in general is hard it does not rule out  
 356 the possibility that one can compute this value exactly for a *fixed* value of  $p$ . Indeed, it is  
 357 easy to check that one can compute  $\tilde{Q}(p, \dots, p)$  exactly in linear time for  $p \in \{0, 1\}$ . In this  
 358 section, we show that these two are the only possibilities:

359 ► **Theorem 3.6.** *Fix  $p \in (0, 1)$ . Then assuming Conjecture 3.2 is true, any algorithm that*  
 360 *computes  $\tilde{Q}_G^3(p, \dots, p)$  from  $G$  exactly has to run in time  $\Omega(m^{1+\epsilon_0})$ , where  $\epsilon_0$  is as defined*  
 361 *in Conjecture 3.2.*

362 The above shows the hardness for a very specific query polynomial but it is easy to come  
 363 up with an infinite family of hard query polynomials by ‘embedding’  $\tilde{Q}_G^3$  into an infinite  
 364 family of trivial query polynomials. Unlike Theorem 3.4 the above result does not show that  
 365 computing  $\tilde{Q}_G^3(p, \dots, p)$  for a fixed  $p \in (0, 1)$  is  $\#W[1]$ -hard. However, in Sec. 4 we show  
 366 that if we are willing to compute an approximation that this problem (and indeed solving  
 367 our problem for a much more general setting) is in linear time.

368 We will prove the above result by the following reduction:

369 ► **Theorem 3.7.** *Fix  $p \in (0, 1)$ . Let  $G$  be a graph on  $m$  edges. If we can compute  $\tilde{Q}_G^3(p, \dots, p)$*   
 370 *exactly in  $T(m)$  time, then we can exactly compute  $\#(G, \mathfrak{A})$  in  $O(T(m) + m)$  time.*

---

<sup>5</sup> Technically,  $Q_G^k(\mathbf{X})$  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.

371 The following result immediately implies Theorem 3.7:

372 ▶ **Lemma 3.8.** *Fix  $p \in (0, 1)$ . Given  $\tilde{Q}_{G^{(\ell)}}^3(p, \dots, p)$  for  $\ell \in [2]$ , we can compute in  $O(m)$*   
 373 *time a vector  $\mathbf{b} \in \mathbb{R}^3$  such that*

$$374 \begin{pmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{pmatrix} \cdot \begin{pmatrix} \#(G, \mathfrak{A}) \\ \#(G, \mathfrak{B}) \end{pmatrix} = \mathbf{b},$$

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

## 376 4 $1 \pm \epsilon$ Approximation Algorithm

377 In Sec. 3, we showed that computing the expected multiplicity of a compressed lineage  
 378 polynomial for TIDB (even just based on project-join queries), and by extension BIDB (or  
 379 any  $\mathbb{N}[\mathbf{X}]$ -PDB) is unlikely to be possible in linear time (Theorem 3.4), even if all tuples have  
 380 the same probability (Theorem 3.6). Given this, we now design an approximation algorithm  
 381 for our problem that runs in *linear time*.<sup>6</sup> The following approximation algorithm applies  
 382 to BIDB, though our bounds are more meaningful for a non-trivial subclass of BIDBs that  
 383 contains both TIDBs, as well as the PDBench benchmark [1].

### 384 4.1 Preliminaries and some more notation

385 We now introduce useful definitions and notation related to circuits and polynomials. All  
 386 proofs and missing pseudocode can be found in Appendix C.

387 ▶ **Definition 4.1** (Variables in a monomial). *Given a monomial  $v$ , we use  $\text{VAR}(v)$  to denote*  
 388 *the set of variables in  $v$ .*

389 For example the monomial  $XY$  has  $\text{VAR}(XY) = \{X, Y\}$ .

390 ▶ **Definition 4.2** ( $E(\mathcal{C})$ ). *The logical view of  $E(\mathcal{C})$  is a list of tuples  $(\mathbf{v}, \mathbf{c})$ , where  $\mathbf{v}$  is a set of*  
 391 *variables and  $\mathbf{c}$  is in  $\mathbb{R}$ .  $E(\mathcal{C})$  has the following recursive definition ( $\circ$  is list concatenation).*

$$392 E(\mathcal{C}) = \begin{cases} E(\mathcal{C}_L) \circ E(\mathcal{C}_R) & \text{if } \mathcal{C}.type = + \\ \{(\mathbf{v}_L \cup \mathbf{v}_R, \mathbf{c}_L \cdot \mathbf{c}_R) \mid (\mathbf{v}_L, \mathbf{c}_L) \in E(\mathcal{C}_L), (\mathbf{v}_R, \mathbf{c}_R) \in E(\mathcal{C}_R)\} & \text{if } \mathcal{C}.type = \times \\ \text{List}[(\emptyset, \mathcal{C}.val)] & \text{if } \mathcal{C}.type = \text{NUM} \\ \text{List}[(\{\mathcal{C}.val\}, 1)] & \text{if } \mathcal{C}.type = \text{VAR}. \end{cases}$$

393 For further explanation, please refer to Example C.2.

394 ▶ **Definition 4.3** ( $|\mathcal{C}|(\mathbf{X})$ ). *For any circuit  $\mathcal{C}$ , the corresponding positive circuit, denoted  $|\mathcal{C}|$ ,*  
 395 *is obtained from  $\mathcal{C}$  as follows. For each leaf node  $\ell$  of  $\mathcal{C}$  where  $\ell.type$  is NUM, update  $\ell.value$*   
 396 *to  $|\ell.value|$ .*

397 Please see Example C.3 for an illustration.

398 ▶ **Definition 4.4** ( $\text{SIZE}(\cdot)$ ). *The function  $\text{SIZE}$  takes a circuit  $\mathcal{C}$  as input and outputs the*  
 399 *number of gates (nodes) in  $\mathcal{C}$ .*

400 ▶ **Definition 4.5** ( $\text{DEPTH}(\cdot)$ ). *The function  $\text{DEPTH}$  has circuit  $\mathcal{C}$  as input and outputs the*  
 401 *number of levels in  $\mathcal{C}$ .*

<sup>6</sup> For a very broad class of circuits: please see the discussion after Lemma 4.11 for more.

402 ► **Definition 4.6** ( $\text{DEG}(\cdot)$ ).<sup>7</sup>  $\text{DEG}(\mathcal{C})$  is defined recursively as follows:

$$403 \quad \text{DEG}(\mathcal{C}) = \begin{cases} \max(\text{DEG}(\mathcal{C}_L), \text{DEG}(\mathcal{C}_R)) & \text{if } \mathcal{C}.\text{type} = + \\ \text{DEG}(\mathcal{C}_L) + \text{DEG}(\mathcal{C}_R) + 1 & \text{if } \mathcal{C}.\text{type} = \times \\ 0 & \text{otherwise.} \end{cases}$$

404 Finally, we will need the following notation for the complexity of multiplying large integers:

405 ► **Definition 4.7** ( $\overline{\mathcal{M}}(\cdot, \cdot)$ ).<sup>8</sup> In a RAM model of word size of  $W$ -bits,  $\overline{\mathcal{M}}(M, W)$  denotes  
406 the complexity of multiplying two integers represented with  $M$ -bits. (We will assume that for  
407 input of size  $N$ ,  $W = O(\log N)$ ).

## 408 4.2 Our main result

409 ► **Theorem 4.8.** Let  $\mathcal{C}$  be a circuit for a UCQ over BIDB and define  $Q(\mathbf{X}) = \text{POLY}(\mathcal{C})$  and  
410 let  $k = \text{DEG}(\mathcal{C})$ . Then an estimate  $\mathcal{E}$  of  $\tilde{Q}(p_1, \dots, p_n)$  can be computed in time

$$411 \quad O\left(\left(\text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta} \cdot |\mathcal{C}|^2(1, \dots, 1) \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C})}{(\epsilon')^2 \cdot \tilde{Q}^2(p_1, \dots, p_n)}\right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right)$$

412 such that

$$413 \quad P\left(\left|\mathcal{E} - \tilde{Q}(p_1, \dots, p_n)\right| > \epsilon' \cdot \tilde{Q}(p_1, \dots, p_n)\right) \leq \delta. \quad (3)$$

414 To get linear runtime results from Theorem 4.8, we will need to define another parameter  
415 modeling the (weighted) number of monomials in  $\mathbf{E}(\mathcal{C})$  to be ‘canceled’ when it is modded  
416 with  $\mathcal{B}$  (Definition 2.5).

417 ► **Definition 4.9** (Parameter  $\gamma$ ). Given an expression tree  $\mathcal{C}$ , define

$$418 \quad \gamma(\mathcal{C}) = \frac{\sum_{(v, c) \in \mathbf{E}(\mathcal{C})} |c| \cdot \mathbb{1}(v \bmod \mathcal{B} \equiv 0)}{|\mathcal{C}|(1, \dots, 1)}$$

419 We next present a few corollaries of Theorem 4.8.

420 ► **Corollary 4.10.** Let  $Q(\mathbf{X})$  be as in Theorem 4.8 and let  $\gamma = \gamma(\mathcal{C})$ . Further let it be the  
421 case that  $p_i \geq p_0$  for all  $i \in [n]$ . Then an estimate  $\mathcal{E}$  of  $\tilde{Q}(p_1, \dots, p_n)$  satisfying Eq. (3) can  
422 be computed in time

$$423 \quad O\left(\left(\text{SIZE}(\mathcal{C}) + \frac{\log \frac{1}{\delta} \cdot k \cdot \log k \cdot \text{DEPTH}(\mathcal{C})}{(\epsilon')^2 \cdot (1 - \gamma)^2 \cdot p_0^{2k}}\right) \cdot \overline{\mathcal{M}}(\log(|\mathcal{C}|(1, \dots, 1)), \log(\text{SIZE}(\mathcal{C})))\right)$$

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

426 The restriction on  $\gamma$  is satisfied by any TIDB (where  $\gamma = 0$ ) as well as for all three queries  
427 of the PDBench BIDB benchmark (see Appendix C.11 for experimental results).

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

<sup>7</sup> Note that the degree of  $\text{POLY}(|\mathcal{C}|)$  is always upper bounded by  $\text{deg}(\mathcal{C})$  and the latter can be strictly larger (e.g. consider the case when  $\mathcal{C}$  multiplies two copies of the constant 1— here we have  $\text{deg}(\mathcal{C}) = 1$  but degree of  $\text{POLY}(|\mathcal{C}|)$  is 0).

<sup>8</sup> We note that when doing arithmetic operations on the RAM model for input of size  $N$ , we have that  $\overline{\mathcal{M}}(O(\log N), O(\log N)) = O(1)$ . More generally we have  $\overline{\mathcal{M}}(N, O(\log N)) = O(N \log N \log \log N)$ .

---

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


---

**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:  $N \leftarrow \left\lceil \frac{2 \log \frac{2}{\epsilon}}{\epsilon^2} \right\rceil$ 
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:    $(M, \text{sgn}_i) \leftarrow \text{SAMPLEMONOMIAL}(\mathcal{C}_{\text{mod}})$  ▷ SAMPLEMONOMIAL is Algorithm 3. Note
   that  $\text{sgn}_i$  is the sign of the monomial's coefficient and not the coefficient itself
6:   if  $M$  has at most one variable from each block then
7:      $Y_i \leftarrow \prod_{X_j \in \text{VAR}(M)} p_j$ 
8:      $Y_i \leftarrow Y_i \times \text{sgn}_i$ 
9:      $\text{acc} \leftarrow \text{acc} + Y_i$  ▷ 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  $\text{acc}$ 

```

---

429 ▶ **Lemma 4.11.** *For any circuit  $\mathcal{C}$  with  $\text{DEG}(\mathcal{C}) = k$ , we have  $|\mathcal{C}|(1, \dots, 1) \leq 2^{2^k \cdot \text{SIZE}(\mathcal{C})}$ .*  
 430 *Further, under either of the following conditions:*

- 431 1.  $\mathcal{C}$  is a tree,  
 432 2.  $\mathcal{C}$  encodes the run of the algorithm in [24] on an FAQ query,

433 *we have  $|\mathcal{C}|(1, \dots, 1) \leq \text{SIZE}(\mathcal{C})^{O(k)}$ .*

434 Note that the above implies that with the assumption  $p_0 > 0$  and  $\gamma < 1$  are absolute  
 435 constants from Corollary 4.10, then the runtime there simplifies to  $O_k \left( \frac{1}{\epsilon^2} \cdot \text{SIZE}(\mathcal{C})^2 \cdot \log \frac{1}{\delta} \right)$   
 436 for general circuits  $\mathcal{C}$  and to  $O_k \left( \frac{1}{\epsilon^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta} \right)$  for the case when  $\mathcal{C}$  satisfies the specific  
 437 conditions in Lemma 4.11. In Appendix C.4 we argue that these conditions are very general  
 438 and encompass many interesting scenarios, including query evaluation under  $\mathcal{RA}^+$  or FAQ.

### 439 4.3 Approximating $\tilde{Q}$

440 The algorithm (APPROXIMATE $\tilde{Q}$  detailed in Algorithm 1) to prove Theorem 4.8 follows from  
 441 the following observation. Given a query polynomial  $Q(\mathbf{X}) = \text{POLY}(\mathcal{C})$  for circuit  $\mathcal{C}$  over  
 442  $BIDB$ , we can exactly represent  $\tilde{Q}(\mathbf{X})$  as follows:

$$443 \quad \tilde{Q}(X_1, \dots, X_n) = \sum_{(\mathbf{v}, \mathbf{c}) \in \mathbf{E}(\mathcal{C})} \mathbb{1}(\mathbf{v} \bmod \mathcal{B} \neq 0) \cdot \mathbf{c} \cdot \prod_{X_i \in \text{VAR}(\mathbf{v})} X_i \quad (4)$$

444 Given the above, the algorithm is a sampling based algorithm for the above sum: we sample  
 445 (via SAMPLEMONOMIAL)  $(\mathbf{v}, \mathbf{c}) \in \mathbf{E}(\mathcal{C})$  with probability proportional to  $|\mathbf{c}|$  and compute  
 446  $Y = \mathbb{1}(\mathbf{v} \bmod \mathcal{B} \neq 0) \cdot \prod_{X_i \in \text{VAR}(\mathbf{v})} p_i$ . Taking  $N$  samples and computing the average of  $Y$   
 447 gives us our final estimate. ONEPASS is used to compute the sampling probabilities needed  
 448 in SAMPLEMONOMIAL (details are in Appendix C).

449 **5 More on Circuits and Moments**

450 We formalize our claim from Sec. 1 that a linear approximation algorithm for our problem  
 451 implies that PDB queries (under bag semantics) can be answered (approximately) in the  
 452 same runtime as deterministic queries under reasonable assumptions. Lastly, we generalize  
 453 our result for expectation to other moments.

454 **The cost model.** So far our analysis of APPROXIMATE $\tilde{Q}$  has been in terms of the size  
 455 of the lineage circuits. We now show that this model corresponds to the behavior of a  
 456 deterministic database by proving that for any  $\mathcal{RA}^+$  query  $Q$ , we can construct a compressed  
 457 circuit for  $Q$  and BIDB  $\mathbf{D}$  of size (and in runtime) linear in that of a general class of query  
 458 processing algorithms for the same query  $Q$  on a deterministic database  $D$ . We assume a  
 459 linear relationship between input sizes  $|\mathbf{D}|$  and  $|D|$  (i.e.,  $\exists c, D \in \mathbf{D}$  s.t.  $|\mathbf{D}| \leq c \cdot |D|$ ).<sup>9</sup> We  
 460 adopt a minimalistic compute-bound model of query evaluation drawn from the worst-case  
 461 optimal join literature [28, 26].

$$\begin{aligned}
 \text{cost}(R, D) &= |R| & \text{cost}(\sigma Q, D) &= \text{cost}(Q, D) & \text{cost}(\pi Q, D) &= \text{cost}(Q, D) + |Q(D)| \\
 \text{cost}(Q \cup Q', D) &= \text{cost}(Q, D) + \text{cost}(Q', D) + |Q(D)| + |Q'(D)| \\
 \text{cost}(Q_1 \bowtie \dots \bowtie Q_n, D) &= \text{cost}(Q_1, D) + \dots + \text{cost}(Q_n, D) + |Q_1(D)| \bowtie \dots \bowtie |Q_n(D)|
 \end{aligned}$$

462 Under this model a query  $Q$  evaluated over database  $D$  has runtime  $O(\text{cost}(Q, D))$ . We  
 463 assume that full table scans are used for every base relation access. We can model index  
 464 scans by treating an index scan query  $\sigma_\theta(R)$  as a base relation.

465 It can be verified that worst-case optimal join algorithms [28, 26], as well as query  
 466 evaluation via factorized databases [30] (and work on FAQs [24]) can be modeled as select-  
 467 union-project-join queries (though these queries can be data dependent).<sup>10</sup> Further, it can be  
 468 verified that the above cost model on the corresponding SPJU join queries correctly captures  
 469 their runtime.

470 We are now ready to formally state our claim from Sec. 1:

471 **► Corollary 5.1.** *Given an SPJU query  $Q$  over a TIDB  $\mathbf{D}$  and let  $D_{max}$  denote the world  
 472 containing all tuples of  $\mathbf{D}$ , we can compute a  $(1 \pm \epsilon)$ -approximation of the expectation for  
 473 each output tuple in  $Q(\mathbf{D})$  with probability at least  $1 - \delta$  in time*

$$O_k \left( \frac{1}{\epsilon^2} \cdot \text{cost}(Q, D_{max}) \cdot \log \frac{1}{\delta} \cdot \log(n) \right)$$

474 **Proof.** This follows from Lemma D.1 (Appendix D.1.2) and Corollary 4.10 (where the latter  
 475 is used with  $\delta$  being substituted<sup>11</sup> with  $\frac{\delta}{n^k}$ ).

476 **Higher Moments.** We make a simple observation to conclude the presentation of our  
 477 results. So far we have only focused on the expectation of  $Q$ . In addition, we could e.g. prove  
 478 bounds of probability of the multiplicity being at least 1. Progress can be made on this as  
 479 follows: For any positive integer  $m$  we can compute the  $m$ -th moment of the multiplicities,

<sup>9</sup> This is a reasonable assumption because each block of a BIDB represents entities with uncertain attributes. In practice there is often a limited number of alternatives for each block (e.g., which of five conflicting data sources to trust). Note that all TIDBs trivially fulfill this condition (i.e.,  $c = 1$ ).

<sup>10</sup> This claim can be verified by e.g. simply looking at the *Generic-Join* algorithm in [28] and *factorize* algorithm in [30].

<sup>11</sup> Recall that Corollary 4.10 is stated for a single output tuple so to get the required guarantee for all (at most  $n^k$ ) output tuples of  $Q$  we get at most  $\frac{\delta}{n^k}$  probability of failure for each output tuple and then just a union bound over all output tuples.

483 allowing us to e.g. use Chebyshev inequality or other high moment based probability bounds  
484 on the events we might be interested in. We leave further investigations for future work.

## 485 **6** Related Work

486 **Probabilistic Databases** (PDBs) have been studied predominantly for set semantics.  
487 Approaches for probabilistic query processing (i.e., computing marginal probabilities of  
488 tuples), fall into two broad categories. *Intensional* (or *grounded*) query evaluation computes  
489 the *lineage* of a tuple and then the probability of the lineage formula. It has been shown  
490 that computing the marginal probability of a tuple is #P-hard [36] (by reduction from  
491 weighted model counting). The second category, *extensional* query evaluation, is in PTIME,  
492 but is limited to certain classes of queries. Dalvi et al. [11] and Olteanu et al. [17] proved  
493 dichotomies for UCQs and two classes of queries with negation, respectively. Amarilli et al.  
494 investigated tractable classes of databases for more complex queries [2]. Another line of work,  
495 studies which structural properties of lineage formulas lead to tractable cases [23, 31, 33]. In  
496 this paper we focus on intensional query evaluation with polynomials.

497 Many data models have been proposed for encoding PDBs more compactly than as sets of  
498 possible worlds. These include tuple-independent databases [37] (TIDBs), block-independent  
499 databases (BIDBs) [32], and *PC-tables* [20]. Fink et al. [15] study aggregate queries over  
500 a probabilistic version of the extension of K-relations for aggregate queries proposed in [3]  
501 (*pvc-tables*). As an extension of K-relations, this approach supports bags. In contrast,  
502 we study a less general data model ( $\mathbb{N}[\mathbf{X}]$ -PDBs) and query class, but provide a linear  
503 time approximation algorithm and provide new insights into the complexity of computing  
504 expectations while [15] computes probabilities for individual output annotations.

505 Several techniques for approximating tuple probabilities have been proposed in related  
506 work [16, 12, 29, 9], relying on Monte Carlo sampling, e.g., [9], or a branch-and-bound  
507 paradigm [29]. Our approximation algorithm is also based on sampling.

508 **Compressed Encodings** are used for Boolean formulas (e.g, various types of circuits  
509 including OBDDs [22]) and polynomials (e.g., factorizations [30]) some of which have been  
510 utilized for probabilistic query processing, e.g., [22]. Compact representations for which  
511 probabilities can be computed in linear time include OBDDs, SDDs, d-DNNF, and FBDD.  
512 [13] studies circuits for absorptive semirings while [34] studies circuits that include negation  
513 (expressed as the monus operation). Algebraic Decision Diagrams [5] (ADDs) generalize  
514 BDDs to variables with more than two values. Chen et al. [7] introduced the generalized  
515 disjunctive normal form. Appendix E covers more related work on fine-grained complexity.

## 516 **7** Conclusions and Future Work

517 We have studied the problem of calculating the expectation of lineage polynomials over  
518 BIDBs. This problem has a practical application in probabilistic databases over multisets,  
519 where it corresponds to calculating the expected multiplicity of a query result tuple. While  
520 the expectation of a polynomial can be calculated in linear time for polynomials in SOP  
521 form, the problem is #W[1]-hard for factorized polynomials (proven through a reduction  
522 from the problem of counting k-matchings). We prove that it is possible to approximate the  
523 expectation of a lineage polynomial in linear time UCQs over TIDBs and BIDBs (under the  
524 assumption that there are few cancellations). Interesting directions for future work include  
525 development of a dichotomy for bag PDBs and approximations for more general data models.

## 526 — References —

- 527 1 pdbench. <http://pdbench.sourceforge.net/>. Accessed: 2020-12-15.
- 528 2 Antoine Amarilli, Pierre Bourhis, and Pierre Senellart. Probabilities and provenance via tree  
529 decompositions. *PODS*, 2015.
- 530 3 Yael Amerdamer, Daniel Deutch, and Val Tannen. Provenance for aggregate queries. In  
531 *PODS*, pages 153–164, 2011.
- 532 4 Lyublena Antova, Thomas Jansen, Christoph Koch, and Dan Olteanu. Fast and simple  
533 relational processing of uncertain data.
- 534 5 R. Iris Bahar, Erica A. Frohm, Charles M. Gaona, Gary D. Hachtel, Enrico Macii, Abelardo  
535 Pardo, and Fabio Somenzi. Algebraic decision diagrams and their applications. In *IEEE CAD*,  
536 1993.
- 537 6 Peter Bürgisser, Michael Clausen, and Mohammad Amin Shokrollahi. *Algebraic complexity*  
538 *theory*, volume 315. Springer, 1997.
- 539 7 Hubie Chen and Martin Grohe. Constraint satisfaction with succinctly specified relations. *J.*  
540 *Comput. Syst. Sci.*, 76(8):847–860, 2010.
- 541 8 Radu Curticapean. Counting matchings of size  $k$  is  $w[1]$ -hard. In *ICALP*, volume 7965, pages  
542 352–363, 2013.
- 543 9 N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. *VLDB*, 16(4):544,  
544 2007.
- 545 10 Nilesh Dalvi and Dan Suciu. The dichotomy of conjunctive queries on probabilistic structures.  
546 In *PODS*, pages 293–302, 2007.
- 547 11 Nilesh Dalvi and Dan Suciu. The dichotomy of probabilistic inference for unions of conjunctive  
548 queries. *JACM*, 59(6):30, 2012.
- 549 12 Maarten Van den Heuvel, Peter Ivanov, Wolfgang Gatterbauer, Floris Geerts, and Martin  
550 Theobald. Anytime approximation in probabilistic databases via scaled dissociations. In  
551 *SIGMOD*, pages 1295–1312, 2019.
- 552 13 Daniel Deutch, Tova Milo, Sudeepa Roy, and Val Tannen. Circuits for datalog provenance. In  
553 *ICDT*, pages 201–212, 2014.
- 554 14 Su Feng, Aaron Huber, Boris Glavic, and Oliver Kennedy. Uncertainty annotated databases -  
555 a lightweight approach for approximating certain answers. In *SIGMOD*, 2019.
- 556 15 Robert Fink, Larisa Han, and Dan Olteanu. Aggregation in probabilistic databases via  
557 knowledge compilation. *PVLDB*, 5(5):490–501, 2012.
- 558 16 Robert Fink, Jiewen Huang, and Dan Olteanu. Anytime approximation in probabilistic  
559 databases. *VLDBJ*, 22(6):823–848, 2013.
- 560 17 Robert Fink and Dan Olteanu. Dichotomies for queries with negation in probabilistic databases.  
561 *TODS*, 41(1):4:1–4:47, 2016.
- 562 18 Jörg Flum and Martin Grohe. *Parameterized Complexity Theory*. Texts in Theoretical  
563 Computer Science. An EATCS Series. Springer, 2006. doi:10.1007/3-540-29953-X.
- 564 19 Todd J. Green, Gregory Karvounarakis, and Val Tannen. Provenance semirings. In *PODS*,  
565 pages 31–40, 2007.
- 566 20 Todd J Green and Val Tannen. Models for incomplete and probabilistic information. In *EDBT*,  
567 pages 278–296. 2006.
- 568 21 Tomasz Imieliński and Witold Lipski Jr. Incomplete information in relational databases.  
569 *JACM*, 31(4):761–791, 1984.
- 570 22 Abhay Kumar Jha and Dan Suciu. Probabilistic databases with markovviews. *PVLDB*,  
571 5(11):1160–1171, 2012.
- 572 23 Batya Kenig, Avigdor Gal, and Ofer Strichman. A new class of lineage expressions over  
573 probabilistic databases computable in  $p$ -time. In *SUM*, volume 8078, pages 219–232, 2013.
- 574 24 Mahmoud Abo Khamis, Hung Q. Ngo, and Atri Rudra. Faq: Questions asked frequently. In  
575 *PODS*, pages 13–28, 2016.
- 576 25 Tsvi Kopelowitz and Virginia Vassilevska Williams. Towards optimal set-disjointness and  
577 set-intersection data structures. In *ICALP*, volume 168, pages 74:1–74:16, 2020.



- 578 26 Hung Q. Ngo. Worst-case optimal join algorithms: Techniques, results, and open problems.  
579 In *PODS*, 2018.
- 580 27 Hung Q. Ngo, Ely Porat, Christopher Ré, and Atri Rudra. Worst-case optimal join algorithms.  
581 *J. ACM*, 65(3):16:1–16:40, 2018.
- 582 28 Hung Q. Ngo, Christopher Ré, and Atri Rudra. Skew strikes back: new developments in the  
583 theory of join algorithms. *SIGMOD Rec.*, 42(4):5–16, 2013.
- 584 29 Dan Olteanu, Jiewen Huang, and Christoph Koch. Approximate confidence computation in  
585 probabilistic databases. In *ICDE*, pages 145–156, 2010.
- 586 30 Dan Olteanu and Maximilian Schleich. Factorized databases. *SIGMOD Rec.*, 45(2):5–16, 2016.
- 587 31 Sudeepa Roy, Vittorio Perduca, and Val Tannen. Faster query answering in probabilistic  
588 databases using read-once functions. In *ICDT*, 2011.
- 589 32 C. Ré and D. Suciu. Materialized views in probabilistic databases: for information exchange  
590 and query optimization. In *VLDB*, pages 51–62, 2007.
- 591 33 Prithviraj Sen, Amol Deshpande, and Lise Getoor. Read-once functions and query evaluation  
592 in probabilistic databases. *PVLDB*, 3(1):1068–1079, 2010.
- 593 34 Pierre Senellart. Provenance and probabilities in relational databases. *SIGMOD Record*,  
594 46(4):5–15, 2018.
- 595 35 Dan Suciu, Dan Olteanu, Christopher Ré, and Christoph Koch. *Probabilistic Databases*.  
596 Morgan & Claypool Publishers, 2011.
- 597 36 Leslie G. Valiant. The complexity of enumeration and reliability problems. *SIAM J. Comput.*,  
598 8(3):410–421, 1979.
- 599 37 Guy Van den Broeck and Dan Suciu. Query processing on probabilistic data: A survey. 2017.
- 600 38 Virginia Vassilevska Williams. Some open problems in fine-grained complexity. *SIGACT News*,  
601 49(4):29–35, 2018. doi:10.1145/3300150.3300158.

## 602 **8 Acknowledgements**

603 We thank Virginia Williams for showing us Eq. (17), which greatly simplified our earlier  
604 proof of Lemma 3.8, and for graciously allowing us to use it.

605 **A** Missing details from Section 2

606 **A.1**  $\mathcal{K}$ -relations and  $\mathbb{N}[\mathbf{X}]$ -PDBs

607 We use  $K$ -relations to model bags. A  $K$ -relation [19] is a relation whose tuples are annotated  
 608 with elements from a commutative semiring  $\mathcal{K} = (K, \oplus_{\mathcal{K}}, \otimes_{\mathcal{K}}, \mathbf{0}_{\mathcal{K}}, \mathbf{1}_{\mathcal{K}})$ . A commutative  
 609 semiring is a structure with a domain  $K$  and associative and commutative binary operations  
 610  $\oplus_{\mathcal{K}}$  and  $\otimes_{\mathcal{K}}$  such that  $\otimes_{\mathcal{K}}$  distributes over  $\oplus_{\mathcal{K}}$ ,  $\mathbf{0}_{\mathcal{K}}$  is the identity of  $\oplus_{\mathcal{K}}$ ,  $\mathbf{1}_{\mathcal{K}}$  is the identity  
 611 of  $\otimes_{\mathcal{K}}$ , and  $\mathbf{0}_{\mathcal{K}}$  annihilates all elements of  $K$  when combined by  $\otimes_{\mathcal{K}}$ . Let  $\mathcal{U}$  be a countable  
 612 domain of values. Formally, an  $n$ -ary  $\mathcal{K}$ -relation over  $\mathcal{U}$  is a function  $R : \mathcal{U}^n \rightarrow K$  with  
 613 finite support  $\text{supp}(R) = \{t \mid R(t) \neq \mathbf{0}_{\mathcal{K}}\}$ . A  $\mathcal{K}$ -database is a set of  $\mathcal{K}$ -relations. It will be  
 614 convenient to also interpret a  $\mathcal{K}$ -database as a function from tuples to annotations. Thus,  
 615  $R(t)$  (resp.,  $D(t)$ ) denotes the annotation associated by  $\mathcal{K}$ -relation  $R$  ( $\mathcal{K}$ -database  $D$ ) to  $t$ .

616 For completeness, we briefly review the semantics for  $\mathcal{RA}^+$  queries over  $\mathcal{K}$ -relations [19]  
 617 illustrated in Fig. 2. In Fig. 2, we use  $[\![\cdot]\!]_D$  to denote the result of evaluating query  $Q$   
 618 over  $\mathcal{K}$ -database  $D$ , assume that tuples are of appropriate arity, use  $\text{sch}(R)$  to denote the  
 619 attributes of  $R$ , and use  $\pi_A(t)$  to denote the projection of tuple  $t$  on a list of attributes  $A$ .  
 620 Furthermore,  $\theta(t)$  denotes the (Boolean) result of evaluating condition  $\theta$  over  $t$ .

621 Consider the semiring  $\mathbb{N} = (\mathbb{N}, +, \times, 0, 1)$  of natural numbers.  $\mathbb{N}$ -databases model bag  
 622 semantics by annotating each tuple with its multiplicity. A probabilistic  $\mathbb{N}$ -database ( $\mathbb{N}$ -PDB)  
 623 is a PDB where each possible world is an  $\mathbb{N}$ -database. We study the problem of computing  
 624 statistical moments for query results over such databases. Specifically, given a probabilistic  
 625  $\mathbb{N}$ -database  $\mathcal{D} = (\Omega, \mathbf{P})$ , query  $Q$ , and possible result tuple  $t$ , we use  $Q(\mathcal{D})(t)$  for  $D \in \Omega$  as  
 626 input in RHS of Eq. (1) to compute the expected multiplicity of  $t$ . Note that the tables of  
 627 Fig. 1 have an implicit 1  $\mathbb{N}$ -valued annotation for each tuple in tables *OnTime* and *Route*.  
 628 Intuitively, the expectation of  $Q(\mathcal{D})(t)$  is the number of duplicates of  $t$  we expect to find in  
 629 result of query  $Q$ .

630 Let  $\mathbb{N}[\mathbf{X}]$  denote the set of polynomials over variables  $\mathbf{X} = (X_1, \dots, X_n)$  with natural  
 631 number coefficients and exponents. Consider now the semiring  $(\mathbb{N}[\mathbf{X}], +, \cdot, 0, 1)$  whose domain  
 632 is  $\mathbb{N}[\mathbf{X}]$ , with the standard addition and multiplication of polynomials. We will use  $\mathbb{N}[\mathbf{X}]$ -PDB  
 633  $\mathbf{D}$ , defined as the tuple  $(\Omega_{\mathbb{N}[\mathbf{X}]}, \mathbf{P})$ , where  $\mathbb{N}[\mathbf{X}]$ -database  $\Omega_{\mathbb{N}[\mathbf{X}]}$  is paired with probability  
 634 distribution  $\mathbf{P}$ . We denote by  $Q_t$  the annotation of tuple  $t$  in the result of  $Q$  on an  
 635 implicit  $\mathbb{N}[\mathbf{X}]$ -PDB (i.e.,  $Q_t = Q(\mathbf{D})(t)$  for some  $\mathbf{D}$ ) and as before, interpret it as a function  
 636  $Q_t : \{0, 1\}^{|\mathbf{X}|} \rightarrow \mathbb{N}$  from vectors of variable assignments to the corresponding value of the  
 637 annotating polynomial.  $\mathbb{N}[\mathbf{X}]$ -PDBs and a function *Mod* (which transforms an  $\mathbb{N}[\mathbf{X}]$ -PDB to  
 638 an equivalent  $\mathbb{N}$ -PDB) are both formalized next.

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

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

$$646 \quad \text{Mod}(Q(M)) = Q(\text{Mod}(M))$$

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

$$649 \quad \text{Mod}(M) = \mathcal{D}$$

As mentioned above we will use  $\mathbb{N}[\mathbf{X}]$ -databases paired with a probability distribution as a representation system. We refer to such databases as  $\mathbb{N}[\mathbf{X}]$ -PDBs and use bold symbols to distinguish them from possible worlds (which are  $\mathbb{N}$ -databases). Formally, an  $\mathbb{N}[\mathbf{X}]$ -PDB is an  $\mathbb{N}[\mathbf{X}]$ -database  $\Omega_{\mathbb{N}[\mathbf{X}]}$  and a probability distribution  $\mathbf{P}$  over assignments  $\varphi$  of the variables  $\mathbf{X} = \{X_1, \dots, X_n\}$  occurring in annotations of  $\Omega_{\mathbb{N}[\mathbf{X}]}$  to  $\{0, 1\}$ . Note that an assignment  $\varphi : \mathbf{X} \rightarrow \{0, 1\}^n$  can be represented as a vector  $\mathbf{w} \in \{0, 1\}^n$  where  $\mathbf{w}[i]$  records the value assigned to variable  $X_i$ . Thus, from now on we will solely use such vectors which we refer to as *world vectors* and implicitly understand them to represent assignments. Given an assignment  $\varphi$  we use  $\varphi(\mathbf{D})$  to denote the semiring homomorphism  $\mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$  that applies the assignment  $\varphi$  to all variables of a polynomial and evaluates the resulting expression in  $\mathbb{N}$ .

**Definition A.2** ( $\mathbb{N}[\mathbf{X}]$ -PDBs). An  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$  over variables  $\mathbf{X} = \{X_1, \dots, X_n\}$  is a tuple  $(\Omega_{\mathbb{N}[\mathbf{X}]}, \mathbf{P})$  where  $D$  is an  $\mathbb{N}[\mathbf{X}]$ -database and  $\mathbf{P}$  is a probability distribution over  $\mathbf{w} \in \{0, 1\}^n$ . We use  $\varphi_{\mathbf{w}}$  to denote the assignment corresponding to  $\mathbf{w} \in \{0, 1\}^n$ . The  $\mathbb{N}$ -PDB  $Mod(\mathbf{D}) = (\Omega, \mathbf{P}')$  encoded by  $\mathbf{D}$  is defined as:

$$\Omega = \{\varphi_{\mathbf{w}}(\mathbf{D}) \mid \mathbf{w} \in \{0, 1\}^n\}$$

$$\forall D \in \Omega : P'(D) = \sum_{\mathbf{w} \in \{0, 1\}^n : \varphi_{\mathbf{w}}(\mathbf{D}) = D} P(\mathbf{w})$$

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

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

Importantly, as the following proposition shows, any finite  $\mathbb{N}$ -PDB can be encoded as an  $\mathbb{N}[\mathbf{X}]$ -PDB and  $\mathbb{N}[\mathbf{X}]$ -PDBs are closed under positive relational algebra queries, the class of queries we are interested in in this work.

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

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

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

The probability distribution  $\mathbf{P}'$  assigns all world vectors zero probability except for  $|\Omega|$  world vectors (representing the possible worlds)  $\mathbf{w}_i$ . All elements of  $\mathbf{w}_i$  are zero except for the position corresponding to variables  $X_i$  which is set to 1. Unfolding definitions it is trivial to show that  $Mod(\mathbf{D}) = \mathcal{D}$ . Thus,  $\mathbb{N}[\mathbf{X}]$  are a complete representation system. The closure under  $\mathcal{RA}^+$  queries follows from the fact that an assignment  $\mathbf{X} \rightarrow \{0, 1\}$  is a semiring homomorphism and that semiring homomorphisms commute with queries over  $\mathcal{K}$ -relations.

Now let us consider computing the expected multiplicity of a tuple  $t$  in the result of a query  $Q$  over an  $\mathbb{N}$ -PDB  $\mathcal{D}$  using the annotation of  $t$  in the result of evaluating  $Q$  over an

690  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$  for which  $Mod(\mathbf{D}) = \mathcal{D}$ . The expectation of the polynomial  $Q = Q(\mathbf{D})(t)$   
 691 based on the probability distribution of  $\mathbf{D}$  over the variables in  $\mathbf{D}$  is:

$$692 \quad \mathbb{E}_{\mathbf{W} \sim \mathbf{P}} [Q(\mathbf{W})] = \sum_{\mathbf{w} \in \{0,1\}^n} \varphi_{\mathbf{w}}(Q(\mathbf{D})(t)) \cdot P(\mathbf{w}) \quad (5)$$

693 Since  $\mathbb{N}[\mathbf{X}]$ -PDBs  $\mathbf{D}$  are a complete representation system for  $\mathbb{N}$ -PDBs which are closed  
 694 under  $\mathcal{RA}^+$ , computing the expectation of the multiplicity of a tuple  $t$  in the result of an  
 695  $\mathcal{RA}^+$  query over the  $\mathbb{N}$ -PDB  $Mod(\mathbf{D})$ , is the same as computing the expectation of the  
 696 polynomial  $Q(\mathbf{D})(t)$ .  $\blacktriangleleft$

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

698 Two important subclasses of  $\mathbb{N}[\mathbf{X}]$ -PDBs that are of interest to us are the bag versions of  
 699 tuple-independent databases (TIDBs) and block-independent databases (BIDBs). Under set  
 700 semantics, a TIDB is a deterministic database  $D$  where each tuple  $t$  is assigned a probability  
 701  $p_t$ . The set of possible worlds represented by a TIDB  $D$  is all subsets of  $D$ . The probability  
 702 of each world is the product of the probabilities of all tuples that exist with one minus  
 703 the probability of all tuples of  $D$  that are not part of this world, i.e., tuples are treated  
 704 as independent random events. In a BIDB, we also assign each tuple a probability, but  
 705 additionally partition  $D$  into blocks. The possible worlds of a BIDB  $D$  are all subsets of  $D$   
 706 that contain at most one tuple from each block. Note then that the tuples sharing the same  
 707 block are disjoint, and the sum of the probabilities of all the tuples in the same block  $b$  is 1.  
 708 The probability of such a world is the product of the probabilities of all tuples present in the  
 709 world. For bag TIDBs and BIDBs, we define the probability of a tuple to be the probability  
 710 that the tuple exists with multiplicity at least 1.

711 As already noted above, in this work, we define TIDBs and BIDBs as subclasses of  
 712  $\mathbb{N}[\mathbf{X}]$ -PDBs. In this work, we consider one further deviation from the standard: We use bag  
 713 semantics for queries. Even though tuples cannot occur more than once in the input TIDB  
 714 or BIDB, they can occur with a multiplicity larger than one in the result of a query. Since  
 715 in TIDBs and BIDBs, there is a one-to-one correspondence between tuples in the database  
 716 and variables, we can interpret a vector  $\mathbf{w} \in \{0,1\}^n$  as denoting which tuples exist in the  
 717 possible world  $\varphi_{\mathbf{w}}(\mathbf{D})$  (the ones where  $\mathbf{w}[j] = 1$ ). For BIDBs specifically, note that at  
 718 most one of the bits corresponding to tuples in each block will be set (i.e., for any pair of  
 719 bits  $w_j, w_{j'}$  that are part of the same block  $b_i \supseteq \{t_{i,j}, t_{i,j'}\}$ , at most one of them will be set).  
 720 Denote the vector  $\mathbf{p}$  to be a vector whose elements are the individual probabilities  $p_i$  of each  
 721 tuple  $t_i$ . Let  $\mathbf{P}(\mathbf{p})$  denote the distribution induced by  $\mathbf{p}$ .

$$722 \quad \mathbb{E}_{\mathbf{W} \sim \mathbf{P}(\mathbf{p})} [Q(\mathbf{W})] = \sum_{\substack{\mathbf{w} \in \{0,1\}^n \\ s.t. w_j, w_{j'} = 1 \rightarrow \exists b_i \supseteq \{t_{i,j}, t_{i,j'}\}}} Q(\mathbf{w}) \prod_{\substack{j \in [n] \\ s.t. w_j = 1}} p_j \prod_{\substack{j \in [n] \\ s.t. w_j = 0}} (1 - p_i) \quad (6)$$

724 Recall that tuple blocks in a TIDB always have size 1, so the outer summation of eq. (6) is  
 725 over the full set of vectors.

## 726 A.3 Proof of Proposition 2.1

727 **Proof.** We need to prove for  $\mathbb{N}$ -PDB  $\mathcal{D} = (\Omega, \mathbf{P})$  and  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D} = (D', \mathbf{P}')$  where  
 728  $Mod(\mathbf{D}) = \mathcal{D}$  that  $\mathbb{E}_{D \sim \mathbf{P}} [Q(D)(t)] = \mathbb{E}_{\mathbf{W} \sim \mathbf{P}'} [Q_t(\mathbf{W})]$  By expanding  $Q_t$  and the expectation

729 we have:

$$730 \quad \mathbb{E}_{\mathbf{W} \sim \mathbf{P}'} [Q_t(\mathbf{W})] = \sum_{\mathbf{w} \in \{0,1\}^n} P'(\mathbf{w}) \cdot Q(\mathbf{D})(t)(\mathbf{w})$$

731 From  $Mod(\mathbf{D}) = \mathcal{D}$ , we have that the range of  $\varphi_{\mathbf{w}(\mathbf{D})}$  is  $\Omega$ , so

$$732 \quad = \sum_{D \in \Omega} \sum_{\mathbf{w} \in \{0,1\}^n : \varphi_{\mathbf{w}(\mathbf{D})} = D} P'(\mathbf{w}) \cdot Q(\mathbf{D})(t)(\mathbf{w})$$

733 In the inner sum,  $\varphi_{\mathbf{w}(\mathbf{D})} = D$ , so by distributivity of  $+$  over  $\times$

$$734 \quad = \sum_{D \in \Omega} Q(D)(t) \sum_{\mathbf{w} \in \{0,1\}^n : \varphi_{\mathbf{w}(\mathbf{D})} = D} P'(\mathbf{w})$$

735 From the definition of  $P$ , given  $Mod(\mathbf{D}) = \mathcal{D}$ , we get

$$736 \quad = \sum_{D \in \Omega} Q(D)(t) \cdot P(D) = \mathbb{E}_{D \sim \mathbf{P}} [Q(D)(t)]$$

737

#### 742 A.4 Lemma A.4

743 ► **Lemma A.4.** If  $Q(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \{0, \dots, B\}^n} q_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n X_i^{d_i}$  then  $\tilde{Q}(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \eta} q_{\mathbf{d}} \cdot$

$$744 \quad \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n X_i$$

745 **Proof.** Follows by the construction of  $\tilde{Q}$  in definition 2.6. ◀

#### 746 A.5 Proposition A.5

747 Note the following fact:

748 ► **Proposition A.5.** For any BIDB-lineage polynomial  $Q(X_1, \dots, X_n)$  and all  $\mathbf{w} \in \eta$ , it holds  
749 that  $Q(\mathbf{w}) = \tilde{Q}(\mathbf{w})$ .

750 **Proof.** Note that any  $Q$  in factorized form is equivalent to its SMB expansion. For each  
751 term in the expanded form, further note that for all  $b \in \{0, 1\}$  and all  $e \geq 1$ ,  $b^e = b$ . ◀

#### 752 A.6 Proof for Lemma 2.8

753 **Proof.** Let  $Q$  be the generalized polynomial, i.e., the polynomial of  $n$  variables with highest  
754 degree =  $B$ :

$$755 \quad Q(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \{0, \dots, B\}^n} q_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n X_i^{d_i}$$

756 . Then, in expectation we have

$$757 \quad \mathbb{E}_{\mathbf{W}} [Q(\mathbf{W})] = \sum_{\mathbf{d} \in \eta} q_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[ \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n w_i^{d_i} \right] \quad (7)$$

$$758 \quad = \sum_{\mathbf{d} \in \eta} q_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n \mathbb{E}_{\mathbf{w}} [w_i^{d_i}] \quad (8)$$

$$759 \quad = \sum_{\mathbf{d} \in \eta} q_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n \mathbb{E}_{\mathbf{w}} [w_i] \quad (9)$$

$$760 \quad = \sum_{\mathbf{d} \in \eta} q_{\mathbf{d}} \cdot \prod_{\substack{i=1 \\ s.t. d_i \geq 1}}^n p_i \quad (10)$$

$$761 \quad = \tilde{Q}(p_1, \dots, p_n) \quad (11)$$

763 In steps eq. (7) and eq. (8), by linearity of expectation (recall the variables are independent,  
764 or the monomial expectation is 0), the expectation can be pushed all the way inside of the  
765 product. In eq. (9), note that  $w_i \in \{0, 1\}$  which further implies that for any exponent  $e \geq 1$ ,  
766  $w_i^e = w_i$ . Next, in eq. (10) the expectation of a tuple is indeed its probability.

767 Finally, observe Eq. (11) by construction in Lemma A.4, that  $\tilde{Q}(p_1, \dots, p_n)$  is exactly  
768 the product of probabilities of each variable in each monomial across the entire sum. ◀

## 769 A.7 Proof For Corollary 2.9

770 **Proof.** Note that lemma 2.8 shows that  $\mathbb{E}[Q] = \tilde{Q}(p_1, \dots, p_n)$ . Therefore, if  $Q$  is already  
771 in SMB form, one only needs to compute  $Q(p_1, \dots, p_n)$  ignoring exponent terms (note that  
772 such a polynomial is  $\tilde{Q}(p_1, \dots, p_n)$ ), which indeed has  $O(\text{SMB}(|Q|))$  computations. ◀

## 773 B Missing details from Section 3

774 We use Lemma 3.5 to prove Theorem 3.4:

### 775 B.1 Proof of Theorem 3.4

776 **Proof.** For the sake of contradiction, let us assume we can solve our problem in  $f(k) \cdot m^c$  time  
777 for some absolute constant  $c$ . Then given a graph  $G$  we can compute the query polynomial or  
778 rather, expression tree representation of  $\tilde{Q}_G^k$  (in the obvious way) in  $O(km)$  time. Then after  
779 we run our algorithm on  $\tilde{Q}_G^k$ , we get  $\tilde{Q}_G^k(p_i, \dots, p_i)$  for  $0 \leq i \leq 2k$  in additional  $f(k) \cdot m^c$   
780 time. Lemma 3.5 then computes the number of  $k$ -matchings in  $G$  in  $O(k^3)$  time. Thus,  
781 overall we have an algorithm for computing the number of  $k$ -matchings in time

$$782 \quad O(km) + f(k) \cdot m^c + O(k^3) \leq (O(k^3) + f(k)) \cdot m^{c+1}$$

$$783 \quad \leq (O(k^3) + f(k)) \cdot n^{2c+2},$$

785 which contradicts Theorem 3.1. ◀

### 786 B.2 Proof of Lemma 3.5

787 **Proof.** We first argue that  $\tilde{Q}_G^k(p, \dots, p) = \sum_{i=0}^{2k} c_i \cdot p^i$ . First, since  $Q_G(\mathbf{X})$  has degree 2, it  
788 follows that  $Q_G^k(\mathbf{X})$  has degree  $2k$ . By definition,  $\tilde{Q}_G^k(\mathbf{X})$  sets every exponent  $e > 1$  to  $e = 1$ ,  
789 which means that  $\text{DEG}(\tilde{Q}_G^k) \leq \text{DEG}(Q_G^k) = 2k$ . Thus, if we think of  $p$  as a variable, then

790  $\tilde{Q}_G^k(p, \dots, p)$  is a univariate polynomial of degree at most  $\text{DEG}(\tilde{Q}_G^k) \leq 2k$ . Thus, we can  
 791 write

$$792 \quad \tilde{Q}_G^k(p, \dots, p) = \sum_{i=0}^{2k} c_i p^i$$

793 We note that  $c_i$  is *exactly* the number of monomials in the SMB expansion of  $Q_G^k(\mathbf{X})$  composed  
 794 of  $i$  distinct variables.<sup>12</sup>

795 Given that we then have  $2k + 1$  distinct values of  $\tilde{Q}_G^k(p, \dots, p)$  for  $0 \leq i \leq 2k$ , it follows  
 796 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  $(p_i^0 \dots p_i^{2k})$ ,  
 797  $\mathbf{c}$  is the coefficient vector  $(c_0, \dots, c_{2k})$ , and  $\mathbf{b}$  is the vector such that  $\mathbf{b}[i] = \tilde{Q}_G^k(p_i, \dots, p_i)$ .  
 798 In other words, matrix  $\mathbf{M}$  is the Vandermonde matrix, from which it follows that we have  
 799 a matrix with full rank (the  $p_i$ 's are distinct), and we can solve the linear system in  $O(k^3)$   
 800 time (e.g., using Gaussian Elimination) to determine  $\mathbf{c}$  exactly. Thus, after  $O(k^3)$  work, we  
 801 know  $\mathbf{c}$  and in particular,  $c_{2k}$  exactly. Next, we show why we can compute  $\#(G, \mathfrak{I} \dots \mathfrak{I}^k)$   
 802 from  $c_{2k}$  in  $O(1)$  additional time. We claim that  $c_{2k}$  is  $k! \cdot \#(G, \mathfrak{I} \dots \mathfrak{I}^k)$ . This can be seen  
 803 intuitively by looking at the original factorized representation

$$804 \quad Q_G^k(\mathbf{X}) = \sum_{(i_1, j_1), \dots, (i_k, j_k) \in E} X_{i_1} X_{j_1} \dots X_{i_k} X_{j_k},$$

805 where across each of the  $k$  products, an arbitrary  $k$ -matching can be selected  $\prod_{i=1}^k i = k!$   
 806 times. Indeed, note that each  $k$ -matching  $(i_1, j_1) \dots (i_k, j_k)$  in  $G$  corresponds to the monomial  
 807  $\prod_{\ell=1}^k X_{i_\ell} X_{j_\ell}$  in  $Q_G^k(\mathbf{X})$ , with distinct indexes. Second, the only surviving monomials  
 808  $\prod_{\ell=1}^k X_{i_\ell} X_{j_\ell}$  of degree exactly  $2k$  in  $\tilde{Q}_G^k(\mathbf{X})$  must have that all of  $i_1, j_1, \dots, i_k, j_k$  are distinct  
 809 in  $Q_G^k(\mathbf{X})$ . By the last two statements, only monomials composed of  $2k$  distinct variables in  
 810  $Q_G^k(\mathbf{X})$  (and hence of degree  $2k$  in  $\tilde{Q}_G^k(\mathbf{X})$ ) correspond to a  $k$ -matching in  $G$ .

811 Notice that each of the  $k!$  permutations of an arbitrary monomial maps to the same  
 812 distinct  $k$ -matching in  $G$ , and this implies a  $k!$  to 1 mapping between degree  $2k$  monomials  
 813 in  $\tilde{Q}_G^k(\mathbf{X})$  and  $k$ -matchings in  $G$ . It then follows that  $c_{2k} = k! \cdot \#(G, \mathfrak{I} \dots \mathfrak{I}^k)$ . Thus, simply  
 814 dividing  $c_{2k}$  by  $k!$  gives us  $\#(G, \mathfrak{I} \dots \mathfrak{I}^k)$ , as needed.  $\blacktriangleleft$

### 815 B.3 Subgraph Notation and $O(1)$ Closed Formulas

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

- 818 ■ Single Edge ( $\mathfrak{I}$ )
- 819 ■ 2-path ( $\mathfrak{A}$ )
- 820 ■ 2-matching ( $\mathfrak{I}\mathfrak{I}$ )
- 821 ■ 3-star ( $\mathfrak{A}\mathfrak{A}$ )—this is the graph that results when all three edges share exactly one common  
 822 endpoint. The remaining endpoint for each edge is disconnected from any endpoint of  
 823 the remaining two edges.
- 824 ■ Disjoint Two-Path ( $\mathfrak{I} \mathfrak{A}$ )—this subgraph consists of a two-path and a remaining disjoint  
 825 edge.

<sup>12</sup>Since  $\tilde{Q}_G^k(\mathbf{X})$  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.

826 For any graph  $G$ , the following formulas for  $\#(G, H)$  compute their respective patterns  
827 exactly in  $O(m)$  time, with  $d_i$  representing the degree of vertex  $i$  (proofs are in Appendix B.4):

$$828 \quad \#(G, \mathfrak{!}) = m, \quad (12)$$

$$829 \quad \#(G, \mathfrak{!}) = \sum_{i \in V} \binom{d_i}{2} \quad (13)$$

$$830 \quad \#(G, \mathfrak{!!}) = \sum_{(i,j) \in E} \frac{m - d_i - d_j + 1}{2} \quad (14)$$

$$831 \quad \#(G, \mathfrak{!}) = \sum_{i \in V} \binom{d_i}{3} \quad (15)$$

$$832 \quad \#(G, \mathfrak{!}) + 3\#(G, \mathfrak{!!!}) = \sum_{(i,j) \in E} \binom{m - d_i - d_j + 1}{2} \quad (16)$$

$$833 \quad \#(G, \mathfrak{!!}) + 3\#(G, \mathfrak{!}) = \sum_{(i,j) \in E} (d_i - 1) \cdot (d_j - 1) \quad (17)$$

834  
835

#### 836 B.4 Proofs of Eq. (12)-Eq. (17)

837 The proofs for Eq. (12), Eq. (13) and Eq. (15) are immediate.

838 **Proof of Eq. (14).** For edge  $(i, j)$  connecting arbitrary vertices  $i$  and  $j$ , finding all other  
839 edges in  $G$  disjoint to  $(i, j)$  is equivalent to finding all edges that are not connected to either  
840 vertex  $i$  or  $j$ . The number of such edges is  $m - d_i - d_j + 1$ , where we add 1 since edge  $(i, j)$   
841 is removed twice when subtracting both  $d_i$  and  $d_j$ . Since the summation is iterating over  
842 all edges such that a pair  $((i, j), (k, \ell))$  will also be counted as  $((k, \ell), (i, j))$ , division by 2  
843 then eliminates this double counting. Note that  $m$  and  $d_i$  for all  $i \in V$  can be computed in  
844 one pass over the set of edges by simply maintaining counts for each quantity. Finally, the  
845 summation is also one traversal through the set of edges where each operation is either a  
846 lookup ( $O(1)$  time) or an addition operation (also  $O(1)$ ) time. ◀

847 **Proof of Eq. (16).** Eq. (16) is true for similar reasons. For edge  $(i, j)$ , it is necessary to find  
848 two additional edges, disjoint or connected. As in our argument for Eq. (14), once the number  
849 of edges disjoint to  $(i, j)$  have been computed, then we only need to consider all possible  
850 combinations of two edges from the set of disjoint edges, since it doesn't matter if the two  
851 edges are connected or not. Note, the factor 3 of  $\mathfrak{!!!}$  is necessary to account for the triple  
852 counting of 3-matchings. It is also the case that, since the two path in  $\mathfrak{!}$   $\mathfrak{!}$  is connected, that  
853 there will be no double counting by the fact that the summation automatically disconnects  
854 the current edge, meaning that a two matching at the current vertex will not be counted. The  
855 sum over all such edge combinations is precisely then  $\#(G, \mathfrak{!}) + 3\#(G, \mathfrak{!!!})$ . Note that  
856 all  $d_i$  and  $d_i - 3$  factorials can be computed in  $O(m)$  time, and then each combination  $\binom{n}{3}$   
857 can be performed with constant time operations, yielding the claimed  $O(m)$  run time. ◀

858 **Proof of Eq. (17).** To compute  $\#(G, \mathfrak{!!})$ , note that for an arbitrary edge  $(i, j)$ , a 3-path  
859 exists for edge pair  $(i, \ell)$  and  $(j, k)$  where  $i, j, k, \ell$  are distinct. Further, the quantity  $(d_i -$   
860  $1) \cdot (d_j - 1)$  represents the number of 3-edge subgraphs with middle edge  $(i, j)$  and outer  
861 edges  $(i, \ell), (j, k)$  such that  $\ell \neq j$  and  $k \neq i$ . When  $k = \ell$ , the resulting subgraph is a triangle,  
862 and when  $k \neq \ell$ , the subgraph is a 3-path. Summing over all edges  $(i, j)$  gives Eq. (17) by  
863 observing that each triangle is counted thrice, while each 3-path is counted just once. For



864 reasons similar to Eq. (14), all  $d_i$  can be computed in  $O(m)$  time and each summand can  
865 then be computed in  $O(1)$  time, yielding an overall  $O(m)$  run time. ◀

## 866 B.5 Proof of Theorem 3.6

867 **Proof.** For the sake of contradiction, assume that for any  $G$ , we can compute  $\tilde{Q}_G^3(p, \dots, p)$   
868 in  $o(m^{1+\epsilon_0})$  time. Let  $G$  be the input graph. It is easy to see that one can compute the  
869 expression tree for  $Q_G^3(\mathbf{X})$  in  $O(m)$  time. Then by Theorem 3.7 we can compute  $\#(G, \mathfrak{A})$   
870 in further time  $o(m^{1+\epsilon_0}) + O(m)$ . Thus, the overall, reduction takes  $o(m^{1+\epsilon_0}) + O(m) =$   
871  $o(m^{1+\epsilon_0})$  time, which violates Conjecture 3.2. ◀

## 872 B.6 Tools to prove Lemma 3.8

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

875 ▶ **Lemma B.1.** *For any  $p$ , we have:*

$$876 \quad \tilde{Q}_G^3(p, \dots, p) = \#(G, \mathfrak{I})p^2 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{II})p^4 + 6\#(G, \mathfrak{B})p^3 \\ 877 \quad \quad \quad + 6\#(G, \mathfrak{AA})p^4 + 6\#(G, \mathfrak{I\mathfrak{I}})p^4 + 6\#(G, \mathfrak{I\mathfrak{A}})p^5 + 6\#(G, \mathfrak{III})p^6. \quad (18)$$

### 879 B.6.1 Proof for Lemma B.1

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

$$881 \quad Q_G^3(\mathbf{X}) = \sum_{(i_1, j_1), (i_2, j_2), (i_3, j_3) \in E} \prod_{\ell=1}^3 X_{i_\ell} X_{j_\ell}.$$

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

886 **CASE 1:**  $e_1 = e_2 = e_3$  (all edges are the same). There are exactly  $m = \#(G, \mathfrak{I})$  such  
887 triples, each with a  $p^2$  factor in  $\tilde{Q}_G^3(p, \dots, p)$ .

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

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

900 Since  $p$  is fixed, Lemma B.1 gives us one linear equation in  $\#(G, \mathfrak{A})$  and  $\#(G, \mathfrak{III})$  (we  
901 can handle the other counts due to equations (12)-(17)). However, we need to generate  
902 one more independent linear equation in these two variables. Towards this end we generate  
903 another graph related to  $G$ :

904 ► **Definition B.2.** For  $\ell > 1$ , let graph  $G^{(\ell)}$  be a graph generated from an arbitrary graph  
 905  $G^{(1)}$ , by replacing every edge  $e$  of  $G^{(1)}$  with a  $\ell$ -path, such that all inner vertexes of an  $\ell$ -path  
 906 replacement edge are disjoint from the inner vertexes of any other  $\ell$ -path replacement edge.

907 Next, we relate the various sub-graph counts in  $G^{(2)}$  to  $G^{(1)}$  ( $G$ ).

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

$$909 \quad \# \left( G^{(2)}, \text{⚡⚡⚡} \right) = 8 \cdot \# \left( G^{(1)}, \text{⚡⚡⚡} \right) + 6 \cdot \# \left( G^{(1)}, \text{⚡} \text{ ⚡} \right) \\ 910 \quad \quad \quad + 4 \cdot \# \left( G^{(1)}, \text{⚡} \right) + 4 \cdot \# \left( G^{(1)}, \text{⚡} \right) + 2 \cdot \# \left( G^{(1)}, \text{⚡} \right). \\ 911$$

912 ► **Lemma B.4.** For  $\ell > 1$  and any graph  $G^{(\ell)}$ ,  $\# \left( G^{(\ell)}, \text{⚡} \right) = 0$ .

### 913 B.7 Proof of Theorem 3.7

914 **Proof.** We can compute  $G^{(2)}$  from  $G^{(1)}$  in  $O(m)$  time. Additionally, if in time  $O(T(m))$ , we  
 915 have  $\tilde{Q}_{G^{(e)}}^3(p, \dots, p)$  for  $\ell \in [2]$ , then the theorem follows by Lemma 3.8. ◀ In other words,  
 916 if Theorem 3.7 holds, then so must Theorem 3.6.

### 917 B.8 Proofs for Lemma B.3, Lemma B.4, and Lemma 3.8

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

919 ► **Definition B.5.** For  $\ell > 1$ , we use  $E_\ell$  to denote the set of edges in  $G^{(\ell)}$ . For any graph  
 920  $G^{(\ell)}$ , its edges are denoted by the a pair  $(e, b)$ , such that  $b \in \{0, \dots, \ell - 1\}$  and  $e \in E_1$ , where  
 921  $(e, 0), \dots, (e, \ell - 1)$  is the  $\ell$ -path that replaces the edge  $e$ .

922 ► **Definition B.6** ( $E_S^{(\ell)}$ ). Given an arbitrary subgraph  $S^{(1)}$  of  $G^{(1)}$ , let  $E_S^{(1)}$  denote the set of  
 923 edges in  $S^{(1)}$ . Define then  $E_S^{(\ell)}$  for  $\ell > 1$  as the set of edges in the generated subgraph  $S^{(\ell)}$   
 924 (i.e. when we apply Definition B.2 to  $S^{(1)}$ ).

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

927 ► **Definition B.7.** Let  $\binom{E}{t}$  denote the set of subsets in  $E$  with exactly  $t$  edges. In a similar  
 928 manner,  $\binom{E}{\leq t}$  is used to mean the subsets of  $E$  with  $t$  or fewer edges.

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

931 ► **Definition B.8.** Let  $f_\ell : \binom{E_\ell}{3} \mapsto \binom{E_1}{\leq 3}$  be defined as follows. For any element  $s \in \binom{E_\ell}{3}$  such  
 932 that  $s = \{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}$ , define:

$$933 \quad f_\ell(\{(e_1, b_1), (e_2, b_2), (e_3, b_3)\}) = \{e_1, e_2, e_3\}.$$

934 ► **Definition B.9** ( $f_\ell^{-1}$ ). For an arbitrary subgraph  $S^{(1)}$  of  $G^{(1)}$  with at most  $m \leq 3$  edges, the  
 935 inverse function  $f_\ell^{-1} : \binom{E_1}{\leq 3} \mapsto 2^{\binom{E_\ell}{3}}$  takes  $E_S^{(1)}$  and outputs the set of all elements  $s \in \binom{E_S^{(\ell)}}{3}$   
 936 such that  $f_\ell(s) = E_S^{(1)}$ .

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

941 We first note that  $f_\ell$  is well-defined:

942 ► **Lemma B.10.**  $f_\ell$  is a function.

943 **Proof.** Note that  $f_\ell$  is properly defined. For any  $S \in \binom{E_\ell}{3}$ ,  $|f(S)| \leq 3$ , since it has to be the  
 944 case that any subset of 3 edges in  $E_\ell$  will map to at most three edges in  $E_1$ . All mappings  
 945 are in the required range. Then, since for any  $b \in \{0, \dots, \ell - 1\}$  the map  $(e, b) \mapsto e$  is a  
 946 function and has exactly one mapping, which implies that  $f_\ell$  is a function. ◀

947 We are now ready to prove the structural lemmas. Note that  $f_\ell$  maps subsets of three  
 948 edges in  $G^{(\ell)}$  to a subset of at most three edges in  $E_1$ . To prove the structural lemmas, we  
 949 will use the map  $f_\ell^{-1}$ . In particular, to count the number of occurrences of  $\mathfrak{A}$ ,  $\mathfrak{B}$ ,  $\mathfrak{C}$  in  
 950  $G^{(\ell)}$  we count for each  $S \in \binom{E_1}{\leq 3}$ , how many of  $\mathfrak{A}/\mathfrak{B}/\mathfrak{C}$  subgraphs appear in  $f_\ell^{-1}(S)$ .

### 951 B.8.1 Proof of Lemma B.3

952 **Proof.** For each subset  $E_S^{(1)} \in \binom{E_1}{\leq 3}$ , we count the number of 3-matchings in the 3-edge  
 953 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  $E_S^{(1)}$   
 954 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  $s \in$   
 955  $\{(e_1, 0), (e_1, 1), (e_2, 0), (e_2, 1), (e_3, 0), (e_3, 1)\}$  such that  $f_\ell(s) = \{e_1, e_2, e_3\}$ .

956 We do a case analysis based on the subgraph  $S^{(1)}$  induced by  $E_S^{(1)}$  (denoted  $E_S^{(1)} \equiv S^{(1)}$ ):

#### 957 ■ 3-matching ( $\mathfrak{C}$ )

958 When  $S^{(1)}$  is isomorphic to  $\mathfrak{C}$ , it is the case that edges in  $E_S^{(2)}$  are *not* disjoint only for the  
 959 pairs  $(e_i, 0), (e_i, 1)$  for  $i \in \{1, 2, 3\}$ . All choices for  $b_1, b_2, b_3 \in \{0, 1\}$ ,  $(e_1, b_1), (e_2, b_2), (e_3, b_3)$   
 960 will compose a 3-matching. One can see that we have a total of two possible choices for  $b_i$   
 961 for each edge  $e_i$  in  $G^{(1)}$  yielding  $2^3 = 8$  possible 3-matchings in  $f_2^{-1}(E_S^{(1)})$ .

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

963 For  $S^{(1)}$  isomorphic to  $\mathfrak{A}$  edges  $e_2, e_3$  form a 2-path with  $e_1$  being disjoint. This means  
 964 that  $(e_2, 0), (e_2, 1), (e_3, 0), (e_3, 1)$  form a 4-path while  $(e_1, 0), (e_1, 1)$  is its own disjoint 2-path.  
 965 We can only pick either  $(e_1, 0)$  or  $(e_1, 1)$  for  $f_2^{-1}(E_S^{(1)})$ , and then we need to pick a 2-matching  
 966 from  $e_2$  and  $e_3$ . Note that the four path allows there to be 3 possible 2 matchings, specifically,

$$967 \{(e_2, 0), (e_3, 0)\}, \{(e_2, 0), (e_3, 1)\}, \{(e_2, 1), (e_3, 1)\}.$$

968 Since these two selections can be made independently, there are  $2 \cdot 3 = 6$  *distinct*  
 969 3-matchings in  $f_2^{-1}(E_S^{(1)})$ .

#### 970 ■ 3-star ( $\mathfrak{B}$ )

971 When  $S^{(1)}$  is isomorphic to  $\mathfrak{B}$ , the inner edges  $(e_i, 1)$  of  $E_S^{(2)}$  are all connected, and the  
 972 outer edges  $(e_i, 0)$  are all disjoint. Note that for a valid 3 matching it must be the case that  
 973 at most one inner edge can be part of the set of disjoint edges. For the case of when exactly  
 974 one inner edge is chosen, there exist 3 possibilities, based on which inner edge is chosen.  
 975 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  
 976  $j' \neq i, j' \neq j$ . The remaining possible 3-matching occurs when all 3 outer edges are chosen.  
 977 Thus, there are four 3-matchings in  $f_2^{-1}(E_S^{(1)})$ .

#### 978 ■ 3-path ( $\mathfrak{B}$ )

979 When  $S^{(1)}$  is isomorphic to  $\mathfrak{I}$  it is the case that all edges beginning with  $e_1$  and ending with  $e_3$   
 980 are successively connected. This means that the edges of  $E_S^{(2)}$  form a 6-path. For a 3-matching  
 981 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  
 982  $j = i + 1$ . There are four such possibilities:  $\{(e_1, 0), (e_2, 0), (e_3, 0)\}$ ,  $\{(e_1, 0), (e_2, 0), (e_3, 1)\}$ ,  
 983  $\{(e_1, 0), (e_2, 1), (e_3, 1)\}$ ,  $\{(e_1, 1), (e_2, 1), (e_3, 1)\}$ , a total of four 3-matchings in  $f_2^{-1}(E_S^{(1)})$ .

984 ■ Triangle ( $\mathfrak{A}$ )

985 For  $S^{(1)}$  isomorphic to  $\mathfrak{A}$ , note that it is the case that the edges in  $E_S^{(2)}$  are connected in a  
 986 successive manner, but this time in a cycle, such that  $(e_1, 0)$  and  $(e_3, 1)$  are also connected.  
 987 While this is similar to the discussion of the three path above, the first and last edges are  
 988 not disjoint, since they are connected. This rules out both subsets of  $(e_1, 0), (e_2, 0), (e_3, 1)$   
 989 and  $(e_1, 0), (e_2, 1), (e_3, 1)$ , yielding two 3-matchings.

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

991 ■ 2-matching ( $\mathfrak{II}$ ), 2-path ( $\mathfrak{A}$ ), 1 edge ( $\mathfrak{I}$ )

992 When  $|E_S^{(1)}| = 2$ , we can only pick one from each of two pairs,  $\{(e_1, 0), (e_1, 1)\}$  and  
 993  $\{(e_2, 0), (e_2, 1)\}$ . This implies that a 3-matching cannot exist in  $f_2^{-1}(E_S^{(1)})$ . The same  
 994 argument holds for  $|E_S^{(1)}| = 1$ , where we can only pick one edge from the pair  $\{(e_1, 0), (e_1, 1)\}$ .  
 995 Trivially, no 3-matching exists in  $f_2^{-1}(E_S^{(1)})$ .

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

## 999 B.8.2 Proof of Lemma B.4

1000 **Proof.** The number of triangles in  $G^{(\ell)}$  for  $\ell \geq 2$  will always be 0 for the simple fact that all  
 1001 cycles in  $G^{(\ell)}$  will have at least six edges. ◀

## 1002 B.8.3 Proof of Lemma 3.8

1003 **Proof.** The proof consists of two parts. First we need to show that a vector  $\mathbf{b}$  satisfying the  
 1004 linear system exists and further can be computed in  $O(m)$  time. Second we need to show  
 1005 that  $\#(G, \mathfrak{A}), \#(G, \mathfrak{II})$  can indeed be computed in time  $O(1)$ .

1006 The lemma claims that for  $\mathbf{M} = \begin{pmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{pmatrix}$ ,  $\mathbf{x} = \begin{pmatrix} \#(G, \mathfrak{A}) \\ \#(G, \mathfrak{II}) \end{pmatrix}$

1007 satisfies the system  $\mathbf{M} \cdot \mathbf{x} = \mathbf{b}$ .

1008 To prove the first step, we use Lemma B.1 to derive the following equality (dropping the  
 1009 superscript and referring to  $G^{(1)}$  as  $G$ ):

$$1010 \quad \#(G, \mathfrak{I})p^2 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{II})p^4 + 6\#(G, \mathfrak{A})p^3 + 6\#(G, \mathfrak{A})p^4 \\ 1011 \quad + 6\#(G, \mathfrak{II})p^4 + 6\#(G, \mathfrak{I} \mathfrak{A})p^5 + 6\#(G, \mathfrak{III})p^6 = \tilde{Q}_G^3(p, \dots, p) \quad (19)$$

$$1012 \quad \#(G, \mathfrak{A}) + \#(G, \mathfrak{II})p + \#(G, \mathfrak{I} \mathfrak{A})p^2 + \#(G, \mathfrak{III})p^3 \\ 1013 \quad = \frac{\tilde{Q}_G^3(p, \dots, p)}{6p^3} - \frac{\#(G, \mathfrak{I})}{6p} - \#(G, \mathfrak{A}) - \#(G, \mathfrak{II})p - \#(G, \mathfrak{A})p \quad (20)$$

$$1014 \quad \#(G, \mathfrak{A})(1 - 3p) - \#(G, \mathfrak{III})(3p^2 - p^3) =$$

$$\begin{aligned}
1015 \quad & \frac{\tilde{Q}_G^3(p, \dots, p)}{6p^3} - \frac{\#(G, \mathfrak{I})}{6p} - \#(G, \mathfrak{A}) - \#(G, \mathfrak{II})p - \#(G, \mathfrak{AA})p \\
1016 \quad & - [\#(G, \mathfrak{I} \mathfrak{A})p^2 + 3\#(G, \mathfrak{III})p^2] - [\#(G, \mathfrak{II})p + 3\#(G, \mathfrak{A})p] \\
1017 \quad & \tag{21}
\end{aligned}$$

1018 Eq. (19) is the result of Lemma B.1. We obtain the remaining equations through standard  
1019 algebraic manipulations.

1020 Note that the LHS of Eq. (21) is indeed the product  $\mathbf{M}[1] \cdot \mathbf{x}[1]$ . Further note that this  
1021 product is equal to the RHS of Eq. (21), where every term is computable in  $O(m)$  time (by  
1022 equations (12)-(17)). We set  $\mathbf{b}[1]$  to the RHS of Eq. (21).

1023 We follow the same process in deriving an equality for  $G^{(2)}$ . Replacing occurrences of  
1024  $G$  with  $G^{(2)}$ , we obtain Eq. (21) for  $G^{(2)}$ . Substituting identities from Lemma B.3 and  
1025 Lemma B.4 we obtain

$$\begin{aligned}
1026 \quad & 0 - (8\#(G, \mathfrak{III}) + 6\#(G, \mathfrak{I} \mathfrak{A}) + 4\#(G, \mathfrak{AA}) + 4\#(G, \mathfrak{II}) + 2\#(G, \mathfrak{A})(3p^2 - p^3)) = \\
1027 \quad & \frac{\tilde{Q}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{I})}{6p} - \#(G^{(2)}, \mathfrak{A}) - \#(G^{(2)}, \mathfrak{II})p - \#(G^{(2)}, \mathfrak{AA})p \\
1028 \quad & - [\#(G^{(2)}, \mathfrak{I} \mathfrak{A})p^2 + 3\#(G^{(2)}, \mathfrak{III})p^2] - [\#(G^{(2)}, \mathfrak{II})p + 3\#(G^{(2)}, \mathfrak{A})p] \\
& \tag{22} \\
1029 \quad & (10\#(G, \mathfrak{A}) + 10G\mathfrak{III})(3p^2 - p^3) = \\
1030 \quad & \frac{\tilde{Q}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{I})}{6p} - \#(G^{(2)}, \mathfrak{A}) - \#(G^{(2)}, \mathfrak{II})p - \#(G^{(2)}, \mathfrak{AA})p \\
1031 \quad & - [\#(G^{(2)}, \mathfrak{II})p + 3\#(G^{(2)}, \mathfrak{A})p] - [\#(G^{(2)}, \mathfrak{I} \mathfrak{A})p^2 - 3\#(G^{(2)}, \mathfrak{III})p^2] \\
1032 \quad & + (4\#(G, \mathfrak{AA}) + [6\#(G, \mathfrak{I} \mathfrak{A}) + 18\#(G, \mathfrak{III})] + [4\#(G, \mathfrak{II}) + 12\#(G, \mathfrak{A})]) (3p^2 - p^3) \\
1033 \quad & \tag{23}
\end{aligned}$$

1034 As in the previous equality derivation for  $G$ , note that the LHS of Eq. (23) is the same as  
1035  $\mathbf{M}[2] \cdot \mathbf{x}[2]$ . The RHS of Eq. (23) has terms all computable (by equations (12)-(17)) in  $O(m)$   
1036 time. Setting  $\mathbf{b}[2]$  to the RHS then completes the proof of step 1.

1037 Note that if  $\mathbf{M}$  has full rank then one can compute  $\#(G, \mathfrak{A})$  and  $\#(G, \mathfrak{III})$  in  $O(1)$   
1038 using Gaussian elimination.

1039 To show that  $\mathbf{M}$  indeed has full rank, we will show that  $\text{Det}(\mathbf{M}) \neq 0$  for every  $p \in (0, 1)$ .

1040 Let  $\mathbf{M} =$

$$\begin{aligned}
1041 \quad & \begin{vmatrix} 1 - 3p & -(3p^2 - p^3) \\ 10(3p^2 - p^3) & 10(3p^2 - p^3) \end{vmatrix} = (1 - 3p) \cdot 10(3p^2 - p^3) + 10(3p^2 - p^3) \cdot (3p^2 - p^3) \\
1042 \quad & = 10(3p^2 - p^3) \cdot (1 - 3p + 3p^2 - p^3) = 10(3p^2 - p^3) \cdot (-p^3 + 3p^2 - 3p + 1) \\
1043 \quad & = 10p^2(3 - p) \cdot (1 - p)^3 \\
1044 \quad & \tag{24}
\end{aligned}$$

1045 From Eq. (24) it can easily be seen that the roots of  $\text{Det}(\mathbf{M})$  are 0, 1, and 3. Hence there  
1046 are no roots in  $(0, 1)$  and Lemma 3.8 follows.  $\blacktriangleleft$

## 1047 **C** Missing Details from Section 4

1048 In the following definitions and examples, we use the following polynomial as an example:

$$1049 \quad Q(X, Y) = 2X^2 + 3XY - 2Y^2. \tag{25}$$

1050 ► **Definition C.1** (Pure Expansion). *The pure expansion of a polynomial  $Q$  is formed by*  
 1051 *computing all product of sums occurring in  $Q$ , without combining like monomials. The pure*  
 1052 *expansion of  $Q$  generalizes Definition 2.2 by allowing monomials  $m_i = m_j$  for  $i \neq j$ .*

1053 Note that similar in spirit to Definition 2.6,  $\mathbf{E}(\mathcal{C})$  Definition 4.2 reduces all variable exponents  
 1054  $e > 1$  to  $e = 1$ .

1055 In the following, we abuse notation and write  $\mathbf{v}$  to denote the monomial obtained as the  
 1056 products of the variables in the set.

1057 ► **Example C.2** (Example for Definition 4.2). *Consider the factorized representation  $(X +$   
 1058  $2Y)(2X - Y)$  of the polynomial in Eq. (25). Its circuit  $\mathcal{C}$  is illustrated in Fig. 3b. The pure  
 1059 expansion of the product is  $2X^2 - XY + 4XY - 2Y^2$  and the  $\mathbf{E}(\mathcal{C})$  is  $[(X, 2), (XY, -1), (XY, 4), (Y, -2)]$ .*

1060  $\mathbf{E}(\mathcal{C})$  effectively<sup>13</sup> encodes the *reduced* form of  $\text{POLY}(\mathcal{C})$ , decoupling each monomial into a  
 1061 set of variables  $\mathbf{v}$  and a real coefficient  $c$ . However, unlike the constraint on the input to  
 1062 compute  $\tilde{Q}$ , the input circuit  $\mathcal{C}$  does not need to be in SMB/SOP form.

1063 ► **Example C.3** (Example for Definition 4.3). *Using the same factorization from Example C.2,*  
 1064  $\text{POLY}(|\mathcal{C}|) = (X + 2Y)(2X + Y) = 2X^2 + XY + 4XY + 2Y^2 = 2X^2 + 5XY + 2Y^2$ . *Note that*  
 1065 *this is not the same as the polynomial from Eq. (25).*

1066 ► **Definition C.4** (Evaluation). *Given a circuit  $\mathcal{C}$  and a valuation  $\mathbf{a} \in \mathbb{R}^n$ , we define the*  
 1067 *evaluation of  $\mathcal{C}$  on  $\mathbf{a}$  as  $\mathcal{C}(\mathbf{a}) = \text{POLY}(\mathcal{C})(\mathbf{a})$ .*

1068 ► **Definition C.5** (Subcircuit). *A subcircuit of a circuit  $\mathcal{C}$  is a circuit  $\mathcal{S}$  such that  $\mathcal{S}$  is a DAG*  
 1069 *subgraph of the DAG representing  $\mathcal{C}$ . The sink of  $\mathcal{S}$  has exactly one gate  $g$ .*

## 1070 C.1 Proof of Theorem 4.8

1071 In order to prove Theorem 4.8, we will need to argue the correctness of  $\text{APPROXIMATE}\tilde{Q}$ ,  
 1072 which relies on the correctness of auxiliary algorithms  $\text{ONEPASS}$  and  $\text{SAMPLEMONOMIAL}$ .

► **Lemma C.6.** *The  $\text{ONEPASS}$  function completes in time:*

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

1073  $\text{ONEPASS}$  guarantees two post-conditions: *First, for each subcircuit  $\mathcal{S}$  of  $\mathcal{C}$ , we have that*  
 1074  $\mathcal{S}.\text{partial}$  *is set to  $|\mathcal{S}|(1, \dots, 1)$ . Second, when  $\mathcal{S}.\text{type} = +$ ,  $\mathcal{S}.\text{Lweight} = \frac{|S_i|(1, \dots, 1)}{|S|(1, \dots, 1)}$  and*  
 1075 *likewise for  $\mathcal{S}.\text{Rweight}$ .*

1076 To prove correctness of Algorithm 1, we only use the following fact that follows from the  
 1077 above lemma: for the modified circuit  $(\mathcal{C}_{\text{mod}})$ ,  $\mathcal{C}_{\text{mod}}.\text{partial} = |\mathcal{C}|(1, \dots, 1)$ .

► **Lemma C.7.** *The function  $\text{SAMPLEMONOMIAL}$  completes in time*

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

1078 *where  $k = \text{DEG}(\mathcal{C})$ . The function returns every  $(\mathbf{v}, \text{sign}(c))$  for  $(\mathbf{v}, c) \in \mathbf{E}(\mathcal{C})$  with probability*  
 1079  $\frac{|c|}{|\mathcal{C}|(1, \dots, 1)}$ .

1080 With the above two lemmas, we are ready to argue the following result:

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

1081 ► **Theorem C.8.** For any  $\mathcal{C}$  with  $\text{DEG}(\text{poly}(|\mathcal{C}|)) = k$ , algorithm 1 outputs an estimate  $\text{acc}$   
 1082 of  $\tilde{Q}(p_1, \dots, p_n)$  such that

$$1083 \quad P \left( \left| \text{acc} - \tilde{Q}(p_1, \dots, p_n) \right| > \epsilon \cdot |\mathcal{C}|(1, \dots, 1) \right) \leq \delta,$$

1084 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.

1085 Before proving Theorem C.8, we use it to argue our main result, Theorem 4.8.

1086 **Proof.** Set  $\mathcal{E} = \text{APPROXIMATE}\tilde{Q}(\mathcal{C}, (p_1, \dots, p_n), \delta, \epsilon')$ , where

$$1087 \quad \epsilon' = \epsilon \cdot \frac{\tilde{Q}(p_1, \dots, p_n) \cdot (1 - \gamma)}{|\mathcal{C}|(1, \dots, 1)},$$

1088 which achieves the claimed accuracy bound on  $\mathcal{E}$  due to Theorem C.8.

1089 The claim on the runtime follows from Theorem C.8 since

$$1090 \quad \frac{1}{(\epsilon')^2} \cdot \log \left( \frac{1}{\delta} \right) = \frac{\log \frac{1}{\delta}}{\epsilon^2 \left( \frac{\tilde{Q}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)} \right)^2}$$

$$1091 \quad = \frac{\log \frac{1}{\delta} \cdot |\mathcal{C}|^2(1, \dots, 1)}{\epsilon^2 \cdot \tilde{Q}^2(p_1, \dots, p_n)},$$

1093 which completes the proof. ◀

1094 We now return to the proof of Theorem C.8:

## 1095 C.2 Proof of Theorem C.8

1096 **Proof.** Consider now the random variables  $Y_1, \dots, Y_n$ , where each  $Y_i$  is the value of  $Y_{\mathbf{i}}$  after  
 1097 Line 8 is executed. In particular, note that we have

$$1098 \quad Y_i = \mathbb{1}(v \bmod \mathcal{B} \neq 0) \cdot \prod_{X_i \in \text{VAR}(v)} p_i,$$

1099 where the indicator variable handles the check in Line 6 Then for random variable  $Y_i$ , it is  
 1100 the case that

$$1101 \quad \mathbb{E}[Y_i] = \sum_{(v, \mathbf{c}) \in \mathbb{E}(\mathcal{C})} \frac{\mathbb{1}(v \bmod \mathcal{B} \neq 0) \cdot c \cdot \prod_{X_i \in \text{VAR}(v)} p_i}{|\mathcal{C}|(1, \dots, 1)}$$

$$1102 \quad = \frac{\tilde{Q}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)},$$

1104 where in the first equality we use the fact that  $\text{sgn}_1 \cdot |\mathbf{c}| = \mathbf{c}$  and the second equality follows  
 1105 from Eq. (4) with  $X_i$  substituted by  $p_i$ .

1106 Let  $\bar{\mathbf{Y}} = \frac{1}{N} \sum_{i=1}^N Y_i$ . It is also true that

$$1107 \quad \mathbb{E}[\bar{\mathbf{Y}}] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[Y_i] = \frac{\tilde{Q}(p_1, \dots, p_n)}{|\mathcal{C}|(1, \dots, 1)}.$$

1108 Hoeffding's inequality states that if we know that each  $Y_i$  (which are all independent)  
 1109 always lie in the intervals  $[a_i, b_i]$ , then it is true that

$$1110 \quad P(|\bar{\mathbf{Y}} - \mathbb{E}[\bar{\mathbf{Y}}]| \geq \epsilon) \leq 2 \exp \left( - \frac{2N^2\epsilon^2}{\sum_{i=1}^N (b_i - a_i)^2} \right).$$

## 23:32 Bag PDB Queries

1111 Line 5 shows that  $\text{sgn}_i$  has a value in  $\{-1, 1\}$  that is multiplied with  $O(k)$   $p_i \in [0, 1]$ ,  
 1112 which implies the range for each  $Y_i$  is  $[-1, 1]$ . Using Hoeffding's inequality, we then get:

$$1113 \quad P(|\bar{\mathbf{Y}} - \mathbb{E}[\bar{\mathbf{Y}}]| \geq \epsilon) \leq 2 \exp\left(-\frac{2N^2\epsilon^2}{2^2N}\right) = 2 \exp\left(-\frac{N\epsilon^2}{2}\right) \leq \delta,$$

1114 where the last inequality follows from our choice of  $N$  in Line 2.

1115 For the claimed probability bound of  $P\left(|\text{acc} - \tilde{Q}(p_1, \dots, p_n)| > \epsilon \cdot |\mathcal{C}|(1, \dots, 1)\right) \leq \delta$ ,  
 1116 note that in the algorithm,  $\text{acc}$  is exactly  $\bar{\mathbf{Y}} \cdot |\mathcal{C}|(1, \dots, 1)$ . Multiplying the rest of the terms  
 1117 by the same factor yields the said bound.

1118 This concludes the proof for the first claim of theorem C.8. We prove the claim on the  
 1119 runtime next.

### 1120 Run-time Analysis

1121 The runtime of the algorithm is dominated by Line 3 (which by Lemma C.6 takes time  
 1122  $O\left(\text{SIZE}(\mathcal{C}) \cdot \overline{\mathcal{M}}\left(\log\left(|\mathcal{C}|^2(1, \dots, 1)\right), \log(\text{SIZE}(\mathcal{C}))\right)\right)$ ) and the  $N$  iterations of the loop in  
 1123 Line 4. Each iteration's run time is dominated by the call to Line 5 (which by Lemma C.7  
 1124 takes  $O\left(\log k \cdot k \cdot \text{DEPTH}(\mathcal{C}) \cdot \overline{\mathcal{M}}\left(\log\left(|\mathcal{C}|^2(1, \dots, 1)\right), \log(\text{SIZE}(\mathcal{C}))\right)\right)$ ) and Line 6, which  
 1125 by the subsequent argument takes  $O(k \log k)$  time. We sort the  $O(k)$  variables by their block  
 1126 IDs and then check if there is a duplicate block ID or not. Adding up all the times discussed  
 1127 here gives us the desired overall runtime. ◀

### 1128 C.3 Proof of Corollary 4.10

1129 **Proof.** The result follows by first noting that by definition of  $\gamma$ , we have

$$1130 \quad \tilde{Q}(1, \dots, 1) = (1 - \gamma) \cdot |\mathcal{C}|(1, \dots, 1).$$

1131 Further, since each  $p_i \geq p_0$  and  $Q(\mathbf{X})$  (and hence  $\tilde{Q}(\mathbf{X})$ ) has degree at most  $k$ , we have that

$$1132 \quad \tilde{Q}(1, \dots, 1) \geq p_0^k \cdot \tilde{Q}(1, \dots, 1).$$

1133 The above two inequalities implies  $\tilde{Q}(1, \dots, 1) \geq p_0^k \cdot (1 - \gamma) \cdot |\mathcal{C}|(1, \dots, 1)$ . Applying  
 1134 this bound in the runtime bound in Theorem 4.8 gives the first claimed runtime. The final  
 1135 runtime of  $O_k\left(\frac{1}{\epsilon^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta} \cdot \overline{\mathcal{M}}\left(\log\left(|\mathcal{C}|^2(1, \dots, 1)\right), \log(\text{SIZE}(\mathcal{C}))\right)\right)$  follows by noting  
 1136 that  $\text{DEPTH}(\mathcal{C}) \leq \text{SIZE}(\mathcal{C})$  and absorbing all factors that just depend on  $k$ . ◀

### 1137 C.4 Proof of Lemma 4.11

1138 We will prove Lemma 4.11 by considering the three cases separately. We start by considering  
 1139 the case when  $\mathcal{C}$  is a tree:

1140 ► **Lemma C.9.** *Let  $\mathcal{C}$  be a tree (i.e. the sub-circuits corresponding to two children of a node  
 1141 in  $\mathcal{C}$  are completely disjoint). Then we have*

$$1142 \quad |\mathcal{C}|(1, \dots, 1) \leq (\text{SIZE}(\mathcal{C}))^{\text{DEG}(\mathcal{C})+1}.$$

1143 **Proof.** For notational simplicity define  $N = \text{SIZE}(\mathcal{C})$  and  $k = \text{DEG}(\mathcal{C})$ . To prove this result,  
 1144 we prove by induction on  $\text{DEPTH}(\mathcal{C})$  that  $|\mathcal{C}|(1, \dots, 1) \leq N^{k+1}$ . For the base case, we  
 1145 have that  $\text{DEPTH}(\mathcal{C}) = 0$ , and there can only be one node which must contain a coefficient



1146 (or constant) of 1. In this case,  $|\mathcal{C}|(1, \dots, 1) = 1$ , and  $\text{SIZE}(\mathcal{C}) = 1$ , and it is true that  
 1147  $|\mathcal{C}|(1, \dots, 1) = 1 \leq N^{k+1} = 1^1 = 1$ .

1148 Assume for  $\ell > 0$  an arbitrary circuit  $\mathcal{C}$  of  $\text{DEPTH}(\mathcal{C}) \leq \ell$  that it is true that  $|\mathcal{C}|(1, \dots, 1) \leq$   
 1149  $N^{\text{deg}(\mathcal{C})+1}$ .

1150 For the inductive step we consider a circuit  $\mathcal{C}$  such that  $\text{DEPTH}(\mathcal{C}) = \ell + 1$ . The sink can  
 1151 only be either a  $\times$  or  $+$  gate. Consider when sink node is  $\times$ . Let  $k_L, k_R$  denote  $\text{DEG}(\mathcal{C}_L)$  and  
 1152  $\text{DEG}(\mathcal{C}_R)$  respectively. Then note that

$$\begin{aligned}
 1153 \quad |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) \cdot |\mathcal{C}_R|(1, \dots, 1) \\
 1154 &\leq (N-1)^{k_L+1} \cdot (N-1)^{k_R+1} \\
 1155 &= (N-1)^{k+1} \\
 1156 &\leq N^{k+1}.
 \end{aligned} \tag{26}$$

1158 In the above the first inequality follows from the inductive hypothesis (and the fact that the  
 1159 size of either subtree is at most  $N-1$ ) and Eq. (26) follows by nothing that for a  $\times$  gate we  
 1160 have  $k = k_L + k_R + 1$ .

1161 For the case when the sink gate is a  $+$  gate, then for  $N_L = \text{SIZE}(\mathcal{C}_L)$  and  $N_R = \text{SIZE}(\mathcal{C}_R)$   
 1162 we have

$$\begin{aligned}
 1163 \quad |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) + |\mathcal{C}_R|(1, \dots, 1) \\
 1164 &\leq N_L^{k+1} + N_R^{k+1} \\
 1165 &\leq (N-1)^{k+1} \\
 1166 &\leq N^{k+1}.
 \end{aligned} \tag{27}$$

1168 In the above, the first inequality follows from the inductive hypothesis (and the fact that  
 1169  $k_L, k_R \leq k$ ). Note that the RHS of this inequality is maximized when the base and exponent  
 1170 of one of the terms is maximized. The second inequality follows from this fact as well as the  
 1171 fact that since  $\mathcal{C}$  is a tree we have  $N_L + N_R = N-1$  and, lastly, the fact that  $k \geq 0$ . This  
 1172 completes the proof.

1173 The upper bound in Lemma 4.11 for the general case is a simple variant of the above  
 1174 proof (but we present a proof sketch of the bound below for completeness):

1175 **► Lemma C.10.** *Let  $\mathcal{C}$  be a (general) circuit. Then we have*

$$1176 \quad |\mathcal{C}|(1, \dots, 1) \leq 2^{2^{\text{deg}(\mathcal{C}) \cdot \text{SIZE}(\mathcal{C})}}.$$

1177 **Proof Sketch.** We use the same notation as in the proof of Lemma C.9. We will prove by  
 1178 induction on  $\text{DEPTH}(\mathcal{C})$  that  $|\mathcal{C}|(1, \dots, 1) \leq 2^{2^k \cdot N}$ . The base case argument is similar to that  
 1179 in the proof of Lemma C.9. In the inductive case we have that  $N_L, N_R \leq N-1$ .

1180 For the case when the sink node is  $\times$ , we get that

$$\begin{aligned}
 1181 \quad |\mathcal{C}|(1, \dots, 1) &= |\mathcal{C}_L|(1, \dots, 1) \times |\mathcal{C}_R|(1, \dots, 1) \\
 1182 &\leq 2^{2^{k_L} \cdot N_L} \times 2^{2^{k_R} \cdot N_R} \\
 1183 &\leq 2^{2 \cdot 2^{k-1} \cdot (N-1)} \\
 1184 &\leq 2^{2^k N}.
 \end{aligned}$$

1186 In the above the first inequality follows from inductive hypothesis while the second inequality  
 1187 follows from the fact that  $k_L, k_R \leq k-1$  and  $N_L, N_R \leq N-1$ .

1188 Now consider the case when the sink node is  $+$ , we get that

$$\begin{aligned}
1189 \quad |\mathbf{C}|(1, \dots, 1) &= |\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1) \\
1190 \quad &\leq 2^{2^{k_L} \cdot N_L} + 2^{2^{k_R} \cdot N_R} \\
1191 \quad &\leq 2 \cdot 2^{2^k(N-1)} \\
1192 \quad &\leq 2^{2^k N}. \\
1193
\end{aligned}$$

1194 In the above the first inequality follows from the inductive hypothesis while the second  
1195 inequality follows from the facts that  $k_L, k_R \leq k$  and  $N_L, N_R \leq N - 1$ . The final inequality  
1196 follows from the fact that  $k \geq 0$ .  $\blacktriangleleft$

1197 Finally, we consider the case when  $\mathbf{C}$  encodes the run of the algorithm from [24] on an  
1198 FAQ query. We cannot handle the full generality of an FAQ query but we can handle an FAQ  
1199 query that has a “core” join query on  $k$  relations and then a subset of the  $k$  attributes are  
1200 “summed” out (e.g. the sum could be because of projecting out a subset of attributes from  
1201 the join query). While the algorithm [24] essentially figures out when to ‘push in’ the sums,  
1202 in our case since we only care about  $|\mathbf{C}|(1, \dots, 1)$  we will consider the obvious circuit that  
1203 computes the “inner join” using a worst-case optimal join (WCOJ) algorithm like [27] and  
1204 then adding in the addition gates. The basic idea is very simple: we will argue that there  
1205 are at most  $\text{SIZE}(\mathbf{C})^k$  tuples in the join output (each with having a value of 1 in  $|\mathbf{C}|(1, \dots, 1)$ ).  
1206 Then the largest value we can see in  $|\mathbf{C}|(1, \dots, 1)$  is by summing up these at most  $\text{SIZE}(\mathbf{C})^k$   
1207 values of 1. Note that this immediately implies the claimed bound in Lemma 4.11.

1208 We now sketch the argument for the claim about the join query above. First, we note  
1209 that the computation of a WCOJ algorithm like [27] can be expressed as a circuit with  
1210 *multiple* sinks (one for each output tuple). Note that annotation corresponding to  $\mathbf{t}$  in  $\mathbf{C}$  is  
1211 the polynomial  $\prod_{e \in E} R(\pi_e(\mathbf{t}))$  (where  $E$  indexes the set of relations). It is easy to see that  
1212 in this case the value of  $\mathbf{t}$  in  $|\mathbf{C}|(1, \dots, 1)$  will be 1 (by multiplying 1  $k$  times). The claim  
1213 on the number of output tuples follow from the trivial bound of multiplying the input size  
1214 bound (each relation has at most  $n \leq \text{SIZE}(\mathbf{C})$  tuples and hence we get an overall bound of  
1215  $n^k \leq \text{SIZE}(\mathbf{C})^k$ . Note that we did not really use anything about the WCOJ algorithm except  
1216 for the fact that  $\mathbf{C}$  for the join part only is built only of multiplication gates. In fact, we do  
1217 not need the better WCOJ join size bounds either (since we used the trivial  $n^k$  bound). As  
1218 a final remark, we note that we can build the circuit for the join part by running say the  
1219 algorithm from [24] on an FAQ query that just has the join query but each tuple is annotated  
1220 with the corresponding variable  $X_i$  (i.e. the semi-ring for the FAQ query is  $\mathbb{N}[\mathbf{X}]$ ).

## 1221 C.5 OnePass Remarks

1222 Please note that it is *assumed* that the original call to ONEPASS consists of a call on an  
1223 input circuit  $\mathbf{C}$  such that the values of members `partial`, `Lweight` and `Rweight` have been  
1224 initialized to Null across all gates.

1225 The evaluation of  $|\mathbf{C}|(1, \dots, 1)$  can be defined recursively, as follows (where  $\mathbf{C}_L$  and  $\mathbf{C}_R$  are  
1226 the ‘left’ and ‘right’ inputs of  $\mathbf{C}$  if they exist):

$$1227 \quad |\mathbf{C}|(1, \dots, 1) = \begin{cases} |\mathbf{C}_L|(1, \dots, 1) \cdot |\mathbf{C}_R|(1, \dots, 1) & \text{if } \mathbf{C.type} = \times \\ |\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1) & \text{if } \mathbf{C.type} = + \\ |\mathbf{C.val}| & \text{if } \mathbf{C.type} = \text{NUM} \\ 1 & \text{if } \mathbf{C.type} = \text{VAR}. \end{cases} \quad (28)$$

1228  
1229

1230 It turns out that for proof of Lemma C.7, we need to argue that when `C.type = +`, we  
 1231 indeed have

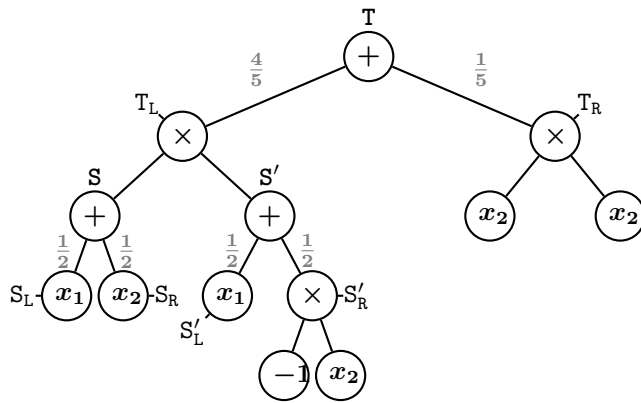
$$1232 \quad \text{C.Lweight} \leftarrow \frac{|\mathbf{C}_L|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)}; \quad (29)$$

$$1233 \quad \text{C.Rweight} \leftarrow \frac{|\mathbf{C}_R|(1, \dots, 1)}{|\mathbf{C}_L|(1, \dots, 1) + |\mathbf{C}_R|(1, \dots, 1)} \quad (30)$$

1235

### 1236 C.6 OnePass Example

1237 ► **Example C.11.** Let  $T$  encode the expression  $(X_1 + X_2)(X_1 - X_2) + X_2^2$ . After one pass,  
 1238 Algorithm 2 would have computed the following weight distribution. For the two inputs of the  
 1239 root  $+$  node  $T$ ,  $T.Lweight = \frac{4}{5}$  and  $T.Rweight = \frac{1}{5}$ . Similarly, let  $S$  denote the left-subtree of  
 1240  $T_L$ ,  $S.Lweight = S.Rweight = \frac{1}{2}$ . This is depicted in Fig. 4.



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

### 1241 C.7 OnePass

### 1242 C.8 Proof of Lemma C.6

1243 **Proof.** We prove the correct computation of `partial`, `Lweight`, `Rweight` values on  $C$  by  
 1244 induction over the number of iterations in line 2 over the topological order `TOPORD` of the  
 1245 input circuit  $C$ . Note that `TOPORD` is the standard definition of a topological ordering over  
 1246 the DAG structure of  $C$ .

1247 For the base case, we have only one gate, which by definition is a source gate and must be  
 1248 either `VAR` or `NUM`. In this case, as per Eq. (28), lines 4 and 6 correctly compute `C.partial`  
 1249 as 1 and `C.val` respectively.

1250 For the inductive hypothesis, assume that ONEPASS correctly computes `S.partial`,  
 1251 `S.Lweight`, and `S.Rweight` for all gates  $g$  in  $C$  with  $k \geq 0$  iterations over `TOPORD`.

1252 We now prove for  $k+1$  iterations that ONEPASS correctly computes the `partial`, `Lweight`,  
 1253 and `Rweight` values for each gate  $g_i$  in  $C$  for  $i \in [k+1]$ . Note that the  $g_{k+1}$  must be in the  
 1254 last ordering of all gates  $g_i$ . It is also the case that  $g_{k+1}$  has two inputs. Finally, note that  
 1255 for `SIZE(C) > 1`, if  $g_{k+1}$  is a leaf node, we are back to the base case. Otherwise  $g_{k+1}$  is an  
 1256 internal node  $g_s.type = +$  or  $g_s.type = \times$ .

---

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

---

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

```

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

```

---

1257 When  $g_{k+1}.\text{type} = +$ , then by line 10  $g_{k+1}.\text{partial} = g_{k+1_L}.\text{partial} + g_{k+1_R}.\text{partial}$ ,  
 1258 a correct computation, as per Eq. (28). Further, lines 11 and 12 compute  $g_{k+1}.L\text{weight} =$   
 1259  $\frac{g_{k+1_L}.\text{partial}}{g_{k+1}.\text{partial}}$  and analogously for  $g_{k+1}.R\text{weight}$ . Note that all values needed for each  
 1260 computation have been correctly computed by the inductive hypothesis.

1261 When  $g_{k+1}.\text{type} = \times$ , then line 8 computes  $g_{k+1}.\text{partial} = g_{k+1_L}.\text{partial} \times g_{k+1_R}.\text{partial}$ ,  
 1262 which indeed is correct, as per Eq. (28).

### 1263 Runtime Analysis

1264 It is known that  $\text{TOPORD}(G)$  is computable in linear time. Next, each of the  $\text{SIZE}(\mathcal{C})$   
 1265 iterations of the loop in Line 2 take  $O(\overline{\mathcal{M}}(\log(|\mathcal{C}(1 \dots, 1)|), \log \text{SIZE}(\mathcal{C})))$  time. It is easy  
 1266 to see that each of all the numbers which the algorithm computes is at most  $|\mathcal{C}|(1, \dots, 1)$ .  
 1267 Hence, by definition each such operation takes  $\overline{\mathcal{M}}(\log(|\mathcal{C}(1 \dots, 1)|), \log \text{SIZE}(\mathcal{C}))$  time, which  
 1268 proves the claimed runtime. ◀

## 1269 C.9 SampleMonomial Remarks

1270 We briefly describe the top-down traversal of `SAMPLEMONOMIAL`. For a parent  $+$  gate, the  
 1271 input to be visited is sampled from the weighted distribution precomputed by `ONEPASS`.  
 1272 When a parent  $\times$  node is visited, both inputs are visited. The algorithm computes two  
 1273 properties: the set of all variable leaf nodes visited, and the product of the signs of visited  
 1274 coefficient leaf nodes. We will assume the `TreeSet` data structure to maintain sets with  
 1275 logarithmic time insertion and linear time traversal of its elements. While we would like to  
 1276 take advantage of the space efficiency gained in using a circuit  $\mathcal{C}$  instead an expression tree  $T$ ,  
 1277 we do not know that such a method exists when computing a sample of the input polynomial  
 1278 representation.

---

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

---

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

```

1:  $\text{vars} \leftarrow \emptyset$ 
2: if  $\mathcal{C}.\text{type} = +$  then ▷ Sample at every + node
3:    $\mathcal{C}_{\text{samp}} \leftarrow$  Sample from left input ( $\mathcal{C}_L$ ) and right input ( $\mathcal{C}_R$ ) w.p.  $\mathcal{C}.\text{Lweight}$  and  $\mathcal{C}.\text{Rweight}$ . ▷ Each call to SAMPLEMONOMIAL uses fresh randomness
4:    $(\mathbf{v}, \mathbf{s}) \leftarrow$  SAMPLEMONOMIAL( $\mathcal{C}_{\text{samp}}$ )
5:   return  $(\mathbf{v}, \mathbf{s})$ 
6: else if  $\mathcal{C}.\text{type} = \times$  then ▷ 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  $\mathcal{C}.\text{type} = \text{numeric}$  then ▷ The leaf is a coefficient
15:   return  $(\{\}, \text{sign}(\mathcal{C}.\text{val}))$ 
16: else if  $\mathcal{C}.\text{type} = \text{var}$  then
17:   return  $(\{\mathcal{C}.\text{val}\}, 1)$ 
18: end if

```

---

1279 The efficiency gains of circuits over trees is found in the capability of circuits to only  
 1280 require space for each *distinct* term in the compressed representation. This saves space  
 1281 in such polynomials containing non-distinct terms multiplied or added to each other, e.g.,  
 1282  $x^4$ . However, to avoid biased sampling, it is imperative to sample from both inputs of a  
 1283 multiplication gate, independently, which is indeed the approach of SAMPLEMONOMIAL.

## 1284 C.10 Proof of Lemma C.7

1285 **Proof.** We first need to show that SAMPLEMONOMIAL indeed returns a monomial  $\mathbf{v}$ ,<sup>14</sup> such  
 1286 that  $(\mathbf{v}, \mathbf{c})$  is in  $E(\mathcal{C})$ , which we do by induction on the depth of  $\mathcal{C}$ .

1287 For the base case, let the depth  $d$  of  $\mathcal{C}$  be 0. We have that the root node is either a  
 1288 constant  $\mathbf{c}$  for which by line 15 we return  $\{\}$ , or we have that  $\mathcal{C}.\text{type} = \text{VAR}$  and  $\mathcal{C}.\text{val} = x$ ,  
 1289 and by line 17 we return  $\{x\}$ . Both cases sample a monomial, and the base case is proven.

1290 For the inductive hypothesis, assume that for  $d \leq k$  for some  $k \geq 0$ , that it is indeed the  
 1291 case that SAMPLEMONOMIAL returns a monomial.

1292 For the inductive step, let us take a circuit  $\mathcal{C}$  with  $d = k + 1$ . Note that each input has  
 1293 depth  $d \leq k$ , and by inductive hypothesis both of them return a valid monomial. Then the  
 1294 root can be either a  $+$  or  $\times$  node. For the case of a  $+$  root node, line 3 of SAMPLEMONOMIAL  
 1295 will choose one of the inputs of the root. By inductive hypothesis it is the case that a  
 1296 monomial in  $E(\mathcal{C})$  is being returned from either input. Then it follows that for the case of  $+$   
 1297 root node a valid monomial is returned by SAMPLEMONOMIAL. When the root is a  $\times$  node,

---

<sup>14</sup>Technically it returns  $\text{VAR}(\mathbf{v})$  but for less cumbersome notation we will refer to  $\text{VAR}(\mathbf{v})$  simply by  $\mathbf{v}$  in this proof.

1298 line 10 computes the set union of the monomials returned by the two inputs of the root, and  
 1299 it is trivial to see by Definition 4.2 that  $v$  is a valid monomial in some  $(v, c) \in E(C)$ .

1300 We will next prove by induction on the depth  $d$  of  $C$  that the  $(v, c) \in E(C)$  is the  $v$   
 1301 returned by `SAMPLEMONOMIAL` with a probability  $\frac{|c|}{|C|(1, \dots, 1)}$ .

1302 For the base case  $d = 0$ , by definition 2.10 we know that the root has to be either a  
 1303 coefficient or a variable. For either case, the probability of the value returned is 1 since there  
 1304 is only one value to sample from. When the root is a variable  $x$  the algorithm correctly returns  
 1305  $(\{x\}, 1)$ . When the root is a coefficient, `SAMPLEMONOMIAL` correctly returns  $(\{\}, \text{sign}(c_i))$ .

1306 For the inductive hypothesis, assume that for  $d \leq k$  and  $k \geq 0$  `SAMPLEMONOMIAL` indeed  
 1307 samples  $v$  in  $(v, c) \in E(C)$  with probability  $\frac{|c|}{|C|(1, \dots, 1)}$ .

1308 We prove now for  $d = k + 1$  the inductive step holds. It is the case that the root of  $C$  has  
 1309 up to two inputs  $C_L$  and  $C_R$ . Since  $C_L$  and  $C_R$  are both depth  $d \leq k$ , by inductive hypothesis,  
 1310 `SAMPLEMONOMIAL` will sample both monomials  $v_L$  in  $(v_L, c_L) \in E(C_L)$  and  $v_R$  in  $(v_R, c_R) \in E(C_R)$ ,  
 1311 from  $C_L$  and  $C_R$  with probability  $\frac{|c_L|}{|C_L|(1, \dots, 1)}$  and  $\frac{|c_R|}{|C_R|(1, \dots, 1)}$ .

1312 The root has to be either a  $+$  or  $\times$  node.

1313 Consider the case when the root is  $\times$ . Note that we are sampling a term from  $E(C)$ .  
 1314 Consider  $(v, c) \in E(C)$ , where  $v$  is the sampled monomial. Notice also that it is the case  
 1315 that  $v = v_L \times v_R$ , where  $v_L$  is coming from  $C_L$  and  $v_R$  from  $C_R$ . The probability that  
 1316 `SAMPLEMONOMIAL` ( $C_L$ ) returns  $v_L$  is  $\frac{|c_{v_L}|}{|C_L|(1, \dots, 1)}$  and  $\frac{|c_{v_R}|}{|C_R|(1, \dots, 1)}$  for  $v_R$ . Since both  $v_L$  and  
 1317  $v_R$  are sampled with independent randomness, the final probability for sample  $v$  is then  
 1318  $\frac{|c_{v_L}| \cdot |c_{v_R}|}{|C_L|(1, \dots, 1) \cdot |C_R|(1, \dots, 1)}$ . For  $(v, c) \in E(C)$ , it is indeed the case that  $|c| = |c_{v_L}| \cdot |c_{v_R}|$  and that  
 1319  $|C|(1, \dots, 1) = |C_L|(1, \dots, 1) \cdot |C_R|(1, \dots, 1)$ , and therefore  $v$  is sampled with correct probability  
 1320  $\frac{|c|}{|C|(1, \dots, 1)}$ .

1321 For the case when  $C.\text{val} = +$ , `SAMPLEMONOMIAL` will sample monomial  $v$  from one of  
 1322 its inputs. By inductive hypothesis we know that any  $v_L$  in  $E(C_L)$  and any  $v_R$  in  $E(C_R)$  will  
 1323 both be sampled with correct probability  $\frac{|c_{v_L}|}{|C_L|(1, \dots, 1)}$  and  $\frac{|c_{v_R}|}{|C_R|(1, \dots, 1)}$ , where either  $v_L$  or  $v_R$  will  
 1324 equal  $v$ , depending on whether  $C_L$  or  $C_R$  is sampled. Assume that  $v$  is sampled from  $C_L$ , and  
 1325 note that a symmetric argument holds for the case when  $v$  is sampled from  $C_R$ . Notice also  
 1326 that the probability of choosing  $C_L$  from  $C$  is  $\frac{|c_L|(1, \dots, 1)}{|C_L|(1, \dots, 1) + |C_R|(1, \dots, 1)}$  as computed by `ONEPASS`.  
 1327 Then, since `SAMPLEMONOMIAL` goes top-down, and each sampling choice is independent  
 1328 (which follows from the randomness in the root of  $C$  being independent from the randomness  
 1329 used in its subtrees), the probability for  $v$  to be sampled from  $C$  is equal to the product of  
 1330 the probability that  $C_L$  is sampled from  $C$  and  $v$  is sampled in  $C_L$ , and

$$\begin{aligned}
 1331 \quad & P(\text{SAMPLEMONOMIAL}(C) = v) = \\
 1332 \quad & P(\text{SAMPLEMONOMIAL}(C_L) = v) \cdot P(\text{SampledChild}(C) = C_L) \\
 1333 \quad & = \frac{|c_v|}{|C_L|(1, \dots, 1)} \cdot \frac{|C_L|(1, \dots, 1)}{|C_L|(1, \dots, 1) + |C_R|(1, \dots, 1)} \\
 1334 \quad & = \frac{|c_v|}{|C|(1, \dots, 1)}, \\
 1335
 \end{aligned}$$

1336 and we obtain the desired result.

### 1337 Run-time Analysis

1338 It is easy to check that except for lines 10 and 3, all lines take  $O(1)$  time. For Line 10, consider  
 1339 an execution of Line 10. We note that we will be adding a given set of variables to some set at  
 1340 most once: since the sum of the sizes of the sets at a given level is at most  $\text{DEG}(C)$ , each gate  
 1341 visited takes  $O(\log \text{DEG}(C))$ . For Line 3, note that we pick  $C_L$  with probability  $\frac{a}{a+b}$  where

1342  $a = \mathbf{C.Lweight}$  and  $b = \mathbf{C.Rweight}$ . We can implement this step by picking a random number  
 1343  $r \in [a+b]$  and then checking if  $r \leq a$ . It is easy to check that  $a+b \leq |\mathbf{C}|(1, \dots, 1)$ . This means  
 1344 we need to add and compare  $\log |\mathbf{C}|(1, \dots, 1)$ -bit numbers, which can certainly be done in  
 1345 time  $\overline{\mathcal{M}}(\log(|\mathbf{C}|(1, \dots, 1)), \log \text{SIZE}(\mathbf{C}))$  (note that this is an over-estimate). Denote  $\text{COST}(\mathbf{C})$   
 1346 (Eq. (31)) to be an upper bound of the number of nodes visited by `SAMPLEMONOMIAL`.  
 1347 Then the runtime is  $O(\text{COST}(\mathbf{C}) \cdot \log \text{DEG}(\mathbf{C}) \cdot \overline{\mathcal{M}}(\log(|\mathbf{C}|(1, \dots, 1)), \log \text{SIZE}(\mathbf{C})))$ .

1348 We now bound the number of recursive calls in `SAMPLEMONOMIAL` by  $O((\text{DEG}(\mathbf{C}) + 1) \cdot$   
 1349  $\text{DEPTH}(\mathbf{C}))$ , which by the above will prove the claimed runtime.

1350 Let  $\text{COST}(\cdot)$  be a function that models an upper bound on the number of gates that can  
 1351 be visited in the run of `SAMPLEMONOMIAL`. We define  $\text{COST}(\cdot)$  recursively as follows.

$$1352 \quad \text{COST}(\mathbf{C}) = \begin{cases} 1 + \text{COST}(\mathbf{C}_L) + \text{COST}(\mathbf{C}_R) & \text{if } \mathbf{C.type} = \times \\ 1 + \max(\text{COST}(\mathbf{C}_L), \text{COST}(\mathbf{C}_R)) & \text{if } \mathbf{C.type} = + \\ 1 & \text{otherwise} \end{cases} \quad (31)$$

1353 First note that the number of gates visited in `SAMPLEMONOMIAL` is  $\leq \text{COST}(\mathbf{C})$ . To  
 1354 show that Eq. (31) upper bounds the number of nodes visited by `SAMPLEMONOMIAL`, note  
 1355 that when `SAMPLEMONOMIAL` visits a gate such that  $\mathbf{C.type} = \times$ , line 8 visits each input  
 1356 of  $\mathbf{C}$ , as defined in (31). For the case when  $\mathbf{C.type} = +$ , line 3 visits exactly one of the  
 1357 input gates, which may or may not be the subcircuit with the maximum number of gates  
 1358 traversed, which makes  $\text{COST}(\cdot)$  an upperbound. Finally, it is trivial to see that when  $\mathbf{C.type}$   
 1359  $\in \{\text{VAR}, \text{NUM}\}$ , i.e., a source gate, that only one gate is visited.

1360 We prove the following inequality holds.

$$1361 \quad 2(\text{DEG}(\mathbf{C}) + 1) \cdot \text{DEPTH}(\mathbf{C}) + 1 \geq \text{COST}(\mathbf{C}) \quad (32)$$

1362 Note that Eq. (32) implies the claimed runtime. We prove Eq. (32) for the number of  
 1363 gates traversed in `SAMPLEMONOMIAL` using induction over  $\text{DEPTH}(\mathbf{C})$ . Recall how degree is  
 1364 defined in Definition 4.6.

1365 For the base case  $\text{DEG}(\mathbf{C}) = \text{DEPTH}(\mathbf{C}) = 0$ ,  $\text{COST}(\mathbf{C}) = 1$ , and it is trivial to see that the  
 1366 inequality  $2\text{DEG}(\mathbf{C}) \cdot \text{DEPTH}(\mathbf{C}) + 1 \geq \text{COST}(\mathbf{C})$  holds.

1367 For the inductive hypothesis, we assume the bound holds for any circuit where  $\ell \geq$   
 1368  $\text{DEPTH}(\mathbf{C}) \geq 0$ . Now consider the case when `SAMPLEMONOMIAL` has an arbitrary circuit  $\mathbf{C}$   
 1369 input with  $\text{DEPTH}(\mathbf{C}) = \ell + 1$ . By definition  $\mathbf{C.type} \in \{+, \times\}$ . Note that since  $\text{DEPTH}(\mathbf{C}) \geq 1$ ,  
 1370  $\mathbf{C}$  must have input(s). Further we know that by the inductive hypothesis the inputs  $\mathbf{C}_i$  for  
 1371  $i \in \{L, R\}$  of the sink gate  $\mathbf{C}$  uphold the bound

$$1372 \quad 2(\text{DEG}(\mathbf{C}_i) + 1) \cdot \text{DEPTH}(\mathbf{C}_i) + 1 \geq \text{COST}(\mathbf{C}_i). \quad (33)$$

1373 It is also true that  $\text{DEPTH}(\mathbf{C}_L) \leq \text{DEPTH}(\mathbf{C}) - 1$  and  $\text{DEPTH}(\mathbf{C}_R) \leq \text{DEPTH}(\mathbf{C}) - 1$ .

1374 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  
 1375  $\text{DEG}(\mathbf{C}) = \text{DEG}(\mathbf{C}_L) + \text{DEG}(\mathbf{C}_R) + 1$ . In either case it is true that  $\text{DEPTH}(\mathbf{C}) = \max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) +$   
 1376  $1$ .

1377 If  $\mathbf{C.type} = \times$ , then, substituting values, the following should hold,

$$1378 \quad 2(\text{DEG}(\mathbf{C}_L) + \text{DEG}(\mathbf{C}_R) + 2) \cdot (\max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) + 1) + 1 \\ 1379 \quad \geq 2(\text{DEG}(\mathbf{C}_L) + 1) \cdot \text{DEPTH}(\mathbf{C}_L) + 2(\text{DEG}(\mathbf{C}_R) + 1) \cdot \text{DEPTH}(\mathbf{C}_R) + 3 \quad (34)$$

$$1380 \quad \geq 1 + \text{COST}(\mathbf{C}_L) + \text{COST}(\mathbf{C}_R) = \text{COST}(\mathbf{C}). \quad (35)$$

## 23:40 Bag PDB Queries

1382 To prove (34), first, the LHS expands to,

$$1383 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) + 4 + 1 \quad (36)$$

1384 where  $\text{DEPTH}_{\max}$  is used to denote the maximum depth of the two input subcircuits. The  
1385 RHS expands to

$$1386 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_R) + 3 \quad (37)$$

1387 Putting Eq. (36) and Eq. (37) together we get

$$1388 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) + 5 \\ 1389 \quad \geq 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_R) + 3 \\ 1390 \quad (38)$$

1391 Since the following is always true,

$$1392 \quad 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 5 \\ 1393 \quad \geq 2\text{DEG}(\mathbf{C}_L) \cdot \text{DEPTH}(\mathbf{C}_L) + 2\text{DEG}(\mathbf{C}_R) \cdot \text{DEPTH}(\mathbf{C}_R) + 2\text{DEPTH}(\mathbf{C}_L) + 2\text{DEPTH}(\mathbf{C}_R) + 3,$$

1395 then it is the case that Eq. (38) is *always* true.

1396 Now to justify (35) which holds for the following reasons. First, the RHS is the result of  
1397 Eq. (31) when  $\mathbf{C.type} = \times$ . The LHS is then produced by substituting the upperbound of  
1398 (33) for each  $\text{COST}(\mathbf{C}_i)$ , trivially establishing the upper bound of (35). This proves Eq. (32)  
1399 for the  $\times$  case.

1400 For the case when  $\mathbf{C.type} = +$ , substituting values yields

$$1401 \quad 2(\max(\text{DEG}(\mathbf{C}_L), \text{DEG}(\mathbf{C}_R)) + 1) \cdot (\max(\text{DEPTH}(\mathbf{C}_L), \text{DEPTH}(\mathbf{C}_R)) + 1) + 1 \\ 1402 \quad \geq \max(2(\text{DEG}(\mathbf{C}_L) + 1) \cdot \text{DEPTH}(\mathbf{C}_L) + 1, 2(\text{DEG}(\mathbf{C}_R) + 1) \cdot \text{DEPTH}(\mathbf{C}_R) + 1) + 1 \quad (39)$$

$$1403 \quad \geq 1 + \max(\text{COST}(\mathbf{C}_L), \text{COST}(\mathbf{C}_R)) = \text{COST}(\mathbf{C}) \quad (40)$$

1405 To prove (39), the LHS expands to

$$1406 \quad 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEG}_{\max} + 2\text{DEPTH}_{\max} + 2 + 1. \quad (41)$$

1407 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  
1408 expanded RHS of (39):

$$1409 \quad 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEPTH}_{\max} + 2 \quad (42)$$

1410 Putting it together we obtain the following for (39):

$$1411 \quad 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEG}_{\max} + 2\text{DEPTH}_{\max} + 3 \\ 1412 \quad \geq 2\text{DEG}_{\max}\text{DEPTH}_{\max} + 2\text{DEPTH}_{\max} + 2, \quad (43)$$

1414 where it can be readily seen that the inequality stand and (43) follows. This proves (39).

1415 Similar to the case of  $\mathbf{C.type} = \times$ , (40) follows by equations (31) and (33).

1416 This proves (32) as desired. ◀



## 1417 C.11 Experimental Results

1418 Recall that by definition of *BIDB*, a query result cannot be derived by a self-join between  
 1419 non-identical tuples belonging to the same block. Note, that by Corollary 4.10,  $\gamma$  must be  
 1420 a constant in order for Algorithm 1 to achieve linear time. We would like to determine  
 1421 experimentally whether queries over *BIDB* instances in practice generate a constant number  
 1422 of cancellations or not. Such an experiment would ideally use a database instance with  
 1423 queries both considered to be typical representations of what is seen in practice.

1424 We ran our experiments using Windows 10 WSL Operating System with an Intel Core i7  
 1425 2.40GHz processor and 16GB RAM. All experiments used the PostgreSQL 13.0 database  
 1426 system.

1427 For the data we used the MayBMS data generator [1] tool to randomly generate uncertain  
 1428 versions of TPC-H tables. The queries computed over the database instance are  $Q_1$ ,  $Q_2$ , and  
 1429  $Q_3$  from [4], all of which are modified versions of TPC-H queries  $Q_3$ ,  $Q_6$ , and  $Q_7$  where all  
 1430 aggregations have been dropped.

1431 As written, the queries disallow *BIDB* cross terms. We first ran all queries, noting the  
 1432 result size for each. Next the queries were rewritten so as not to filter out the cross terms.  
 1433 The comparison of the sizes of both result sets should then suggest in one way or another  
 1434 whether or not there exist many cross terms in practice. As seen, the experimental query  
 1435 results contain little to no cancelling terms. Fig. 5 shows the result sizes of the queries,  
 1436 where column CF is the result size when all cross terms are filtered out, column CI shows  
 1437 the number of output tuples when the cancelled tuples are included in the result, and the  
 1438 last column is the value of  $\gamma$ . The experiments show  $\gamma$  to be in a range between  $[0, 0.1]\%$ ,  
 1439 indicating that only a negligible or constant (compare the result sizes of  $Q_1 < Q_2$  and their  
 1440 respective  $\gamma$  values) amount of tuples are cancelled in practice when running queries over a  
 1441 typical *BIDB* instance. Interestingly, only one of the three queries had tuples that violated  
 1442 the *BIDB* constraint.

1443 To conclude, the results in Fig. 5 show experimentally that  $\gamma$  is negligible in practice for  
 1444 *BIDB* queries. We also observe that (i) tuple presence is independent across blocks, so the  
 1445 corresponding probabilities (and hence  $p_0$ ) are independent of the number of blocks, and (ii)  
 1446 *BIDB*s model uncertain attributes, so block size (and hence  $\gamma$ ) is a function of the “messiness”  
 1447 of a dataset, rather than its size. Thus, we expect the corollary 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*.

## 1448 **D** Circuits

### 1449 D.1 Representing Polynomials with Circuits

#### 1450 D.1.1 Circuits for query plans

1451 We now formalize circuits and the construction of circuits for SPJU queries. As mentioned  
 1452 earlier, we represent lineage polynomials as arithmetic circuits over  $\mathbb{N}$ -valued variables with  $+$ ,  
 1453  $\times$ . A circuit for query  $Q$  and  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$  is a directed acyclic graph  $\langle V_{Q,\mathbf{D}}, E_{Q,\mathbf{D}}, \phi_{Q,\mathbf{D}}, \ell_{Q,\mathbf{D}} \rangle$

1454 with vertices  $V_{Q,\mathbf{D}}$  and directed edges  $E_{Q,\mathbf{D}} \subset V_{Q,\mathbf{D}}^2$ . The sink function  $\phi_{Q,\mathbf{D}} : \mathcal{U}^n \rightarrow V_{Q,\mathbf{D}}$   
 1455 is a partial function that maps the tuples of the  $n$ -ary relation  $Q(\mathbf{D})$  to vertices. We require  
 1456 that  $\phi_{Q,\mathbf{D}}$ 's range be limited to sink vertices (i.e., vertices with out-degree 0). A function  
 1457  $\ell_{Q,\mathbf{D}} : V_{Q,\mathbf{D}} \rightarrow \{+, \times\} \cup \mathbb{N} \cup \mathbf{X}$  assigns a label to each node: Source nodes (i.e., vertices  
 1458 with in-degree 0) are labeled with constants or variables (i.e.,  $\mathbb{N} \cup \mathbf{X}$ ), while the remaining  
 1459 nodes are labeled with the symbol  $+$  or  $\times$ . We require that vertices have an in-degree of  
 1460 at most two. For the specifics on how to construct a circuit to encode the polynomials of  
 1461 all result tuples for a query and  $\mathbb{N}[\mathbf{X}]$ -PDB see Appendix D.1. Note that we can construct  
 1462 circuits for BIDs in time linear in the time required for deterministic query processing over  
 1463 a possible world of the BID under the aforementioned assumption that  $|\mathbf{D}| \leq c \cdot |D|$ .

### 1464 D.1.2 Circuit size vs. runtime

1465 We now connect the size of a circuit (where the size of a circuit is the number of vertices  
 1466 in the corresponding DAG) for a given SPJU query  $Q$  and  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$  to its  $\mathbf{cost}(Q, D)$   
 1467 where  $D$  is one of the possible worlds of  $\mathbf{D}$ . We do this formally by showing that the size  
 1468 of the circuit is asymptotically no worse than the corresponding runtime of a large class of  
 1469 deterministic query processing algorithms.

1470 Each vertex  $v \in V_{Q,\mathbf{D}}$  in the arithmetic circuit for

$$1471 \quad \langle V_{Q,\mathbf{D}}, E_{Q,\mathbf{D}}, \phi_{Q,\mathbf{D}}, \ell_{Q,\mathbf{D}} \rangle$$

1472 encodes a polynomial, realized as

$$1473 \quad \mathbf{lin}(v) = \begin{cases} \sum_{v':(v',v) \in E_{Q,\mathbf{D}}} \mathbf{lin}(v') & \text{if } \ell(v) = + \\ \prod_{v':(v',v) \in E_{Q,\mathbf{D}}} \mathbf{lin}(v') & \text{if } \ell(v) = \times \\ \ell(v) & \text{otherwise} \end{cases}$$

1474 We define the circuit for a select-union-project-join  $Q$  recursively by cases as follows. In  
 1475 each case, let  $\langle V_{Q_i,\mathbf{D}}, E_{Q_i,\mathbf{D}}, \phi_{Q_i,\mathbf{D}}, \ell_{Q_i,\mathbf{D}} \rangle$  denote the circuit for subquery  $Q_i$ .

1476 **Base Relation.** Let  $Q$  be a base relation  $R$ . We define one node for each tuple. Formally,  
 1477 let  $V_{Q,\mathbf{D}} = \{v_t \mid t \in R\}$ , let  $\phi_{Q,\mathbf{D}}(t) = v_t$ , let  $\ell_{Q,\mathbf{D}}(v_t) = R(t)$ , and let  $E_{Q,\mathbf{D}} = \emptyset$ . This  
 1478 circuit has  $|R|$  vertices.

**Selection.** Let  $Q = \sigma_\theta(Q_1)$ . We re-use the circuit for  $Q_1$ . Formally, let  $V_{Q,\mathbf{D}} = V_{Q_1,\mathbf{D}}$ , let  
 $\ell_{Q,\mathbf{D}}(v_0) = 0$ , and let  $\ell_{Q,\mathbf{D}}(v) = \ell_{Q_1,\mathbf{D}}(v)$  for any  $v \in V_{Q_1,\mathbf{D}}$ . Let  $E_{Q,\mathbf{D}} = E_{Q_1,\mathbf{D}}$ , and define

$$\phi_{Q,\mathbf{D}}(t) = \phi_{Q_1,\mathbf{D}}(t) \text{ for } t \text{ s.t. } \theta(t).$$

1479 Dead sinks are iteratively removed, and so this circuit has at most  $|V_{Q_1,\mathbf{D}}|$  vertices.

**Projection.** Let  $Q = \pi_{\mathbf{A}}Q_1$ . We extend the circuit for  $Q_1$  with a new set of sum vertices  
 (i.e., vertices with label  $+$ ) for each tuple in  $Q$ , and connect them to the corresponding sink  
 nodes of the circuit for  $Q_1$ . Naively, let  $V_{Q,\mathbf{D}} = V_{Q_1,\mathbf{D}} \cup \{v_t \mid t \in \pi_{\mathbf{A}}Q_1\}$ , let  $\phi_{Q,\mathbf{D}}(t) = v_t$ ,  
 and let  $\ell_{Q,\mathbf{D}}(v_t) = +$ . Finally let

$$E_{Q,\mathbf{D}} = E_{Q_1,\mathbf{D}} \cup \{(\phi_{Q_1,\mathbf{D}}(t'), v_t) \mid t = \pi_{\mathbf{A}}t', t' \in Q_1, t \in \pi_{\mathbf{A}}Q_1\}$$

1480 This formulation will produce vertices with an in-degree greater than two, a problem that  
 1481 we correct by replacing every vertex with an in-degree over two by an equivalent fan-in tree.

1482 The resulting structure has at most  $|Q_1| - 1$  new vertices. The corrected circuit thus has at  
 1483 most  $|V_{Q_1, \mathbf{D}}| + |Q_1|$  vertices.

1484 **Union.** Let  $Q = Q_1 \cup Q_2$ . We merge graphs and produce a sum vertex for all tuples  
 1485 in both sides of the union. Formally, let  $V_{Q, \mathbf{D}} = V_{Q_1, \mathbf{D}} \cup V_{Q_2, \mathbf{D}} \cup \{v_t \mid t \in Q_1 \cap Q_2\}$ , let  
 1486  $\ell_{Q, \mathbf{D}}(v_t) = +$ , and let

$$1487 \quad E_{Q, \mathbf{D}} = E_{Q_1, \mathbf{D}} \cup E_{Q_2, \mathbf{D}} \cup \{(\phi_{Q_1, \mathbf{D}}(t), v_t), (\phi_{Q_2, \mathbf{D}}(t), v_t) \mid t \in Q_1 \cap Q_2\}$$

$$1488$$

$$1489 \quad \phi_{Q, \mathbf{D}}(t) = \begin{cases} v_t & \text{if } t \in Q_1 \cap Q_2 \\ \phi_{Q_1, \mathbf{D}}(t) & \text{if } t \notin Q_2 \\ \phi_{Q_2, \mathbf{D}}(t) & \text{if } t \notin Q_1 \end{cases}$$

1490 This circuit has  $|V_{Q_1, \mathbf{D}}| + |V_{Q_2, \mathbf{D}}| + |Q_1 \cap Q_2|$  vertices.

1491  **$k$ -ary Join.** Let  $Q = Q_1 \bowtie \dots \bowtie Q_k$ . We merge graphs and produce a multiplication  
 1492 vertex for all tuples resulting from the join. Naively, let  $V_{Q, \mathbf{D}} = V_{Q_1, \mathbf{D}} \cup \dots \cup V_{Q_k, \mathbf{D}} \cup$   
 1493  $\{v_t \mid t \in Q_1 \bowtie \dots \bowtie Q_k\}$ , let

$$1494 \quad E_{Q, \mathbf{D}} = E_{Q_1, \mathbf{D}} \cup \dots \cup E_{Q_k, \mathbf{D}} \cup \{(\phi_{Q_1, \mathbf{D}}(\pi_{sch(Q_1)}t), v_t),$$

$$1495 \quad \dots, (\phi_{Q_k, \mathbf{D}}(\pi_{sch(Q_k)}t), v_t) \mid t \in Q_1 \bowtie \dots \bowtie Q_k\}$$

$$1496$$

$$1497$$

1498 Let  $\ell_{Q, \mathbf{D}}(v_t) = \times$ , and let  $\phi_{Q, \mathbf{D}}(t) = v_t$ . As in projection, newly created vertices will have an  
 1499 in-degree of  $k$ , and a fan-in tree is required. There are  $|Q_1 \bowtie \dots \bowtie Q_k|$  such vertices, so the  
 1500 corrected circuit has  $|V_{Q_1, \mathbf{D}}| + \dots + |V_{Q_k, \mathbf{D}}| + (k - 1)|Q_1 \bowtie \dots \bowtie Q_k|$  vertices.

1501 **► Lemma D.1.** *Given a  $\mathbb{N}[\mathbf{X}]$ -PDB  $\mathbf{D}$  and query plan  $Q$ , the runtime of  $Q$  over  $\mathbf{D}$  has the*  
 1502 *same or better complexity as the size of the lineage of  $Q(\mathbf{D})$ . That is, we have  $|V_{Q, \mathbf{D}}| \leq$*   
 1503  *$(k - 1)\mathbf{cost}(Q)$ , where  $k$  is the maximal degree of any polynomial in  $Q(\mathbf{D})$ .*

1504 The proof is shown in Appendix D.2. We now have all the pieces to argue that using our  
 1505 approximation algorithm, the expected multiplicities of a SPJU query can be computed in  
 1506 essentially the same runtime as deterministic query processing for the same query.

## 1507 D.2 Proof for Lemma D.1

1508 **Proof.** Proof by induction. The base case is a base relation:  $Q = R$  and is trivially true  
 1509 since  $|V_{R, \mathbf{D}}| = |R|$ . For the inductive step, we assume that we have circuits for subplans  
 1510  $Q_1, \dots, Q_n$  such that  $|V_{Q_i, \mathbf{D}}| \leq (k_i - 1)\mathbf{cost}(Q_i, \mathbf{D})$  where  $k_i$  is the degree of  $Q_i$ .

1511 **Selection.** Assume that  $Q = \sigma_\theta(Q_1)$ . In the circuit for  $Q$ ,  $|V_{Q, \mathbf{D}}| = |V_{Q_1, \mathbf{D}}|$  vertices,  
 1512 so from the inductive assumption and  $\mathbf{cost}(Q, \mathbf{D}) = \mathbf{cost}(Q_1, \mathbf{D})$  by definition, we have  
 1513  $|V_{Q, \mathbf{D}}| \leq (k - 1)\mathbf{cost}(Q, \mathbf{D})$ . **Projection.** Assume that  $Q = \pi_{\mathbf{A}}(Q_1)$ . The circuit for  $Q$  has  
 1514 at most  $|V_{Q_1, \mathbf{D}}| + |Q_1|$  vertices.

$$1515 \quad |V_{Q, \mathbf{D}}| \leq |V_{Q_1, \mathbf{D}}| + |Q_1|$$

$$1516$$

1517 (From the inductive assumption)

$$1518 \quad \leq (k - 1)\mathbf{cost}(Q_1, \mathbf{D}) + |Q_1|$$

$$1519$$

1520 (By definition of  $\mathbf{cost}(Q, \mathbf{D})$ )

$$1521 \quad \leq (k - 1)\mathbf{cost}(Q, \mathbf{D}).$$

$$1522$$

## 23:44 Bag PDB Queries

1523 **Union.** Assume that  $Q = Q_1 \cup Q_2$ . The circuit for  $Q$  has  $|V_{Q_1, \mathbf{D}}| + |V_{Q_2, \mathbf{D}}| + |Q_1 \cap Q_2|$   
1524 vertices.

$$1525 \quad |V_{Q, \mathbf{D}}| \leq |V_{Q_1, \mathbf{D}}| + |V_{Q_2, \mathbf{D}}| + |Q_1| + |Q_2|$$

1527 (From the inductive assumption)

$$1528 \quad \leq (k-1)(\mathbf{cost}(Q_1, \mathbf{D}) + \mathbf{cost}(Q_2, \mathbf{D})) + (b_1 + b_2)$$

1530 (By definition of  $\mathbf{cost}(Q, \mathbf{D})$ )

$$1531 \quad \leq (k-1)(\mathbf{cost}(Q, \mathbf{D})).$$

1533  **$k$ -ary Join.** Assume that  $Q = Q_1 \bowtie \dots \bowtie Q_k$ . The circuit for  $Q$  has  $|V_{Q_1, \mathbf{D}}| + \dots +$   
1534  $|V_{Q_k, \mathbf{D}}| + (k-1)|Q_1 \bowtie \dots \bowtie Q_k|$  vertices.

$$1535 \quad |V_{Q, \mathbf{D}}| = |V_{Q_1, \mathbf{D}}| + \dots + |V_{Q_k, \mathbf{D}}| + (k-1)|Q_1 \bowtie \dots \bowtie Q_k|$$

1537 From the inductive assumption and noting  $\forall i : k_i \leq k-1$

$$1538 \quad \leq (k-1)\mathbf{cost}(Q_1, \mathbf{D}) + \dots + (k-1)\mathbf{cost}(Q_k, \mathbf{D}) +$$

$$1539 \quad (k-1)|Q_1 \bowtie \dots \bowtie Q_k|$$

$$1540 \quad \leq (k-1)(\mathbf{cost}(Q_1, \mathbf{D}) + \dots + \mathbf{cost}(Q_k, \mathbf{D}) +$$

$$1541 \quad |Q_1 \bowtie \dots \bowtie Q_k|)$$

1543 (By definition of  $\mathbf{cost}(Q, \mathbf{D})$ )

$$1544 \quad = (k-1)\mathbf{cost}(Q, \mathbf{D}).$$

1546 The property holds for all recursive queries, and the proof holds.  $\blacktriangleleft$

## 1547 **E** Parameterized Complexity

1548 In Sec. 3, we utilized common conjectures from fine-grained complexity theory. The notion of  
1549  $\#W[1]$  – *hard* is a standard notion in *parameterized complexity*, which by now is a standard  
1550 complexity tool in providing data complexity bounds on query processing results [18]. E.g.  
1551 the fact that  $k$ -matching is  $\#W[1]$  – *hard* implies that we cannot have an  $n^{\Omega(1)}$  runtime.  
1552 However, these results do not carefully track the exponent in the hardness result. E.g.  
1553  $\#W[1]$  – *hard* for the general  $k$ -matching problem does not imply anything specific for the  
1554 3-matching problem. Similar questions has led to intense research into the new sub-field  
1555 of *fine-grained complexity* (see [38]), where we care about the exponent in our hardness  
1556 assumptions as well– e.g. Conjecture 3.2 is based on the popular *Triangle detection hypothesis*  
1557 in this area (cf. [25]).