

OK, at this point it is probably easiest for me to draft this.

**Aaron says:** I am *unsure* of this footnote. @atri may need to word smith this one. I don't feel like I entirely understand the purpose of this footnote. E.g., we could have a query that runs deterministically in  $\Omega_k(n)$  worst case time; but this doesn't mean that  $T^*(Q, \mathcal{D})$  doesn't have a worst case lower bound of  $\Omega(n)^{c_0}$ , correct? We would replace  $T_{det}(Q, D, c)$  with  $\Omega_k(n)$ , no? If we replace  $T_{det}(Q, D, c)$  with  $n$ , then this doesn't accurately reflect the worst case lower bound for counting  $k$ -cliques in the first place.

I don't see why you say this.

100 already imply the claimed lower bounds if we were to replace the  $T_{det}(\text{OPT}(Q), D, c)$   
 101 by just  $n$  (indeed these results follow from known lower bounds for deterministic query  
 102 processing). Our contribution is to then identify a family of hard queries where deterministic  
 103 query processing is 'easy' but computing the expected multiplicities is hard.  
 104

105 **Our upper bound results.** We introduce a  $(1 \pm \epsilon)$ -approximation algorithm that  
 106 computes Problem 1.1 in time  $O_\epsilon(T_{det}(\text{OPT}(Q), D, c))$ . This means, when we are okay  
 107 with approximation, that we solve Problem 1.1 in time linear in the size of the deterministic  
 108 query and bag PDBs are deployable in practice. In contrast, known approximation techniques  
 109 ([40, 32]) in set-PDBs need time  $\Omega(T_{det}(\text{OPT}(Q), D, c)^{2k})$  (see Appendix G). Further, our  
 110 approximation algorithm works for a more general notion of bag PDBs beyond  $c$ -TIDBs (see  
 111 Sec. 2.2).

112 **1.1 Polynomial Equivalence**

113 A common encoding of probabilistic databases (e.g., in [30, 29, 5, 2] and many others)  
 114 relies on annotating tuples with lineages or propositional formulas that describe the set of  
 115 possible worlds that the tuple appears in. The bag semantics analog is a provenance/lineage  
 116 polynomial (see Fig. 1)  $\Phi[Q, D, t]$  [27], a polynomial with non-zero integer coefficients and  
 117 exponents, over integer variables  $\mathbf{X}$  encoding input tuple multiplicities.

X should be typeless. you can say when we abstract this poly X get replaced by int. values.

**Aaron says:** This seems confusing since I thought the goal was to have  $\mathbf{X}$  be abstract/typeless.

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

database  $D$  that counts the number of  $k$ -cliques, the results show a deterministic runtime of  $\Omega_k(n)$ , implying our lower bounds would hold.

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

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

122 ► **Problem 1.2** (Expected Multiplicity of Lineage Polynomials). Given an  $\mathcal{RA}^+$  query  $Q$ ,  
 123  $c$ -TIDB  $\mathcal{D}$  and result tuple  $t$ , compute the expected multiplicity of the polynomial  $\Phi[Q, D, t]$   
 124 (i.e.,  $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}} [\Phi[Q, D, t](\mathbf{W})]$ , where  $\mathbf{W} \in \{0, \dots, c\}^D$ ).

125 We note that computing Problem 1.1 is equivalent (yields the same result as) to computing  
 126 Problem 1.2 (see Proposition 2.8).

127 All of our results rely on working with a *reduced* form  $(\tilde{\Phi})$  of the lineage polynomial  $\Phi$ .  
 128 In fact, it turns out that for the 1-TIDB case, computing the expected multiplicity (over  
 129 bag query semantics) is *exactly* the same as evaluating this reduced polynomial over the  
 130 probabilities that define the 1-TIDB. This is also true when the query input(s) is a block  
 131 independent disjoint probabilistic database [40] (bag query semantics with tuple multiplicity  
 132 at most 1), for which the proof of Lemma 1.4 (introduced shortly) holds. Next, we motivate  
 133 this reduced polynomial. Consider the query  $Q_1$  defined as follows over the bag relations of  
 134 Fig. 2:

```
135 SELECT DISTINCT 1 FROM T t1, R r, T t2
136 WHERE t1.Point = r.Point1 AND t2.Point = r.Point2
137
```

It can be verified that  $\Phi(A, B, C, E, X, Y, Z)$  for the sole result tuple of  $Q_1$  is  $AXB + BYE + BZC$ . Now consider the product query  $Q_1^2 = Q_1 \times Q_1$ . The lineage polynomial for  $Q_1^2$  is given by  $\Phi_1^2(A, B, C, E, X, Y, Z)$

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

139 To compute  $\mathbb{E}[\Phi_1^2]$  we can use linearity of expectation and push the expectation through  
 140 each summand. To keep things simple, let us focus on the monomial  $\Phi_1^{(ABX)^2} = A^2 X^2 B^2$   
 141 as the procedure is the same for all other monomials of  $\Phi_1^2$ . Let  $W_X$  be the random  
 142 variable corresponding to a lineage variable  $X$ . Because the distinct variables in the  
 143 product are independent, we can push expectation through them yielding  $\mathbb{E}[W_A^2 W_X^2 W_B^2] =$   
 144  $\mathbb{E}[W_A^2] \mathbb{E}[W_X^2] \mathbb{E}[W_B^2]$ . Since  $W_A, W_B \in \{0, 1\}$  we can further derive  $\mathbb{E}[W_A] \mathbb{E}[W_X^2] \mathbb{E}[W_B]$   
 145 by the fact that for any  $W \in \{0, 1\}$ ,  $W^2 = W$ . Observe that if  $X \in \{0, 1\}$ , then we  
 146 further would have  $\mathbb{E}[W_A] \mathbb{E}[W_X] \mathbb{E}[W_B] = p_A \cdot p_X \cdot p_B$  (denoting  $Pr[W_A = 1] = p_A$ )  
 147  $= \tilde{\Phi}_1^{(ABX)^2}(p_A, p_X, p_B)$  (see *ii*) of Definition 1.3). However, in this example, we get stuck  
 148 with  $\mathbb{E}[W_X^2]$ , since  $W_X \in \{0, 1, 2\}$  and for  $W_X \leftarrow 2$ ,  $W_X^2 \neq W_X$ .

149 Denote the variables of  $\Phi$  to be  $\text{VARS}(\Phi)$ . In the  $c$ -TIDB setting,  $\Phi(\mathbf{X})$  has an equivalent  
 150 reformulation  $(\Phi_R(\mathbf{X}_R))$  that is of use to us, where  $|\mathbf{X}_R| = c \cdot |\mathbf{X}|$ . Given  $X_t \in \text{VARS}(\Phi)$ , by  
 151 definition  $X_t \in \{0, \dots, c\}$ . We can replace  $X_t$  by  $\sum_{j \in [c]} j X_{t,j}$  where the variables  $(X_{t,j})_{j \in [c]}$   
 152 are disjoint and each  $X_{t,j} \in \{0, 1\}$ . Then for any  $\mathbf{W} \in \{0, \dots, c\}^D$  and corresponding  
 153 reformulated world  $\mathbf{W}_R \in \{0, 1\}^{Dc}$ , we set  $\mathbf{W}_{R,t,j} = 1$  for  $\mathbf{W}_t = j$ , while  $\mathbf{W}_{R,t,j'} = 0$  for  
 154 all  $j' \neq j \in [c]$ . By construction then  $\Phi(\mathbf{X}) \equiv \Phi_R(\mathbf{X}_R)$  ( $\mathbf{X}_R = \text{VARS}(\Phi_R)$ ) since for any  
 155 valuation  $X_t \in [c]$  we have the equality  $X_t = j = \sum_{j \in [c]} j X_{t,j}$ .

**Aaron says:** I don't know the rules here, but since we have already (informally) defined  $\mathbf{X}$  to be variables of type integer encoding multiplicities (see todo note above) and thus worlds, it seems that it is fine and natural to refer to valuations of the variables themselves, without having to use  $\mathbf{W}$  necessarily. The point I am trying to get across in the last sentence is, given these semantics and domains, we have an equivalent polynomial. Or is it wrong to use  $\mathbf{X}$  and we should rather say, "for any  $\mathbf{W} \in \{0, \dots, c\}^D$ ,  $\mathbf{W}_R \in \{0, 1\}^{Dc}$  we have that  $\mathbf{W}_t = j = \sum_{j \in [c]} j \cdot \mathbf{W}_{R,t,j}$ ?"

156  
 157 Considering again our example,  
 158

Yeah the above discussion can be done just in terms of  $X$

$$\begin{aligned}
159 \quad \Phi_{1,R}^{(ABX)^2}(A, X, B) &= \Phi_1^{(AXB)^2} \left( \sum_{j_1 \in [c]} j_1 A_{j_1}, \sum_{j_2 \in [c]} j_2 X_{j_2}, \sum_{j_3 \in [c]} j_3 B_{j_3} \right) \\
160 &= \left( \sum_{j_1 \in [c]} j_1 A_{j_1} \right)^2 \left( \sum_{j_2 \in [c]} j_2 X_{j_2} \right)^2 \left( \sum_{j_3 \in [c]} j_3 B_{j_3} \right)^2.
\end{aligned}$$

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

166 ► **Definition 1.3.** For any polynomial  $\Phi((X_t)_{t \in D})$  define the reformulated polynomial  
167  $\Phi_R((X_{t,j})_{t \in D, j \in [c]})$  to be the polynomial  $\Phi_R = \Phi\left(\left(\sum_{j \in [c]} j \cdot X_{t,j}\right)_{t \in D}\right)$  and ii) define the  
168 reduced polynomial  $\tilde{\Phi}((X_{t,j})_{t \in D, j \in [c]})$  to be the polynomial resulting from converting  $\Phi_R$   
169 into the standard monomial basis (SMB),<sup>4</sup> removing all monomials containing the term  
170  $X_{t,j} X_{t,j'}$  for  $t \in D, j \neq j' \in [c]$ , and setting all variable exponents  $e > 1$  to 1.

171 Continuing with the example<sup>5</sup>  $\Phi_1^2(A, B, C, E, X_1, X_2, Y, Z)$  we have

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

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

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

## 182 1.2 Our Techniques

### 183 Lower Bound Proof Techniques.

**Aaron says:** Regarding what follows (in the next paragraph): I think this *may* be misleading (also, technically incorrect since  $\Phi$  is used instead of  $\tilde{\Phi}$ ) since it the lead  $c_{2k}$  of the term in  $\tilde{\Phi}(\mathbf{X})$  with  $2k$  distinct variables. However, technically, since we have that  $\tilde{\Phi}(\mathbf{p})$  is a univariate polynomial, then, indeed this IS an accurate statement, since the term with  $2k$  distinct variables in  $\tilde{\Phi}(\mathbf{p})$  is the term with the highest degree (this assumes for  $d$  distinct edges that  $d \geq k$  for our special graph query; otherwise, there is no  $k$ -matching, and the leading coefficient is not  $c_{2k}$ ). Perhaps we should note this. However, the context is in light of considering the univariate polynomial  $\tilde{\Phi}(\mathbf{p})$ . Perhaps change  $\Phi$  to  $\tilde{\Phi}(p, \dots, p)$ .

184

<sup>4</sup> This is the representation, typically used in set-PDBs, where the polynomial is represented as sum of 'pure' products. See Definition 2.1 for a formal definition.

<sup>5</sup> To save clutter we do not show the full expansion for variables with greatest multiplicity = 1 since e.g. for variable  $A$ , the sum of products itself evaluates to  $1^2 \cdot A^2 = A$ .

*already*  
*make sense. Looks like you did this?*

217 However, systems can directly emit compact, factorized representations of  $\Phi(\mathbf{X})$  (e.g.,  
 218 as a consequence of the standard projection push-down optimization [25]). For example,  
 219 in Fig. 2,  $B(Y + Z)$  is a factorized representation of the SMB-form  $BY + BZ$ . Accordingly,  
 220 this work uses (arithmetic) circuits<sup>6</sup> as the representation system of  $\Phi(\mathbf{X})$ .

221 Given that there exists a representation  $\mathbf{C}^*$  such that  $T_{LC}(Q, D, \mathbf{C}^*) \leq O(T_{det}(\text{OPT}(Q), D, c))$ ,  
 222 we can now focus on the complexity of the EC step. We can represent the factorized lineage  
 223 polynomial by its corresponding arithmetic circuit  $\mathbf{C}$  (whose size we denote by  $|\mathbf{C}|$ ). As  
 224 we also show in Appendix E.2.2, this size is also bounded by  $T_{det}(\text{OPT}(Q), D, c)$  (i.e.,  
 225  $|\mathbf{C}^*| \leq O(T_{det}(\text{OPT}(Q), D, c))$ ). Thus, the question of approximation can be stated as the  
 226 following stronger (since Problem 1.5 has access to *all* equivalent  $\mathbf{C}$  representing  $Q(\mathbf{W})(t)$ ),  
 227 but sufficient condition:

228 ► **Problem 1.6.** *Given one circuit  $\mathbf{C}$  that encodes  $\Phi[Q, D, t]$  for all result tuples  $t$  (one sink  
 229 per  $t$ ) for  $c$ -TIDB  $\mathcal{D}$  and  $\mathcal{RA}^+$  query  $Q$ , does there exist an algorithm that computes a  
 230  $(1 \pm \epsilon)$ -approximation of  $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[Q(\mathbf{W})(t)]$  (for all result tuples  $t$ ) in  $O(|\mathbf{C}|)$  time?*

231 For an upper bound on approximating the expected count, it is easy to check that if all the  
 232 probabilities are constant then  $\Phi(p_1, \dots, p_n)$  (i.e. evaluating the original lineage polynomial  
 233 over the probability values) is a constant factor approximation. For example, using  $Q_1^2$  from  
 234 above (with  $c = 1$ ) and  $p_A$  to denote  $Pr[A = 1]$ , we can see that

**Aaron says:** I changed Problem 1.6 to use  $c$ -TIDB. Correct me if this is wrong.  
 Our results do apply to a more general class of bag-PDB, but the main data model  
 considered in this paper is  $c$ -TIDB.  
 Also, for the example above and worked out in what follows, it might be better flow to  
 keep  $c = 2$  and change what is below.

Sounds good

$$236 \quad \Phi_1^2(\mathbf{p}) = p_A^2 p_X^2 p_B^2 + p_B^2 p_Y^2 p_E^2 + p_B^2 p_Z^2 p_C^2 + 2p_A p_X p_B^2 p_Y p_E + 2p_A p_X p_B^2 p_Z p_C + 2p_B^2 p_Y p_E p_Z p_C$$

$$238 \quad \leq p_A p_X p_B + p_B p_Y p_E + p_B p_Z p_C + 2p_A p_X p_B p_Y p_E + 2p_A p_X p_B p_Z p_C + 2p_B p_Y p_E p_Z p_C = \tilde{\Phi}_1^2(\mathbf{p})$$

239 If we assume that all seven probability values are at least  $p_0 > 0$ , we get that  $\Phi_1^2(\mathbf{p})$  is in the  
 240 range  $[(p_0)^3 \cdot \tilde{\Phi}_1^2(\mathbf{p}), \tilde{\Phi}_1^2(\mathbf{p})]$ , which is *not a tight approximation*. In sec. 4 we demonstrate  
 241 that a  $(1 \pm \epsilon)$  (multiplicative) approximation with competitive performance is achievable.  
 242 To get an  $(1 \pm \epsilon)$ -multiplicative approximation and solve Problem 1.6, using  $\mathbf{C}$  we uniformly  
 243 sample monomials from the equivalent SMB representation of  $\Phi$  (without materializing the  
 244 SMB representation) and ‘adjust’ their contribution to  $\tilde{\Phi}(\cdot)$ .

245 **Applications.** Recent work in heuristic data cleaning [51, 45, 42, 8, 45] emits a PDB when  
 246 insufficient data exists to select the ‘correct’ data repair. Probabilistic data cleaning is a  
 247 crucial innovation, as the alternative is to arbitrarily select one repair and ‘hope’ that queries  
 248 receive meaningful results. Although PDB queries instead convey the trustworthiness of  
 249 results [37], they are impractically slow [19, 18], even in approximation (see Appendix G).  
 250 Bags, as we consider, are sufficient for production use, where bag-relational algebra is already  
 251 the default for performance reasons. Our results show that bag-PDBs can be competitive,  
 252 laying the groundwork for probabilistic functionality in production database engines.

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

<sup>6</sup> An arithmetic circuit is a DAG with variable and/or numeric source nodes and internal, each nodes representing either an addition or multiplication operator.



This setup is not quite right. First of all Prop 2.4 has nothing to do with a circuit. The statement should go as follows:

(i) Let  $C$  be a  $c$ -TIDB Binary-BIDB circuit.  
 (ii) Given the above, the algorithm is a sampling based algorithm for the above sum: we sample (via SAMPLEMONOMIAL)  $(v, c) \in E(C)$  with probability proportional to  $|c|$  and compute  $Y = \mathbb{1}_{\text{isIND}(v_m)} \cdot \prod_{X_i \in v} p_i$ . Repeating the sampling an appropriate number of times and computing the average of  $Y$  gives us our final estimate. ONEPASS is used to compute the sampling probabilities needed in SAMPLEMONOMIAL (details are in Appendix D).  
 (iii) state how the depth & size changes

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

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

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

can be computed in time

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

In particular, if  $p_0 > 0$  and  $\gamma < 1$  are absolute constants then the above runtime simplifies to

$$O_k \left( \left( \frac{1}{(\epsilon')^2} \cdot \text{SIZE}(C) \cdot \log \frac{1}{\delta} \right) \cdot \bar{M}(\log(|C|(1, \dots, 1)), \log(\text{SIZE}(C))) \right).$$

The restriction on  $\gamma$  is satisfied by any 1-TIDB (where  $\gamma = 0$  in the equivalent 1-BIDB of Proposition 2.4) as well as for all three queries of the PDBench BIDB benchmark (see Appendix D.10 for experimental results). Further, we can also argue the following result:

**Lemma 4.8.** Given Binary-BIDB computed from the reduction of Proposition 2.4,  $\gamma(C) \leq 1 - (c+1)^{-(k-1)}$ .

**Proof of Lemma 4.8** Let  $\mathcal{D}' = (\times_{t \in \mathcal{D}'} \{0, c_t\}, \mathcal{D}')$  be the reduced Binary-BIDB and  $\mathcal{D} = (\{0, \dots, c\}^{\mathcal{D}}, \mathcal{D})$  the original  $c$ -TIDB.

By Proposition 2.4,  $\mathcal{D}'$  is a Binary-BIDB. By Definition 2.3, a block  $B_t$  of  $\mathcal{D}'$  has the property that  $\sum_{t \in \mathcal{D}', j \in [c]} p_{t,j} \leq 1$ . Then, if we consider the case of strict inequality, we have an extra possible outcome in block  $B_t$ , the outcome when no tuple is present in a possible world. Let us denote this as  $t_0$ . Then there are at most  $c+1$  disjoint tuples in  $B_t$ . We argue later that the case when  $t_0$  is a possibility produces the worst case  $\gamma$ .

Let  $\Phi'(\mathbf{X})$  be an arbitrary polynomial produced by  $Q(\mathcal{D}')$  with  $\mathbf{X} = (X_{t,j})_{t \in \mathcal{D}', j \in [0,c]}$  the set of variables in  $\mathcal{D}'$ . Let  $m$  be an arbitrary monomial in  $\Phi'(\mathbf{X})$  and  $v_m$  be the set of variables appearing in  $m$ . We define a cross term to be any monomial  $m$  such that there exists  $j \neq j' \in [0, c]$  such that  $X_{t,j}, X_{t,j'} \in v_m$ .

The semantics of Fig. 3 show that a new monomial product can only be generated by the  $\bowtie$  operator of  $\mathcal{RA}^+$  queries. Further, a cross term may only be produced specifically when the join is a self join. The highest number of terms that can be produced by a self join of  $B_t$  is  $(c+1)^k$ , the case for when all tuples join and  $\sum_{t \in \mathcal{D}', j \in [c]} p_{t,c} < 1$  as noted above. For monomials  $m \in \{ \times_{i \in [k], j \in [0,c]} X_{t,j_i} \}$ , there exist exactly  $(c+1)$  non-cross terms, specifically  $X_{t,j}^k$  for  $j \in [0, c]$ . Then there are exactly  $(c+1)^k - (c+1)$  cross terms (cancellations). This implies that  $\gamma(C) = 1 - \frac{(c+1)}{(c+1)^k}$  for this case.

We now show that the case above is indeed the worst case. First, given a self join, it is always the case that  $X_{t,j}^k$  will be in the output since all tuples join with themselves. Then, the most number of cancellations occurs when we have that all  $X_{t,j}$  joins with all  $X_{t,j'}$  for

Specifically state how  $C'$  is built for  $C$  i.e. replace the input gate  $X_t$  by  $\sum_{j=1}^c X_{t,j}$

Argument just in terms of  $C \rightarrow C'$ .  
 Do not bring in  $C \& C'$  terms.

Needs to be re-written as above. Structure as follows:  
 (1) Show the conversion for  $C$  to  $C'$  as above  
 (2) Then argue that

the  $C \rightarrow C'$  conversion follows from redux in Def 2.3 & hence by Prop 2.4,  $C'$  is a valid circuit. (state that the corresponding linear polys have the same exp. num)  
 (3) Rest of the should be just  $C, C', \text{POLY}(C) \& \text{POLY}(C')$  in terms of  $\text{POLY}(C) \& \text{POLY}(C')$

Here's a better way to do the monomial argument:

- (1) First consider  $E(C)$  i.e. set of all expanded monomials for  $c$ -TIDB det  $C$
- (2) Fix a bit monomial  $m$  in  $C \Rightarrow$  then argue

23:18 Bag PDB Queries

578  $j \neq j' \in [0, c]$ . Finally, it is the case that  $c^k - c \leq (c+1)^k - (c+1) = \sum_{i=1}^k \binom{k}{i} c^i - (c-1)$   
579 for  $c, k \in \mathbb{N}$ , which implies that the worst case is when we have the 'extra' tuple  $t_0$  and all  
580 tuples joining, which is exactly the case above, producing the greatest  $\gamma(C)$  ratio.

581 Since the size of any block  $B$  is  $c+1$ , it follows that  $\gamma(C)$  ratio for block  $B_t$  is the same  
582 when taken across all blocks of  $Q(D')$ , since the number of blocks  $n$  cancels out of the ratio  
583 calculations.

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

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

591 **Lemma 4.9.** For any Binary-BIDB circuit  $C$  with  $DEG(C) = k$ , we have  $|C|(1, \dots, 1) \leq$   
592  $2^{2^k \cdot DEPTH(C)}$ . Further, if  $C$  is a tree, then we have  $|C|(1, \dots, 1) \leq SIZE(C)^{O(k)}$ .

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

597 Finally, note that by Proposition E.1 and Lemma E.2 for any  $\mathcal{RA}^+$  query  $Q$ , there exists a  
598 circuit  $C^*$  for  $\Phi[Q, D, t]$  such that  $DEPTH(C^*) \leq O_{|Q|}(\log n)$  and  $SIZE(C) \leq O_k(T_{det}(Q, D_\Omega))$ .  
599 Using this along with Lemma 4.9, Theorem 4.7 and the fact that  $n \leq T_{det}(Q, D_\Omega)$ , we have  
600 the following corollary:

601 **Corollary 4.10.** Let  $Q$  be an  $\mathcal{RA}^+$  query and  $D$  be a Binary-BIDB with  $p_0 > 0$  and  $\gamma < 1$   
602 (where  $p_0, \gamma$  as in Theorem 4.7) are absolute constants. Let  $\Phi(\mathbf{X}) = \Phi[Q, D, t]$  for any result  
603 tuple  $t$  with  $\deg(\Phi) = k$ . Then one can compute an approximation satisfying Eq. (3) in time  
604  $O_{k, |Q|, \epsilon', \delta}(T_{det}(OPT(Q), D, c))$  (given  $Q, D$  and  $p_i$  for each  $i \in [n]$  that defines  $\mathcal{P}$ ).

605 Next, we note that the above result along with Lemma 4.8 answers Problem 1.5 in the  
606 affirmative as follows:

607 **Corollary 4.11.** Let  $Q$  be an  $\mathcal{RA}^+$  query and  $D$  be a  $c$ -TIDB with  $p_0 > 0$  (where  $p_0$   
608 as in Theorem 4.7) is an absolute constant. Let  $\Phi(\mathbf{X}) = \Phi[Q, D, t]$  for any result tuple  
609  $t$  with  $\deg(\Phi) = k$ . Then one can compute an approximation satisfying Eq. (3) in time  
610  $O_{k, |Q|, \epsilon', \delta, c}(T_{det}(OPT(Q), D, c))$  (given  $Q, D$  and  $p_{t,j}$  for each  $t \in D, j \in [c]$  that defines  
611  $\mathcal{P}$ ).

612 **Proof of Corollary 4.11.** The proof follows by Proposition 2.4, Lemma 4.8, and Corollary 4.10.

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

### 5 Related Work

618 **Probabilistic Databases (PDBs)** have been studied predominantly for set semantics.  
619 Approaches for probabilistic query processing (i.e., computing marginal probabilities of

(3) Observe that cancellations can only happen for monomials for each  $x_{t_i}$ .  
Then do the  
(3-1) Argue  
(3-2) that then overall for the

each argve monomial  $x_{t_1} \dots x_{t_k}$  gives us blah cot that if we can show the cancellations among this batch is  $\leq \gamma_0$ , then same bound of  $\gamma_0$  holds overall.

$$\gamma \geq \frac{1 - (c-1)}{1 - T(c+1)} \Rightarrow \frac{1 - (c-1)}{1 - (c+1)}$$

Why do you need this?

Also this goes to appendix?