

Computing expected multiplicities for bag-TIDBs with bounded multiplicities

Su Feng

Boris Glavic

sfeng14@hawk.iit.edu

bglavic@iit.edu

Illinois Institute of Technology, USA

Chicago, New York, USA

Aaron Huber

Oliver Kennedy

Atri Rudra

ahuber@buffalo.edu

okennedy@buffalo.edu

atri@buffalo.edu

University at Buffalo, USA

Buffalo, New York, USA

ABSTRACT

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

CCS CONCEPTS

• Information systems → Information systems applications; Data management systems; • Theory of computation → Probabilistic computation; Complexity classes.

KEYWORDS

probabilistic data model, parameterized complexity, fine-grained complexity, lineage polynomial

ACM Reference Format:

Su Feng, Boris Glavic, Aaron Huber, Oliver Kennedy, and Atri Rudra. 2022. Computing expected multiplicities for bag-TIDBs with bounded multiplicities. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 35 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

This work explores the problem of computing the expectation of the multiplicity of a tuple in the result of a query over a c -TIDB, a type of probabilistic database with bag semantics where the multiplicity of a tuple is a random variable with range $[0, c]$ for some fixed constant c and multiplicities assigned to any two tuples are independent of each other. Formally, a c -TIDB, $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ consists of a set of tuples D and a probability distribution \mathcal{P} over all possible worlds generated by assigning each tuple $t \in D$ a multiplicity in the range $[0, c]$. Any such world can be encoded as a vector from $\{0, \dots, c\}^D$, the set of all vectors of length $n = |D|$ such that each index corresponds to a distinct $t \in D$ storing its multiplicity. A given world $\mathbf{W} \in \{0, \dots, c\}^D$ can be interpreted as follows: for each $t \in D$, \mathbf{W}_t is the multiplicity of t in \mathbf{W} . Given that the multiplicities of tuples are independent events, the probability distribution \mathcal{P} can be expressed compactly by assigning each tuple a (disjoint) probability distribution over $[0, c]$. Let $p_{t,j}$ denote the probability that tuple t is assigned multiplicity j . The probability of a particular world \mathbf{W} is then $\prod_{t \in D} p_{t, \mathbf{W}_t}$.

Allowing for $\leq c$ multiplicities across all tuples gives rise to having $\leq (c + 1)^n$ possible worlds instead of the usual 2^n possible worlds of a 1-TIDB, which (assuming set query semantics), is the same as the traditional set TIDB. In this work, since we are generally considering bag query input, we will only be considering bag query semantics. We denote by $Q(\mathbf{W})(t)$ the multiplicity of t in query Q over possible world $\mathbf{W} \in \{0, \dots, c\}^D$.

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

PROBLEM 1.1. *Given a c -TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, \mathcal{RA}^+ query¹ Q , and result tuple t , compute the expected multiplicity of t : $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}} [Q(\mathbf{W})(t)]$.*

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2022 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

¹An \mathcal{RA}^+ query is a query expressed in positive relational algebra, i.e., using only the relational algebra operators selection (σ), projection (π), natural join (\bowtie) and union (\cup).

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

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.

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

Hardness of Set Query Semantics and Bag Query Semantics. Set query evaluation semantics over 1-TIDBs have been studied extensively, and the data complexity of the problem in general has been shown by Dalvi and Suicu to be #P-hard [13]. For our setting, there exists a trivial polytime algorithm to compute Problem 1.1 for any \mathcal{RA}^+ query over a c -TIDB due to linearity of expectation (see Sec. 1.1). Since we can compute Problem 1.1 in polynomial time, the interesting question that we explore deals with analyzing the hardness of computing expectation using fine-grained analysis and parameterized complexity, where we are interested in the exponent of polynomial runtime.

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

Let $T_{det}(Q, \bar{D}, c)$ (see Sec. 2.4 for further details) denote the runtime for query Q , deterministic database \bar{D} , and multiplicity bound c . This paper considers \mathcal{RA}^+ queries for which order of operations is *explicit*, as opposed to other query languages, e.g. Datalog, UCQ. Thus, since order of operations affects runtime, we denote the optimized \mathcal{RA}^+ query picked by an arbitrary production system as $\text{OPT}(Q) = \min_{Q' \in \mathcal{RA}^+, Q' \equiv Q} T_{det}(Q', \bar{D}, c)$. Then $T_{det}(\text{OPT}(Q), \bar{D}, c)$ is the runtime for the optimized query.²

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

Specifically, depending on what hardness result/conjecture we assume, we get various weaker or stronger versions of *no* as an

²Note that our work applies to any $Q \in \mathcal{RA}^+$, which implies that specific heuristics for choosing an optimized query can be abstracted away, i.e., our work does not consider heuristic techniques.

answer to our question. To make some sense of the other lower bounds in Table 1, we note that it is not too hard to show that $T^*(Q, \mathcal{D}) \leq O\left((T_{det}(\text{OPT}(Q), D, c))^k\right)$, where k is the join width (our notion of join width follows from Definition 2.2 and Fig. 1.) of the query Q over all result tuples t (and the parameter that defines our family of hard queries).

What our lower bound in the third row says is that one cannot get more than a polynomial improvement over essentially the trivial algorithm for Problem 1.1. However, this result assumes a hardness conjecture that is not as well studied as those in the first two rows of the table (see Sec. 3 for more discussion on the hardness assumptions). Further, we note that existing results³ already imply the claimed lower bounds if we were to replace the $T_{det}(\text{OPT}(Q), D, c)$ by just n (indeed these results follow from known lower bounds for deterministic query processing). Our contribution is to then identify a family of hard queries where deterministic query processing is ‘easy’ but computing the expected multiplicities is hard.

Our upper bound results. We introduce a $(1 \pm \epsilon)$ -approximation algorithm that computes Problem 1.1 in time $O_\epsilon(T_{det}(\text{OPT}(Q), D, c))$. This means, when we are okay with approximation, that we solve Problem 1.1 in time linear in the size of the deterministic query and bag PDBs are deployable in practice. In contrast, known approximation techniques ([32, 39]) in set-PDBs need time $\Omega(T_{det}(\text{OPT}(Q), D, c)^{2k})$ (see Appendix F). Further, our approximation algorithm works for a more general notion of bag PDBs beyond c -TIDBs (see Sec. 2.2).

1.1 Polynomial Equivalence

A common encoding of probabilistic databases (e.g., in [1, 4, 29, 30] and many others) relies on annotating tuples with lineages or propositional formulas that describe the set of possible worlds that the tuple appears in. The bag semantics analog is a provenance/lineage polynomial (see Fig. 1) $\Phi[Q, D, t]$ [26], a polynomial with non-zero integer coefficients and exponents, over variables \mathbf{X} encoding input tuple multiplicities. Evaluating a lineage polynomial for a query result tuple t_{out} by, for each tuple t_{in} , assigning the variable $X_{t_{in}}$ encoding the tuple’s multiplicity to the tuple’s multiplicity in the possible world yields the multiplicity of the t_{out} in the query result for this world.

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

PROBLEM 1.2 (EXPECTED MULTIPLICITY OF LINEAGE POLYNOMIALS). *Given an \mathcal{RA}^+ query Q , c -TIDB \mathcal{D} and result tuple t , compute the expected multiplicity of the polynomial $\Phi[Q, D, t]$ (i.e., for $\mathbf{W} \in \{0, \dots, c\}^D$, compute $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$).*

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

³This claim follows from known results for the problem of counting k -cliques, where for a query Q over database D that counts the number of k -cliques. Specifically, a lower bound of the form $\Omega(n^{1+\epsilon_0})$ for some $\epsilon_0 > 0$ follows from the triangle detection hypothesis (this like our result is for $k = 3$). Second, a lower bound of $\omega(n^{C_0})$ for all $C_0 > 0$ under the assumption $\#\mathbf{W}[0] \neq \#\mathbf{W}[1]$ for counting k -clique [23]. Finally, a lower bound of $\Omega(n^{c_0 \cdot k})$ for some $c_0 > 0$ was shown by [10] (under the strong exponential time hypothesis).

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

Table 1: Our lower bounds for a specific hard query Q_{hard} parameterized by k . For $\mathcal{D} = \{0, \dots, c\}^D, \mathcal{P}$ those with ‘Multiple’ in the second column need the algorithm to be able to handle multiple \mathcal{P} , i.e. probability distributions (for a given D). The last column states the hardness assumptions that imply the lower bounds in the first column (ϵ_0, C_0, c_0 are constants that are independent of k).

All of our results rely on working with a *reduced* form ($\tilde{\Phi}$) of the lineage polynomial Φ . In fact, it turns out that for the 1-TIDB case, computing the expected multiplicity (over bag query semantics) is *exactly* the same as evaluating this reduced polynomial over the probabilities that define the 1-TIDB. This is also true when the query input(s) is a block independent disjoint probabilistic database [39] (bag query semantics with tuple multiplicity at most 1), for which the proof of Lemma 1.4 (introduced shortly) holds.

Next, we motivate this reduced polynomial. Consider the query Q_1 defined as follows over the bag relations of Fig. 2:

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

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 $\Phi_1^2(A, B, C, E, X, Y, Z) = A^2X^2B^2 + B^2Y^2E^2 + B^2Z^2C^2 + 2AXB^2YE + 2AXB^2ZC + 2B^2YEZC$.

To compute $\mathbb{E}[\Phi_1^2]$ we can use linearity of expectation and push the expectation through each summand. To keep things simple, let us focus on the monomial $\Phi_1^{(ABX)^2} = A^2X^2B^2$ as the procedure is the same for all other monomials of Φ_1^2 . Let W_X be the random variable corresponding to a lineage variable X . Because the distinct variables in the product are independent, we can push expectation through them yielding $\mathbb{E}[W_A^2 W_X^2 W_B^2] = \mathbb{E}[W_A^2] \mathbb{E}[W_X^2] \mathbb{E}[W_B^2]$. Since $W_A, W_B \in \{0, 1\}$ we can further derive $\mathbb{E}[W_A] \mathbb{E}[W_X^2] \mathbb{E}[W_B]$ by the fact that for any $W \in \{0, 1\}$, $W^2 = W$. Observe that if $X \in \{0, 1\}$, then we further would have $\mathbb{E}[W_A] \mathbb{E}[W_X] \mathbb{E}[W_B] = p_A \cdot p_X \cdot p_B$ (denoting $\Pr[W_A = 1] = p_A$) = $\tilde{\Phi}_1^{(ABX)^2}(p_A, p_X, p_B)$ (see ii) of Definition 1.3). However, in this example, we get stuck with $\mathbb{E}[W_X^2]$, since $W_X \in \{0, 1, 2\}$ and for $W_X \leftarrow 2$, $W_X^2 \neq W_X$.

Denote the variables of Φ to be $\text{VARS}(\Phi)$. In the c -TIDB setting, $\Phi(\mathbf{X})$ has an equivalent reformulation ($\Phi_R(\mathbf{X}_R)$) that is of use to us, where $|\mathbf{X}_R| = c \cdot |\mathbf{X}|$. Given $X_t \in \text{VARS}(\Phi)$ and integer valuation $X_t \in \{0, \dots, c\}$. We can replace X_t by $\sum_{j \in [c]} jX_{t,j}$ where the variables $(X_{t,j})_{j \in [c]}$ are disjoint with integer assignments $X_t \in \{0, 1\}$. Then for any $\mathbf{W} \in \{0, \dots, c\}^D$ and corresponding 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 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 integer valuation $X_t \in [c]$, $X_{t,j} \in \{0, 1\}$ we have the equality $X_t = j = \sum_{j \in [c]} jX_{t,j}$. Considering again our example,

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

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

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

Continuing with the example⁵ $\Phi_1^2(A, B, C, E, X_1, X_2, Y, Z)$ we have $\tilde{\Phi}_1^2(A, B, C, E, X_1, X_2, Y, Z) =$

$$\begin{aligned} &A \left(\sum_{j \in [c]} j^2 X_j \right) B + BYE + BZC + 2A \left(\sum_{j \in [c]} j^2 X_j \right) BYE \\ &+ 2A \left(\sum_{j \in [c]} j^2 X_j \right) BZC + 2BYEZC \\ &= ABX_1 + AB(2)^2 X_2 + BYE + BZC + 2AX_1 BYE + 2A(2)^2 X_2 BYE \\ &+ 2AX_1 BZC + 2A(2)^2 X_2 BZC + 2BYEZC. \end{aligned}$$

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

LEMMA 1.4. For any c -TIDB \mathcal{D} , \mathcal{RA}^+ query Q , and lineage polynomial $\Phi(\mathbf{X}) = \Phi[Q, D, t](\mathbf{X})$, it holds that $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi(\mathbf{W})] = \tilde{\Phi}(\mathbf{p})$, where $\mathbf{p} = (p_t)_{t \in D}$.

⁴This is the representation, typically used in set-PDBs, where the polynomial is re-represented as sum of ‘pure’ products. See Definition 2.1 for a formal definition.

⁵To 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 $1^2 \cdot A^2 = A$.

The proof follows by Proposition 2.4 and Lemma 2.6.

1.2 Our Techniques

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

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

We adopt a two-step intensional model of query evaluation used in set-PDBs, as illustrated in Fig. 2: (i) Lineage Computation (LC): Given input D and Q , output every tuple t that possibly satisfies Q , annotated with its lineage polynomial ($\Phi(\mathbf{X}) = \Phi[Q, D, t](\mathbf{X})$); (ii) Expectation Computation (EC): Given $\Phi(\mathbf{X})$ for each tuple, compute $\mathbb{E}_{W \sim \mathcal{P}}[\Phi(\mathbf{W})]$. Let $T_{LC}(Q, D, C)$ denote the runtime of LC when it outputs C (which is a representation of Φ as an arithmetic circuit – more on this representation in Sec. 2.3). Denote by $T_{EC}(C, \epsilon)$ (recall C is the output of LC) the runtime of EC, which we can leverage Definition 1.3 and Lemma 1.4 to address the next formal objective:

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

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

However, systems can directly emit compact, factorized representations of $\Phi(\mathbf{X})$ (e.g., as a consequence of the standard projection push-down optimization [25]). For example, in Fig. 2, $B(Y + Z)$ is a factorized representation of the SMB-form $BY + BZ$. Accordingly,

this work uses (arithmetic) circuits⁶ as the representation system of $\Phi(\mathbf{X})$.

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

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

For an upper bound on approximating the expected count, it is easy to check that if all the probabilities are constant then (with an additive adjustment) $\Phi(p_1, \dots, p_n)$ (i.e. evaluating the original lineage polynomial over the probability values) is a constant factor approximation. This is illustrated in the following example using Q_1^2 from earlier. To aid in presentation we assume $c = 2$ for variable X and $c = 1$ for all other variables. Let p_A denote $Pr[A = 1]$. In computing $\tilde{\Phi}$, we have some cancellations to deal with:

$$\begin{aligned} \Phi_{1,R}^2(\mathbf{X}) = & A^2(X_1^2 + 4X_1X_2 + 4X_2^2)B^2 + B^2Y^2E^2 + B^2Z^2C^2 + 2AX_1B^2YE \\ & + 2AX_2B^2YE + 2AX_1B^2ZC + 2AX_2B^2ZC + 2B^2YEZC \end{aligned}$$

This then implies

$$\begin{aligned} \tilde{\Phi}^2(\mathbf{X}) = & AX_1B + 4AX_2B + BYE + BZC + 2AX_1BYE + 2AX_2BYE + 2AX_1BZC \\ & + 2AX_2BZC + 2BYEZC \end{aligned}$$

Substituting \mathbf{p} for \mathbf{X} , $\Phi_{1,R}^2(\mathbf{p}) =$

$$\begin{aligned} & p_A^2 p_{X_1}^2 p_B^2 + 4p_A p_{X_1} p_{X_2} p_B^2 + 4p_A^2 p_{X_2}^2 p_B^2 + p_B^2 p_Y^2 p_E^2 + p_B^2 p_Z^2 p_C^2 + 2p_A p_{X_1} p_B^2 p_Y p_E \\ & + 2p_A p_{X_2} p_B^2 p_Y p_E + 2p_A p_{X_1} p_B^2 p_Z p_C + 2p_A p_{X_2} p_B^2 p_Z p_C + 2p_B^2 p_Y p_E p_Z p_C \\ & \leq p_A p_{X_1} p_B + 4p_A^2 p_{X_1} p_{X_2} p_B^2 + 4p_A p_{X_2} p_B + p_B p_Y p_E + p_B p_Z p_C + 2p_A p_{X_1} p_B p_Y p_E \\ & + 2p_A p_{X_2} p_B p_Y p_E + 2p_A p_{X_1} p_B p_Z p_C + 2p_A p_{X_2} p_B p_Z p_C + 2p_B p_Y p_E p_Z p_C \\ & = \tilde{\Phi}_1^2(\mathbf{p}) + 4p_A^2 p_{X_1} p_{X_2} p_B^2. \end{aligned}$$

If we assume that all probability values are at least $p_0 > 0$, then given access to $\Phi_{1,R}^2(\mathbf{p}) - 4p_A^2 p_{X_1} p_{X_2} p_B^2$ we get that $\Phi_{1,R}^2(\mathbf{p}) - 4p_A^2 p_{X_1} p_{X_2} p_B^2$ is in the range $\left((p_0)^3 \cdot (\tilde{\Phi}_1^2(\mathbf{p}), \tilde{\Phi}_1^2(\mathbf{p}))\right)$. Note however, that this is *not a tight approximation*. In sec. 4 we demonstrate that a $(1 \pm \epsilon)$ (multiplicative) approximation with competitive performance is achievable. To get an $(1 \pm \epsilon)$ -multiplicative approximation and solve Problem 1.6, using C we uniformly sample monomials from the equivalent SMB representation of Φ (without materializing the SMB representation) and ‘adjust’ their contribution to $\tilde{\Phi}(\cdot)$.

Applications. Recent work in heuristic data cleaning [7, 42, 45, 45, 51] emits a PDB when insufficient data exists to select the ‘correct’ data repair. Probabilistic data cleaning is a crucial innovation, as the alternative is to arbitrarily select one repair and ‘hope’ that queries receive meaningful results. Although PDB queries instead

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

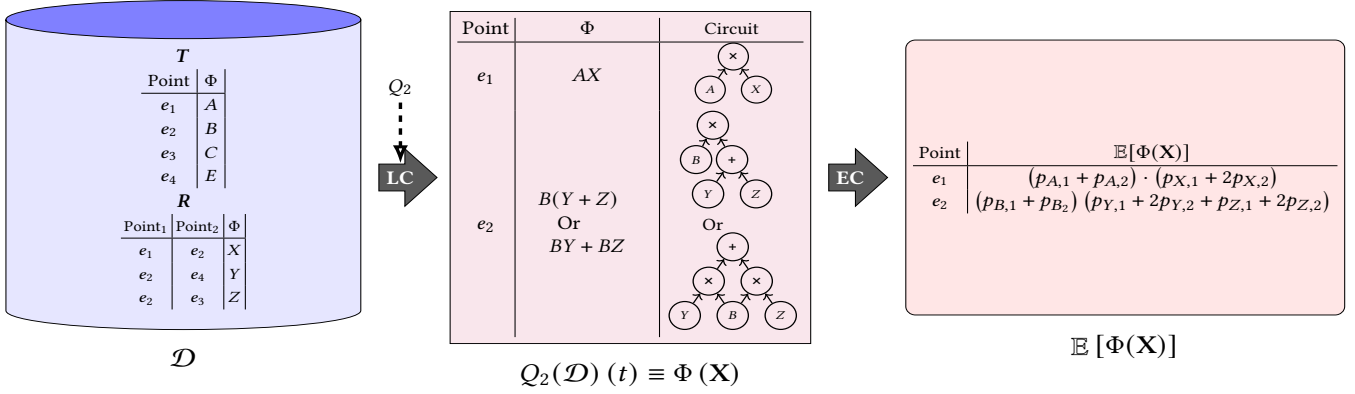


Figure 2: Intensional Query Evaluation Model ($Q_2 = \pi_{\text{Point}}(T \bowtie_{\text{Point}=\text{Point}_1} R)$ where, for table R , $c = 2$, while for T , $c = 1$.)

convey the trustworthiness of results [36], they are impractically slow [18, 19], even in approximation (see Appendix F). Bags, as we consider, are sufficient for production use, where bag-relational algebra is already the default for performance reasons. Our results show that bag-PDBs can be competitive, laying the groundwork for probabilistic functionality in production database engines.

Concurrent Work. In work independent of ours, Grohe, et. al. [28] investigate bag-TIDBs allowing for unbounded multiplicities (which requires them to explicitly address the issue of a succinct representation of the distribution over infinitely many multiplicities). While the authors observe that computing the expected value of an output tuple multiplicity is in polynomial time, no further (fine-grained) analysis of the expected value is considered. The work primarily investigates the query evaluation problem over bag-TIDBs when computing the probability of an output tuple having at most a multiplicity of k , showing that a dichotomy exists for this problem. Our work in contrast assumes a finite bound on the multiplicities where we simply list the finitely many probability values (and hence do not need consider a more succinct representation). Further, our work primarily looks into the fine-grained analysis of computing the expected multiplicity of an output tuple.

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

2 BACKGROUND AND NOTATION

2.1 Polynomial Definition and Terminology

Given an index set S over variables X_t for $t \in S$, a (general) polynomial ϕ over $(X_t)_{t \in S}$ with individual degree $K < \infty$ is formally defined as:

$$\phi((X_t)_{t \in S}) = \sum_{\mathbf{d} \in \{0, \dots, K\}^S} c_{\mathbf{d}} \cdot \prod_{t \in S} X_t^{d_t} \quad \text{where } c_{\mathbf{d}} \in \mathbb{N}. \quad (1)$$

DEFINITION 2.1 (STANDARD MONOMIAL BASIS). *The term $\prod_{t \in S} X_t^{d_t}$ in Eq. (1) is a monomial. A polynomial $\phi(\mathbf{X})$ is in standard monomial basis (SMB) when we keep only the terms with $c_{\mathbf{d}} \neq 0$ from Eq. (1).*

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

DEFINITION 2.2 (DEGREE). *The degree of polynomial $\phi(\mathbf{X})$ is the largest $\sum_{t \in S} d_t$ for all $\mathbf{d} \in \{0, \dots, K\}^S$ such that $c_{(d_1, \dots, d_n)} \neq 0$. We denote the degree of ϕ as $\text{deg}(\phi)$.*

As an example, the degree of the polynomial $X^2 + 2XY^2 + Y^2$ is 3. Product terms in lineage arise only from join operations (Fig. 1), so intuitively, the degree of a lineage polynomial is analogous to the largest number of joins needed to produce a result tuple.

We call a polynomial $\Phi(\mathbf{X})$ a c -TIDB-lineage polynomial (or simply lineage polynomial), if it is clear from context that there exists an \mathcal{RA}^+ query Q , c -TIDB \mathcal{D} , and result tuple t such that $\Phi(\mathbf{X}) = \Phi[Q, \mathcal{D}, t](\mathbf{X})$.

2.2 Binary-BIDB

A block independent database BIDB \mathcal{D}' models a set of worlds each of which consists of a subset of the possible tuples D' , where D' is partitioned into m blocks B_i and all B_i are independent random events. \mathcal{D}' further constrains that all $t \in B_i$ for all $i \in [m]$ of D' be disjoint events. We refer to any monomial that includes $X_t X_{t'}$ for $t \neq t' \in B_i$ as a *cancellation*. We define next a specific construction of BIDB that is useful for our work.

DEFINITION 2.3 (BINARY-BIDB). *Define a Binary-BIDB to be the pair $\mathcal{D}' = (\times_{t \in D'} \{0, c_t\}, \mathcal{P}')$, where D' is the set of possible tuples such that each $t \in D'$ has a multiplicity domain of $\{0, c_t\}$, with $c_t \in \mathbb{N}$. D' is partitioned into m independent blocks B_i , for $i \in [m]$, of disjoint tuples. \mathcal{P}' is characterized by the vector $(p_t)_{t \in D'}$ where for every block B_i , $\sum_{t \in B_i} p_t \leq 1$. Given $W \in \times_{t \in D'} \{0, c_t\}$ and for $i \in [m]$,*

$$\text{let } p_i(W) = \begin{cases} 1 - \sum_{t \in B_i} p_t & \text{if } W_t = 0 \text{ for all } t \in B_i \\ 0 & \text{if there exists } t, t' \in B_i, W_t, W_{t'} \neq 0 \\ p_t & W_t \neq 0 \text{ for the unique } t \in B_i. \end{cases}$$

\mathcal{P}' is the probability distribution across all worlds such that, given $W \in \times_{t \in D'} \{0, c_t\}$, $\text{Pr}[\mathbf{W} = W] = \prod_{i \in [m]} p_i(W)$.⁷

⁷We slightly abuse notation here, denoting a world vector as W rather than \mathbf{W} to distinguish between the random variable and the world instance. When there is no ambiguity, we will denote a world vector as \mathbf{W} .

$$\begin{aligned}
\Phi' \left[\pi_A(Q), \bar{D}', t_j \right] &= \Phi' \left[Q_1 \cup Q_2, \bar{D}', t_j \right] = \\
&\sum_{\substack{t_{j'}, \\ \pi_A(t_{j'})=t_j}} \Phi' \left[Q, \bar{D}', t_{j'} \right] \quad \Phi' \left[Q_1, \bar{D}', t_j \right] + \Phi' \left[Q_2, \bar{D}', t_j \right] \\
\Phi' \left[\sigma_\theta(Q), \bar{D}', t_j \right] &= \Phi' \left[Q_1 \bowtie Q_2, \bar{D}', t_j \right] = \\
\begin{cases} \theta = 1 & \Phi' \left[Q, \bar{D}', t_j \right] \quad \Phi' \left[Q_1, \bar{D}', \pi_{attr(Q_1)}(t_j) \right] \\ \theta = 0 & 0 \quad \cdot \Phi' \left[Q_2, \bar{D}', \pi_{attr(Q_2)}(t_j) \right] \end{cases} \\
\Phi' \left[R, \bar{D}', t_j \right] &= j \cdot X_{t,j}.
\end{aligned}$$

Figure 3: Construction of the lineage (polynomial) for an \mathcal{RA}^+ query Q over \bar{D}' .

Fig. 3 shows the lineage construction of $\Phi'(X)$ given \mathcal{RA}^+ query Q for arbitrary deterministic \bar{D}' . Note that the semantics differ from Fig. 1 only in the base case.

PROPOSITION 2.4 (c-TIDB REDUCTION). *Given c-TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, let $\mathcal{D}' = (\times_{t \in D} \{0, c_t\}, \mathcal{P}')$ be the Binary-BIDB obtained in the following manner: for each $t \in D$, create block $B_t = \{\langle t, j \rangle_{j \in [c]}\}$ of disjoint tuples, for all $j \in [c]$ where $c_{t_j} = j$ for each t_j in \mathcal{D}' . The probability distribution \mathcal{P}' is characterized by the vector $\mathbf{p} = (p_{t,j})_{t \in D, j \in [c]}$. Then, the distributions \mathcal{P} and \mathcal{P}' are equivalent.*

Aaron says: Not entirely sure if I the notation above $c_{t_j} = j$ for each t_j in \mathcal{D}' .

We now define the reduced polynomial $\tilde{\Phi}'$ of a Binary-BIDB.

DEFINITION 2.5 ($\tilde{\Phi}'$). *Given a polynomial $\Phi'(X)$ generated from a Binary-BIDB and let $\tilde{\Phi}'(X)$ denote the reduced form of $\Phi'(X)$ derived as follows: i) compute SMB ($\Phi'(X)$) eliminating all monomials with cross terms $X_t X_{t'}$ for $t \neq t' \in B_i$ and ii) reduce all variable exponents $e > 1$ to 1.*

Then given $\mathbf{W} \in \{0, 1\}^{D'}$ over the reduced Binary-BIDB of Proposition 2.4, the disjoint requirement and the semantics for constructing the lineage polynomial over a Binary-BIDB, $\Phi'(\mathbf{W})$ is of the same structure as the reformulated polynomial $\Phi_R(\mathbf{W})$ of step i) from Definition 1.3, which then implies that $\tilde{\Phi}'$ is the reduced polynomial that results from step ii) of both Definition 1.3 and Definition 2.5, and further that Lemma 1.4 immediately follows for Binary-BIDB polynomials.

LEMMA 2.6. *Given any Binary-BIDB \mathcal{D}' , \mathcal{RA}^+ query Q , and lineage polynomial $\Phi'(X) = \Phi'[Q, \mathcal{D}', t](X)$, it holds that $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}'}[\tilde{\Phi}'(\mathbf{W})] = \tilde{\Phi}'(\mathbf{p})$.*

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

COROLLARY 2.7. *If Φ is a 1-BIDB lineage polynomial already in SMB, then the expectation of Φ , i.e., $\mathbb{E}[\Phi] = \tilde{\Phi}(p_1, \dots, p_n)$ can be computed in $O(|\Phi|)$ time.*

2.2.1 Possible World Semantics. In this section, we show how the traditional possible worlds semantics corresponds to our setup. Readers can safely skip this part without missing anything vital to the results of this paper.

Queries over probabilistic databases are traditionally viewed as being evaluated using the so-called possible world semantics. A general bag-PDB can be defined as the pair $\mathcal{D} = (\Omega, \mathcal{P})$ where Ω is the set of possible worlds represented by \mathcal{D} and \mathcal{P} the probability distribution over Ω . Under the possible world semantics, the result of a query Q over an incomplete database Ω is the set of query answers produced by evaluating Q over each possible world $\omega \in \Omega$: $\{Q(\omega) : \omega \in \Omega\}$. The result of a query is the pair $(Q(\Omega), \mathcal{P}')$ where \mathcal{P}' is a probability distribution that assigns to each possible query result the sum of the probabilities of the worlds that produce this answer: $Pr[\omega \in \Omega] = \sum_{\omega' \in \Omega, Q(\omega')=Q(\omega)} Pr[\omega']$.

Suppose that \mathcal{D}' is a reduced Binary-BIDB from c-TIDB \mathcal{D} as defined by Proposition 2.4. Instead of looking only at the possible worlds of \mathcal{D}' , one can consider the set of all worlds, including those that cannot exist due to, e.g., disjointness. Since $|D| = n$ the all worlds set can be modeled by $\mathbf{W} \in \{0, 1\}^{nc}$, such that $\mathbf{W}_{t,j} \in \mathbf{W}$ represents whether or not the multiplicity of t is j (here and later, especially in Sec. 4, we will rename the variables as X_1, \dots, X_n , where $n' = \sum_{t \in D} |B_t|$).⁸ We can denote a probability distribution over all $\mathbf{W} \in \{0, 1\}^{nc}$ as \mathcal{P}' . When \mathcal{P}' is the one induced from each $p_{t,j}$ while assigning $Pr[\mathbf{W}] = 0$ for any \mathbf{W} with $\mathbf{W}_{t,j}, \mathbf{W}_{t,j'} \neq 0$ for $j \neq j'$, we end up with a bijective mapping from \mathcal{P} to \mathcal{P}' , such that each mapping is equivalent, implying the distributions are equivalent, and thus query results. Appendix A.1.2 has more details.

We now make a meaningful connection between possible world semantics and world assignments on the lineage polynomial.

PROPOSITION 2.8 (EXPECTATION OF POLYNOMIALS). *Given a bag-PDB $\mathcal{D} = (\Omega, \mathcal{P})$, \mathcal{RA}^+ query Q , and lineage polynomial $\Phi[Q, D, t]$ for arbitrary result tuple t , we have (denoting \mathbf{D} as the random variable over Ω): $\mathbb{E}_{\mathbf{D} \sim \mathcal{P}}[Q(\mathbf{D})(t)] = \mathbb{E}_{\mathbf{W} \sim \mathcal{P}'}[\Phi[Q, D, t](\mathbf{W})]$.*

A formal proof of Proposition 2.8 is given in Appendix A.3.⁹

2.3 Formalizing Problem 1.6

We focus on the problem of computing $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$ from now on, assume implicit Q, D, t , and drop them from $\Phi[Q, D, t]$ (i.e., $\Phi(X)$ will denote a polynomial).

Problem 1.6 asks if there exists a linear time approximation algorithm in the size of a given circuit C which encodes $\Phi(X)$. Recall that in this work we represent lineage polynomials via *arithmetic circuits* [8], a standard way to represent polynomials over fields (particularly in the field of algebraic complexity) that we use for polynomials over \mathbb{N} in the obvious way. Since we are specifically using circuits to model lineage polynomials, we can refer to these circuits as lineage circuits. However, when the meaning is clear, we will drop the term lineage and only refer to them as circuits.

DEFINITION 2.9 (CIRCUIT). *A circuit C is a Directed Acyclic Graph (DAG) whose source gates (in degree of 0) consist of elements in either \mathbb{N} or $\mathbf{X} = (X_1, \dots, X_n)$. For each result tuple there exists one sink gate. The internal gates have binary input and are either sum (+) or product (\times) gates. Each gate has the following members: type, input, val, partial, degree, Lweight, and Rweight, where type is the value*

⁸In this example, $|B_t| = c$ for all t .

⁹Although Proposition 2.8 follows, e.g., as an obvious consequence of [30]'s Theorem 7.1, we are unaware of any formal proof for bag-probabilistic databases.

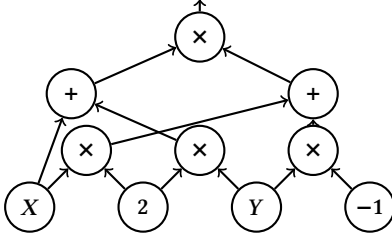


Figure 4: Circuit encoding of $(X + 2Y)(2X - Y)$

type $\{+, \times, \text{VAR}, \text{NUM}\}$ and input the list of inputs. Source gates have an extra member `val` storing the value. C_L (C_R) denotes the left (right) input of C .

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

The circuits in Fig. 2 encode their respective polynomials in column Φ . Note that the circuit C representing AX and the circuit C' representing $B(Y + Z)$ each encode a tree, with edges pointing towards the root.

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

DEFINITION 2.10 (POLY(\cdot)). *POLY(C) maps the sink of circuit C to its corresponding polynomial in SMB. POLY(\cdot) is recursively defined on C as follows, with addition and multiplication following the standard interpretation for polynomials:*

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

C need not encode $\Phi(X)$ in the same, default SMB representation. For instance, C could encode the factorized representation $(X + 2Y)(2X - Y)$ of $\Phi(X) = 2X^2 + 3XY - 2Y^2$, as shown in Fig. 4, while $\text{POLY}(C) = \Phi(X)$ is always the equivalent SMB representation.

DEFINITION 2.11 (CIRCUIT SET). *CSet($\Phi(X)$) is the set of all possible circuits C such that $\text{POLY}(C) = \Phi(X)$.*

The circuit of Fig. 4 is an element of $\text{CSet}(2X^2 + 3XY - 2Y^2)$. One can think of $\text{CSet}(\Phi(X))$ as the infinite set of circuits where for each element C , $\text{POLY}(C) = \Phi(X)$.

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

DEFINITION 2.12 (THE EXPECTED RESULT MULTIPLICITY PROBLEM). *Let \mathcal{D}' be an arbitrary c -TiDB and X be the set of variables annotating tuples in \mathcal{D}' . Fix an \mathcal{RA}^+ query Q and a result tuple t . The EXPECTED RESULT MULTIPLICITY PROBLEM is defined as follows:*

Input: $C \in \text{CSet}(\Phi(X))$ for $\Phi'(X) = \Phi'[Q, \mathcal{D}', t]$

Output: $\mathbb{E}_{W \sim \mathcal{P}}[\Phi'[Q, \mathcal{D}', t](W)]$.

2.4 Relationship to Deterministic Query Runtimes

In Sec. 1, we introduced the structure $T_{det}(\cdot)$ to analyze the runtime complexity of Problem 1.1. To decouple our results from specific join algorithms, we first lower bound the cost of a join.

DEFINITION 2.13 (JOIN COST). *Denote by $T_{join}(R_1, \dots, R_m)$ the runtime of an algorithm for computing the m -ary join $R_1 \bowtie \dots \bowtie R_m$. We require only that the algorithm must enumerate its output, i.e., that $T_{join}(R_1, \dots, R_m) \geq |R_1 \bowtie \dots \bowtie R_m|$. With this definition of $T_{join}(\cdot)$, worst-case optimal join algorithms are handled.*

Worst-case optimal join algorithms [37, 38] and query evaluation via factorized databases [40] (as well as work on FAQs [34]) can be modeled as \mathcal{RA}^+ queries (though the query size is data dependent). For these algorithms, $T_{join}(R_1, \dots, R_n)$ is linear in the AGM bound [5]. Our cost model for general query evaluation follows from the join cost:

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

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

Lemma D.2 and Lemma D.3 show that for any \mathcal{RA}^+ query Q and D , there exists a circuit C^* such that $T_{LC}(Q, D, C^*)$ and $|C^*|$ are both $O(T_{det}(\text{OPT}(Q), D, c))$. Recall we assumed these two bounds when we moved from Problem 1.5 to Problem 1.6. Lastly, we can handle FAQs and factorized databases by allowing for optimization, i.e. $\text{OPT}(Q)$.

3 HARDNESS OF EXACT COMPUTATION

In this section, we will prove the hardness results claimed in Table 1 for a specific (family) of hard instance (Q_{hard}, \mathcal{D}) for Problem 1.2 where \mathcal{D} is a 1-TiDB. Note that this implies hardness for c -TiDBs ($c \geq 1$), showing Problem 1.2 cannot be done in $O(T_{det}(\text{OPT}(Q), D, c))$ runtime. The results also apply to Binary-BiDB and other more general PDBs.

3.1 Preliminaries

Our hardness results are based on (exactly) counting the number of (not necessarily induced) subgraphs in G isomorphic to H . Let $\#(G, H)$ denote this quantity. 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 in its adjacency list representation): $\#(G, \triangle)$ (the number of triangles), $\#(G, \triangle\triangle\triangle)$ (the number of 3-matchings), and the latter's generalization $\#(G, \mathfrak{g} \dots \mathfrak{g}^k)$ (the number of k -matchings). We use $T_{match}(k, G)$ to denote the optimal runtime of computing $\#(G, \mathfrak{g} \dots \mathfrak{g}^k)$ exactly. Our hardness

results in Sec. 3.2 are based on the following hardness results/conjectures:

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

CONJECTURE 3.2 ([12]). *For every $G = (V, E)$, we have $T_{match}(k, G) \geq n^{\Omega(k/\log k)}$.*

We note that the above conjecture is somewhat non-standard. In particular, the best known algorithm to compute $\#(G, \mathfrak{g} \cdots \mathfrak{g}^k)$ takes time $\Omega(|V|^{k/2})$ [11]. What the above conjecture is saying is that one can only hope for a polynomial improvement over the state of the art algorithm to compute $\#(G, \mathfrak{g} \cdots \mathfrak{g}^k)$.

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

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

The so called *Triangle detection hypothesis* (cf. [35]), which states that detecting the presence of triangles in G takes time $\Omega(|E|^{4/3})$, implies that in Conjecture 3.3 we can take $\epsilon_0 \geq \frac{1}{3}$.

All of our hardness results rely on a simple lineage 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 lineage polynomial has a variable X_i for every i in $[n]$. Consider the polynomial $\Phi_G(\mathbf{X}) = \sum_{(i,j) \in E} X_i \cdot X_j$. The hard polynomial for our problem will be a suitable power $k \geq 3$ of the polynomial above:

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

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

Returning to Fig. 2, it can be seen that $\Phi_G^k(\mathbf{X})$ is the lineage polynomial from query Q_{hard}^k , which we define next (Q_2 from Sec. 1 is the same query with $k = 2$). Let us alias

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

as R . The query Q_{hard}^k then becomes

```
SELECT COUNT(*) FROM R JOIN R JOIN ... JOIN R
                        k times
```

Consider again the c -TIDB instance \mathcal{D} of Fig. 2 and, for our hard instance, let $c = 1$. \mathcal{D} generalizes to one compatible to Definition 3.4 as follows. Relation T has n tuples corresponding to each vertex for i in $[n]$, each with probability p and R has tuples corresponding to the edges E (each with probability of 1).¹⁰ In other words, this instance

¹⁰Technically, $\Phi_G^k(\mathbf{X})$ should have variables corresponding to tuples in R as well, but since they always are present with probability 1, we drop those. Our argument also works when all the tuples in R also are present with probability p but to simplify notation we assign probability 1 to edges.

\mathcal{D} contains the set of n unary tuples in T (which corresponds to V) and m binary tuples in R (which corresponds to E). Note that this implies that Φ_G^k is indeed a 1-TIDB lineage polynomial.

Next, we note that the runtime for answering Q_{hard}^k on deterministic database D , as defined above, is $O_k(m)$ (i.e. deterministic query processing is 'easy' for this query):

LEMMA 3.5. *Let Q_{hard}^k and D be as defined above. Then $T_{det}(Q_{hard}^k, D)$ is $O_k(m)$.*

3.2 Multiple Distinct p Values

We are now ready to present our main hardness result.

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

Note that the second row of Table 1 follows from Proposition 2.8, Theorem 3.6, Lemma 3.5, and Theorem 3.1 while the third row is proved by Proposition 2.8, Theorem 3.6, Lemma 3.5, and Conjecture 3.2. Since Conjecture 3.2 is non-standard, the latter hardness result should be interpreted as follows. Any substantial polynomial improvement for Problem 1.2 (over the trivial algorithm that converts Φ into SMB and then uses Corollary 2.7 for EC) would lead to an improvement over the state of the art *upper* bounds on $T_{match}(k, G)$. Finally, note that Theorem 3.6 needs one to be able to compute the expected multiplicities over $(2k + 1)$ distinct values of p_i , each of which corresponds to distinct \mathcal{P} (for the same D), which explain the 'Multiple' entry in the second column in the second and third row in Table 1. Next, we argue how to get rid of this latter requirement.

3.3 Single p value

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

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

Note that Proposition 2.8 and Theorem 3.7 above imply the hardness result in the first row of Table 1. We note that Theorem 3.1 and Conjecture 3.2 (and the lower bounds in the second and third row of Table 1) need k to be large enough (in particular, we need a family of hard queries). But the above Theorem 3.7 (and the lower bound in first row of Table 1) holds for $k = 3$ (and hence for a fixed query).

4 $1 \pm \epsilon$ APPROXIMATION ALGORITHM

We showed in Sec. 3 that a runtime of $O(T_{det}(\text{OPT}(Q), D, c))$ cannot be achieved for Problem 1.2. In light of this, we desire to produce

an approximation algorithm that runs in time $O(T_{det}(\text{OPT}(Q), D, c))$. We do this by showing the result via circuits, such that our approximation algorithm for this problem runs in $O(|C|)$ for a very broad class of circuits, (thus affirming Problem 1.6); see the discussion after Lemma 4.9 for more. The following approximation algorithm applies to bag query semantics over both c -TIDB lineage polynomials and general BIDB lineage polynomials in practice, where for the latter we note that a 1-TIDB is equivalently a BIDB (blocks are size 1). Our experimental results (see Appendix C.11) which use queries from the PDBench benchmark [41] show a low γ (see Definition 4.6) supporting the notion that our bounds hold for general BIDB in practice.

Corresponding proofs and pseudocode for all formal statements and algorithms can be found in Appendix C.

4.1 Preliminaries and some more notation

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

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

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

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

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

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

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

DEFINITION 4.4 ($\text{DEG}(\cdot)$). ¹¹ *$\text{DEG}(C)$ is defined recursively as follows:*

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

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

¹¹Note that the degree of $\text{POLY}(|C|)$ is always upper bounded by $\text{DEG}(C)$ and the latter can be strictly larger (e.g. consider the case when C multiplies two copies of the constant 1— here we have $\text{deg}(C) = 1$ but degree of $\text{POLY}(|C|)$ is 0).

DEFINITION 4.5 ($\overline{M}(\cdot, \cdot)$). ¹² *In a RAM model of word size of W -bits, $\overline{M}(M, W)$ denotes the complexity of multiplying two integers represented with M -bits. (We will assume that for input of size N , $W = O(\log N)$.)*

Finally, to get linear runtime results, we will need to define another parameter modeling the (weighted) number of monomials in $E(C)$ that need to be ‘canceled’ when monomials with dependent variables are removed (Sec. 2.2). Let $\text{ISIND}(\cdot)$ be a boolean function returning true if monomial v_m is composed of independent variables and false otherwise; further, let $\mathbb{1}_\theta$ also be a boolean function returning true if θ evaluates to true.

DEFINITION 4.6 (PARAMETER γ). *Given a Binary-BIDB circuit C define*

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

4.2 Our main result

We solve Problem 1.6 for any fixed $\epsilon > 0$ in what follows.

Algorithm Idea. Our approximation algorithm (`APPROXIMATE $\tilde{\Phi}$` pseudo code in Appendix C.1) is based on the following observation. Given a lineage polynomial $\Phi(X) = \text{POLY}(C)$ for circuit C over Binary-BIDB (recall that all c -TIDB can be reduced to Binary-BIDB by Proposition 2.4), we have:

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

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 C).

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

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

can be computed in time

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

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

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

¹²We note that when doing arithmetic operations on the RAM model for input of size N , we have that $\overline{M}(O(\log N), O(\log N)) = O(1)$. More generally we have $\overline{M}(N, O(\log N)) = O(N \log N \log \log N)$.

The restriction on γ 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 C.11 for experimental results). Further, we can also argue the following result, recalling from Sec. 1 for c -TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, where D is the set of possible tuples across all possible worlds of \mathcal{D} .

LEMMA 4.8. *Given \mathcal{RA}^+ query Q and c -TIDB \mathcal{D} , let C be the circuit computed by $Q(D)$. Then, for the reduced Binary-BIDB \mathcal{D}' there exists an equivalent circuit C' obtained from $Q(D')$, such that $\gamma(C') \leq 1 - c^{-(k-1)}$ with $\text{SIZE}(C') \leq \text{SIZE}(C) + O(nc)$ and $\text{DEPTH}(C') = \text{DEPTH}(C) + O(\log c)$.*

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

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

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

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

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

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

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

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

PROOF OF COROLLARY 4.11. By Lemma 4.8 and Corollary 4.10, the proof follows. \square

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

5 RELATED WORK

Probabilistic Databases (PDBs) have been studied predominantly for set semantics. Approaches for probabilistic query processing (i.e., computing marginal probabilities of tuples), fall into two broad categories. *Intensional* (or *grounded*) query evaluation computes the *lineage* of a tuple and then the probability of the lineage formula. It has been shown that computing the marginal probability of a tuple is #P-hard [48] (by reduction from weighted model counting). The second category, *extensional* query evaluation, is in PTIME, but is limited to certain classes of queries. Dalvi et al. [15] and Olteanu et al. [22] proved dichotomies for UCQs and two classes of queries with negation, respectively. Amarilli et al. investigated tractable classes of databases for more complex queries [2]. Another line of work studies which structural properties of lineage formulas lead to tractable cases [33, 43, 46]. In this paper we focus on intensional query evaluation with polynomials.

Many data models have been proposed for encoding PDBs more compactly than as sets of possible worlds. These include tuple-independent databases [49] (TIDBs), block-independent databases (BIDBs) [44], and *PC-tables* [27]. Fink et al. [20] study aggregate queries over a probabilistic version of the extension of K-relations for aggregate queries proposed in [3] (*pvc-tables*) that supports bags, and has runtime complexity linear in the size of the lineage. However, this lineage is encoded as a tree; the size (and thus the runtime) are still superlinear in $T_{\text{det}}(Q, D, c)$. The runtime bound is also limited to a specific class of (hierarchical) queries, suggesting the possibility of a generalization of [15]'s dichotomy result to bag-PDBs.

Several techniques for approximating tuple probabilities have been proposed in related work [14, 16, 21, 39], relying on Monte Carlo sampling, e.g., [14], or a branch-and-bound paradigm [39]. Our approximation algorithm is also based on sampling.

Compressed Encodings are used for Boolean formulas (e.g. various types of circuits including OBDDs [31]) and polynomials (e.g., factorizations [40]) some of which have been utilized for probabilistic query processing, e.g., [31]. Compact representations for which probabilities can be computed in linear time include OBDDs, SDDs, d-DNNF, and FBDD. [17] studies circuits for absorptive semirings while [47] studies circuits that include negation (expressed as the monus operation). Algebraic Decision Diagrams [6] (ADDs) generalize BDDs to variables with more than two values. Chen et al. [9] introduced the generalized disjunctive normal form. Appendix G covers more related work on fine-grained complexity.

6 CONCLUSIONS AND FUTURE WORK

We have studied the problem of calculating the expected multiplicity of a bag-query result tuple, a problem that has a practical application in probabilistic databases over multisets. We show that under various parameterized complexity hardness results/conjectures computing the expected multiplicities exactly is not possible in time linear in the corresponding deterministic query processing

time. We prove that it is possible to approximate the expectation of a lineage polynomial in linear time in the deterministic query processing over TIDBs and BIDBs (assuming that there are few cancellations). Interesting directions for future work include development of a dichotomy for bag PDBs. While we can handle higher moments (this follows fairly easily from our existing results— see Appendix E), more general approximations are an interesting area for exploration, including those for more general data models.

ACKNOWLEDGMENTS

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

REFERENCES

- [1] Parag Agrawal, Omar Benjelloun, Anish Das Sarma, Chris Hayworth, Shubha U. Nabar, Tomoe Sugihara, and Jennifer Widom. 2006. Trio: A System for Data, Uncertainty, and Lineage. In *VLDB*. 1151–1154.
- [2] Antoine Amarilli, Pierre Bourhis, and Pierre Senellart. 2015. Probabilities and provenance via tree decompositions. *PODS* (2015).
- [3] Yael Amsterdamer, Daniel Deutch, and Val Tannen. 2011. Provenance for Aggregate Queries. In *PODS*. 153–164.
- [4] Lyublena Antova, Thomas Jansen, Christoph Koch, and Dan Olteanu. 2008. Fast and Simple Relational Processing of Uncertain Data. In *2008 IEEE 24th International Conference on Data Engineering*. 983–992. <https://doi.org/10.1109/ICDE.2008.4497507>
- [5] Albert Atserias, Martin Grohe, and Dániel Marx. 2013. Size Bounds and Query Plans for Relational Joins. *SIAM J. Comput.* 42, 4 (2013), 1737–1767. <https://doi.org/10.1137/110859440>
- [6] R. Iris Bahar, Erica A. Frohm, Charles M. Gaona, Gary D. Hachtel, Enrico Macii, Abelardo Pardo, and Fabio Somenzi. 1993. Algebraic decision diagrams and their applications. In *IEEE CAD*.
- [7] George Beskales, Ihab F. Ilyas, and Lukasz Golab. 2010. Sampling the Repairs of Functional Dependency Violations under Hard Constraints. *Proc. VLDB Endow.* 3, 1 (2010), 197–207.
- [8] Peter Bürgisser, Michael Clausen, and Mohammad Amin Shokrollahi. 1997. *Algebraic complexity theory*. Vol. 315. Springer.
- [9] Hubie Chen and Martin Grohe. 2010. Constraint Satisfaction With Succinctly Specified Relations. *J. Comput. Syst. Sci.* 76, 8 (2010), 847–860.
- [10] Jianer Chen, Xiuzhen Huang, Iyad A. Kanj, and Ge Xia. 2006. Strong computational lower bounds via parameterized complexity. *J. Comput. System Sci.* 72, 8 (2006), 1346–1367. <https://doi.org/10.1016/j.jcss.2006.04.007>
- [11] Radu Curticapean. 2013. Counting Matchings of Size k Is $W[1]$ -Hard. In *ICALP*, Vol. 7965. 352–363.
- [12] Radu Curticapean and Dániel Marx. 2014. Complexity of counting subgraphs: only the boundedness of the vertex-cover number counts. *CoRR abs/1407.2929* (2014). [arXiv:1407.2929](http://arxiv.org/abs/1407.2929) <http://arxiv.org/abs/1407.2929>
- [13] Nilesh Dalvi and Dan Suciu. 2007. The Dichotomy of Conjunctive Queries on Probabilistic Structures. In *PODS*. 293–302.
- [14] N. Dalvi and D. Suciu. 2007. Efficient query evaluation on probabilistic databases. *VLDB* 16, 4 (2007), 544.
- [15] Nilesh Dalvi and Dan Suciu. 2012. The dichotomy of probabilistic inference for unions of conjunctive queries. *JACM* 59, 6 (2012), 30.
- [16] Maarten Van den Heuvel, Peter Ivanov, Wolfgang Gatterbauer, Floris Geerts, and Martin Theobald. 2019. Anytime Approximation in Probabilistic Databases via Scaled Dissociations. In *SIGMOD*. 1295–1312.
- [17] Daniel Deutch, Tova Milo, Sudeepa Roy, and Val Tannen. 2014. Circuits for Datalog Provenance. In *ICDT*. 201–212.
- [18] Su Feng, Boris Glavic, Aaron Huber, and Oliver Kennedy. 2021. Efficient Uncertainty Tracking for Complex Queries with Attribute-level Bounds. In *SIGMOD*.
- [19] Su Feng, Aaron Huber, Boris Glavic, and Oliver Kennedy. 2019. Uncertainty Annotated Databases - A Lightweight Approach for Approximating Certain Answers. In *SIGMOD*.
- [20] Robert Fink, Larisa Han, and Dan Olteanu. 2012. Aggregation in probabilistic databases via knowledge compilation. *PVLDB* 5, 5 (2012), 490–501.
- [21] Robert Fink, Jiewen Huang, and Dan Olteanu. 2013. Anytime approximation in probabilistic databases. *VLDBJ* 22, 6 (2013), 823–848.
- [22] Robert Fink and Dan Olteanu. 2016. Dichotomies for Queries with Negation in Probabilistic Databases. *TODS* 41, 1 (2016), 4:1–4:47.
- [23] Jörg Flum and Martin Grohe. 2002. The Parameterized Complexity of Counting Problems. In *Proceedings of the 43rd Symposium on Foundations of Computer Science (FOCS '02)*. IEEE Computer Society, USA, 538.
- [24] Jörg Flum and Martin Grohe. 2006. *Parameterized Complexity Theory*. Springer. <https://doi.org/10.1007/3-540-29953-X>
- [25] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. 2009. *Database systems - the complete book (2. ed.)*. Pearson Education.
- [26] Todd J. Green, Gregory Karvounarakis, and Val Tannen. 2007. Provenance semirings. In *PODS*. 31–40.
- [27] Todd J Green and Val Tannen. 2006. Models for incomplete and probabilistic information. In *EDBT*. 278–296.
- [28] Martin Grohe, Peter Lindner, and Christoph Standke. 2022. Probabilistic Query Evaluation with Bag Semantics. <https://doi.org/10.48550/ARXIV.2201.11524>
- [29] T. Imieliński and W. Lipski. 1989. Incomplete Information in Relational Databases.
- [30] Tomasz Imieliński and Witold Lipski Jr. 1984. Incomplete Information in Relational Databases. *JACM* 31, 4 (1984), 761–791.
- [31] Abhay Kumar Jha and Dan Suciu. 2012. Probabilistic Databases With Markovviews. *PVLDB* 5, 11 (2012), 1160–1171.
- [32] Richard M. Karp, Michael Luby, and Neal Madras. 1989. Monte-Carlo Approximation Algorithms for Enumeration Problems. *J. Algorithms* 10, 3 (1989), 429–448.
- [33] Batya Kenig, Avigdor Gal, and Ofer Strichman. 2013. A New Class of Lineage Expressions over Probabilistic Databases Computable in P-Time. In *SUM*, Vol. 8078. 219–232.
- [34] Mahmoud Abokhamis, Hung Q. Ngo, and Atri Rudra. 2016. FAQ: Questions Asked Frequently. In *PODS*. 13–28.
- [35] Tsvi Kopelowitz and Virginia Vassilevska Williams. 2020. Towards Optimal Set-Disjointness and Set-Intersection Data Structures. In *ICALP*, Vol. 168. 74:1–74:16.
- [36] Poonam Kumari, Said Achmiz, and Oliver Kennedy. 2016. Communicating Data Quality in On-Demand Curation. In *QDB*.
- [37] Hung Q. Ngo. 2018. Worst-Case Optimal Join Algorithms: Techniques, Results, and Open Problems. In *PODS*.
- [38] Hung Q. Ngo, Christopher Ré, and Atri Rudra. 2013. Skew strikes back: new developments in the theory of join algorithms. *SIGMOD Rec.* 42, 4 (2013), 5–16.
- [39] Dan Olteanu, Jiewen Huang, and Christoph Koch. 2010. Approximate confidence computation in probabilistic databases. In *ICDE*. 145–156.
- [40] Dan Olteanu and Maximilian Schleich. 2016. Factorized Databases. *SIGMOD Rec.* 45, 2 (2016), 5–16.
- [41] pdbench utility 2008. *pdbench*. <http://pdbench.sourceforge.net/> Accessed: 2020-12-15.
- [42] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *Proc. VLDB Endow.* 10, 11 (2017), 1190–1201.
- [43] Sudeepa Roy, Vittorio Perduca, and Val Tannen. 2011. Faster query answering in probabilistic databases using read-once functions. In *ICDT*.
- [44] C. Ré and D. Suciu. 2007. Materialized views in probabilistic databases: for information exchange and query optimization. In *VLDB*. 51–62.
- [45] Christopher De Sa, Alexander Ratner, Christopher Ré, Jaeho Shin, Feiran Wang, Sen Wu, and Ce Zhang. 2017. Incremental knowledge base construction using DeepDive. *VLDB J.* 26, 1 (2017), 81–105.
- [46] Prithviraj Sen, Amol Deshpande, and Lise Getoor. 2010. Read-Once Functions and Query Evaluation in Probabilistic Databases. *PVLDB* 3, 1 (2010), 1068–1079.
- [47] Pierre Senellart. 2018. Provenance and Probabilities in Relational Databases. *SIGMOD Record* 46, 4 (2018), 5–15.
- [48] Leslie G. Valiant. 1979. The Complexity of Enumeration and Reliability Problems. *SIAM J. Comput.* 8, 3 (1979), 410–421.
- [49] Guy Van den Broeck and Dan Suciu. 2017. Query Processing on Probabilistic Data: A Survey. (2017).
- [50] Virginia Vassilevska Williams. 2018. Some Open Problems in Fine-Grained Complexity. *SIGACT News* 49, 4 (2018), 29–35. <https://doi.org/10.1145/3300150.3300158>
- [51] Ying Yang, Niccolò Meneghetti, Ronny Fehling, Zhen Hua Liu, Dieter Gawlick, and Oliver Kennedy. 2015. Lenses: An On-Demand Approach to ETL. *PVLDB* 8, 12 (2015), 1578–1589.

A MISSING DETAILS FROM SECTION 2

A.1 Background details for proof of Proposition 2.8

A.1.1 \mathcal{K} -relations and $\mathbb{N}[\mathbf{X}]$ -encoded PDBs. We can use \mathcal{K} -relations to model bags. A \mathcal{K} -relation [26] is a relation whose tuples are annotated with elements from a commutative semiring $\mathcal{K} = \{K, \oplus_{\mathcal{K}}, \otimes_{\mathcal{K}}, \mathbb{0}_{\mathcal{K}}, \mathbb{1}_{\mathcal{K}}\}$. A commutative semiring is a structure with a domain K and associative and commutative binary operations $\oplus_{\mathcal{K}}$ and $\otimes_{\mathcal{K}}$ such that $\otimes_{\mathcal{K}}$ distributes over $\oplus_{\mathcal{K}}$, $\mathbb{0}_{\mathcal{K}}$ is the identity of $\oplus_{\mathcal{K}}$, $\mathbb{1}_{\mathcal{K}}$ is the identity of $\otimes_{\mathcal{K}}$, and $\mathbb{0}_{\mathcal{K}}$ annihilates all elements of K when combined by $\otimes_{\mathcal{K}}$. Let \mathcal{U} be a countable domain of values. Formally, an n -ary \mathcal{K} -relation R over \mathcal{U} is a function $R : \mathcal{U}^n \rightarrow K$ with finite support $\text{supp}(R) = \{t \mid R(t) \neq \mathbb{0}_{\mathcal{K}}\}$. A \mathcal{K} -database is defined similarly, where we view the \mathcal{K} -database (relation) as a function mapping tuples to their respective annotations. \mathcal{RA}^+ query semantics over \mathcal{K} -relations are analogous to the lineage construction semantics of Fig. 1, with the exception of replacing $+$ with $\oplus_{\mathcal{K}}$ and \cdot with $\otimes_{\mathcal{K}}$.

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

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

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

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

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

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

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

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

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

$$\begin{aligned} \Omega &= \{\psi_{\mathbf{W}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) \mid \mathbf{W} \in \{0, 1\}^n\} \\ \forall \mathcal{D} \in \Omega : \text{Pr}(\mathcal{D}) &= \sum_{\mathbf{W} \in \{0, 1\}^n : \psi_{\mathbf{W}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]) = \mathcal{D}} \text{Pr}(\mathbf{W}) \end{aligned}$$

For instance, consider a $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}$ consisting of a single tuple $t_1 = (1)$ annotated with $X_1 + X_2$ with probability distribution $\text{Pr}([0, 0]) = 0$, $\text{Pr}([0, 1]) = 0$, $\text{Pr}([1, 0]) = 0.3$ and $\text{Pr}([1, 1]) = 0.7$. This $\mathbb{N}[\mathbf{X}]$ -encoded 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 $\mathbb{N}[\mathbf{X}]$ -encoded PDB and $\mathbb{N}[\mathbf{X}]$ -encoded PDBs are closed under \mathcal{RA}^+ [26].

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

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

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

The probability distribution \mathcal{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 $\text{Mod}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = \mathcal{D}$. Thus, $\mathbb{N}[\mathbf{X}]$ -encoded PDBs are a complete representation system.

Since $\mathbb{N}[\mathbf{X}]$ is the free object in the variety of semirings, Birkhoff's HSP theorem implies that any assignment $\mathbf{X} \rightarrow \mathbb{N}$, which includes as a special case the assignments $\psi_{\mathbf{W}}$ used here, uniquely extends to the semiring homomorphism alluded to above, $\psi_{\mathbf{W}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) : \mathbb{N}[\mathbf{X}] \rightarrow \mathbb{N}$.

For a polynomial $\psi_{\mathbf{W}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ substitutes variables based on \mathbf{W} and then evaluates the resulting expression in \mathbb{N} . For instance, consider the polynomial $\mathcal{D}_{\mathbb{N}[\mathbf{X}]}(t) = \Phi = X + Y$ and assignment $\mathbf{W} := X = 0, Y = 1$. We get $\psi_{\mathbf{W}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]}) = 0 + 1 = 1$. Closure under \mathcal{RA}^+ queries follows from this and from [26]'s Proposition 3.5, which states that semiring homomorphisms commute with queries over \mathcal{K} -relations. \square

A.1.2 TIDBs and BIDBs in the $\mathbb{N}[\mathbf{X}]$ -encoded PDB model. Two important subclasses of $\mathbb{N}[\mathbf{X}]$ -encoded PDBs that are of interest to us are the bag versions of tuple-independent databases (TIDBs) and block-independent databases (BIDBs). Under set semantics, a TIDB is a deterministic database D where each tuple t is assigned a probability p_t . The set of possible worlds represented by a TIDB D is all subsets of D . The probability of each world is the product of the probabilities of all tuples that exist with one minus the probability of all tuples of D that are not part of this world, i.e., tuples are treated as independent random events. In a BIDB, we also assign each tuple a probability, but additionally partition D into blocks. The possible worlds of a BIDB D are all subsets of D that contain at most one tuple from each block. Note then that the tuples sharing the same block are disjoint, and the sum of the probabilities of all the tuples in the same block B is at most 1. The probability of such a world is the product of the probabilities of all tuples present in the world and the product of the probabilities that no tuple is present in each block B for which no tuple exists in that world. For bag TIDBs and BIDBs, we define the probability of a tuple to be the probability that the tuple exists with multiplicity at least 1.

In this work, we define TIDBs and BIDBs as subclasses of $\mathbb{N}[\mathbf{X}]$ -encoded PDBs defined over variables \mathbf{X} (Definition A.2) where \mathbf{X} can be partitioned into blocks that satisfy the conditions of a TIDB or BIDB (stated formally in Sec. 2.2). In this work, we consider one further deviation from the standard: We use bag semantics for queries. Even though tuples cannot occur more than once in the input TIDB or BIDB, they can occur with a multiplicity larger than one in the result of a query. Since in TIDBs and BIDBs, there is a one-to-one correspondence between tuples in the database and variables, we can interpret a vector $\mathbf{W} \in \{0, 1\}^n$ as denoting which tuples exist in the possible world $\psi_{\mathbf{W}}(\mathcal{D}_{\mathbb{N}[\mathbf{X}]})$ (the ones where $W_i = 1$). For BIDBs specifically, note that at most one of the bits corresponding to tuples in each of the m blocks will be set (i.e., for any pair of bits $W_j, W_{j'}$ that are part of the same block $B_i \supseteq \{t_j, t_{j'}\}$, at most one of them will be set). Denote the vector \mathbf{p} to be a vector whose elements are the individual probabilities p_i of each tuple t_i . Given PDB \mathcal{D} where \mathcal{P} is the distribution induced by \mathbf{p} , which we will denote $\mathcal{P}(\mathbf{p})$.

$$\mathbb{E}_{\mathbf{W} \sim \mathcal{P}(\mathbf{p})} [\Phi(\mathbf{W})] = \sum_{\substack{\mathbf{W} \in \{0,1\}^n \\ \text{s.t. } W_j, W_{j'}=1 \rightarrow \nexists B_i \supseteq \{t_j, t_{j'}\}}} \Phi(\mathbf{W}) \prod_{\substack{j \in [n] \\ \text{s.t. } W_j=1}} p_j \prod_{\substack{i \in [m] \text{ s.t.} \\ \forall t_j \in B_i, W_j=0}} \left(1 - \sum_{t_j \in B_i} p_j\right) \quad (5)$$

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

A.2 Proof of Proposition 2.4

PROOF OF PROPOSITION 2.4. We first need to prove that any c -TIDB \mathcal{D} can be reduced to the Binary-BIDB created by Proposition 2.4. By definition, any $t \in D$ can be present $c' \in [c]$ times in the possible worlds it appears in, with a disjoint probability distribution across the multiplicities $[c]$. Then the construction of Proposition 2.4 of the block of tuples $\{\langle t, j \rangle_{j \in [c]}\}$ for each $t \in D$ indeed encodes the disjoint behavior across multiplicities. \mathcal{D} further requires that all $t \in D$ are independent, a property which is enforced by independence constraint across all blocks of tuples in Definition 2.3. Then the construction of \mathcal{D}' in Proposition 2.4 is an equivalent representation of \mathcal{D} .

Next we need to show that the distributions over \mathcal{D} and \mathcal{D}' are equivalent. The distribution \mathcal{P} is a distribution disjoint across the set of multiplicities $[c]$ and independent across all $t \in D$. By definition, Proposition 2.4 creates \mathcal{D}' by producing a block of disjoint tuples $t_j = \langle t, j \rangle$ for each $t \in D$ and $j \in [c]$, where $\mathbf{p} = (p_{t,j})_{t \in D, j \in [c]}$. Since the probability vector $\mathbf{p}_{\mathcal{D}} = (p_{t,j})_{t \in D, j \in [c]}$ and each $p_{t,j}, p_{t,j'}$ for $j \neq j' \in [c]$ are disjoint, the distributions are hence the same.

A.3 Proof of Proposition 2.8

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

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

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

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

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

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

By distributivity of $+$ over \times

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

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

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

□

A.4 Proposition A.4

Note the following fact:

PROPOSITION A.4. For any BIDB-lineage polynomial $\Phi(X_1, \dots, X_n)$ and all \mathbf{W} such that $Pr[\mathbf{W}] > 0$, it holds that $\Phi(\mathbf{W}) = \tilde{\Phi}(\mathbf{W})$.

PROOF. Note that any Φ in factorized form is equivalent to its SMB expansion. For each term in the expanded form, further note that for all $b \in \{0, 1\}$ and all $e \geq 1$, $b^e = b$. By definition (see Definition 2.5), $\tilde{\Phi}(\mathbf{X})$ is the SMB expansion of $\Phi(\mathbf{X})$ followed by reducing every exponent $e > 1$ to 1 and eliminating all cross terms for the BIDB case. Note that it must be that no cross terms exist in $\Phi(\mathbf{X})$, since by the proposition statement, $Pr[\mathbf{W}] > 0$. Thus, since all monomials are indeed the same, it follows that $\Phi(\mathbf{W}) = \tilde{\Phi}(\mathbf{W})$. □

A.5 Proof for Lemma 2.6

PROOF. Let Φ be a polynomial of n variables with highest degree = K , defined as follows:

$$\Phi(X_1, \dots, X_n) = \sum_{\mathbf{d} \in \{0, \dots, K\}^n} c_{\mathbf{d}} \cdot \prod_{i=1}^n X_i^{d_i},$$

where D' has n tuples, we can equivalently write $\prod_{t \in D'} X_t^{d_t}$ for the product term.

Let the boolean function $\text{ISIND}(\cdot)$ take \mathbf{d} as input and return true if there does not exist any dependent variables in the monomial encoded by \mathbf{d} , i.e., for any block $B \in D'$, $\nexists t, t' \in B \mid \mathbf{d}_t, \mathbf{d}_{t'} \geq 1$.

Then, given Binary-BIDB \mathcal{D} , query Q , and polynomial $\Phi(\mathbf{W}) = \Phi[Q, D, t]$, in expectation we have

$$\mathbb{E}_{\mathbf{W}} [\Phi(\mathbf{W})] = \mathbb{E}_{\mathbf{W}} \left[\sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{t \in D'} W_t^{d_t} + \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \neg \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{t \in D'} W_t^{d_t} \right] \quad (6)$$

$$= \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[\prod_{t \in D'} W_t^{d_t} \right] + \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \neg \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[\prod_{t \in D'} W_t^{d_t} \right] \quad (7)$$

$$= \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \mathbb{E}_{\mathbf{W}} \left[\prod_{t \in D'} W_t^{d_t} \right] \quad (8)$$

$$= \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{t \in D'} \mathbb{E}_{\mathbf{W}} [W_t^{d_t}] \quad (9)$$

$$= \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{t \in D'} \mathbb{E}_{\mathbf{W}} [W_t] \quad (10)$$

$$= \sum_{\substack{\mathbf{d} \in \{0, \dots, K\}^n \\ \wedge \text{ISIND}(\mathbf{d})}} c_{\mathbf{d}} \cdot \prod_{t \in D'} p_t \quad (11)$$

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

Eq. (6) is the result of substituting in the definition of Φ given above. Then we arrive at eq. (7) by linearity of expectation. Next, eq. (8) is the result of the independence constraint of BIDBs, specifically that any monomial composed of dependent variables, i.e., variables from the same block B , has a probability of 0. Eq. (9) is obtained by the fact that all variables in each monomial are independent, which allows for the expectation to be pushed through the product. In eq. (10), recall the lineage construction semantics of Fig. 3 for a Binary-BIDB, where the annotation of a tuple t with multiplicity j is written as $j^{d_t} \cdot X_t^{d_t}$ such that $\text{DOM}(\mathbf{W}_t) = \{0, 1\}$. Then $c_{\mathbf{d}}$ absorbs all such j^{d_t} factors. Since $W_t \in \{0, 1\}$ it is the case that for any exponent $e \geq 1$, $W_t^e = W_t$. Next, in eq. (11) the expectation of a tuple is indeed its probability.

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

A.6 Proof For Corollary 2.7

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

B MISSING DETAILS FROM SECTION 3

B.1 Lemma B.1

LEMMA B.1. *Assuming that each $v \in V$ has degree ≥ 1 ,¹³ the PDB relations encoding the edges for Φ_G^k of Definition 3.4 can be computed in $O(m)$ time.*

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

B.2 Proof of Lemma 3.5

PROOF. By the recursive definition of $T_{det}(\cdot, \cdot)$ (see Sec. 2.4), we have the following equation for our hard query Q when $k = 1$, (we denote this as Q^1).

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

We argue that $T_{join}(D.V, D.E, D.V)$ is at most $O(m)$ by noting that there exists an algorithm that computes $D.V \bowtie D.E \bowtie D.V$ in the same runtime¹⁴. Then by the assumption of Lemma B.1 (each $v \in V$ has degree ≥ 1), the sum of the first three terms is $O(m)$. We then obtain that $T_{det}(Q^1, D) = O(m) + O(m) = O(m)$. For $Q^k = Q_1^1 \times \dots \times Q_k^1$, we have the recurrence $T_{det}(Q^k, D) = T_{det}(Q_1^1, D) + \dots + T_{det}(Q_k^1, D) + T_{join}(Q_1^1, \dots, Q_k^1)$. Since Q^1 outputs a count, computing the join $Q_1^1 \bowtie \dots \bowtie Q_k^1$ is just multiplying k numbers, which takes $O(k)$ time. Thus, we have

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

as desired. \square

B.3 Lemma B.2

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

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

¹³This is WLOG, since any vertex with degree 0 can be dropped without affecting the result of our hard query.

¹⁴Indeed the trivial algorithm that computes the obvious pair-wise joins has the claimed runtime. That is, we first compute $D.V \bowtie D.E$, which takes $O(m)$ (assuming $D.V$ is stored in hash map) since tuples in $D.V$ can only filter tuples in $D.E$. The resulting subset of tuples in $D.E$ are then again joined (on the right) with $D.V$, which by the same argument as before also takes $O(m)$ time, as desired.

B.4 Proof of Lemma B.2

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

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

We note that c_i is *exactly* the number of monomials in the SMB expansion of $\Phi_G^k(\mathbf{X})$ composed of i distinct variables.¹⁵

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

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

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

where a unique k -matching in the multi-set of product terms has k distinct (i_ℓ, j_ℓ) index pairs. Further, any monomial composed of such a distinct set of $2k$ variables will be produced $k!$ times in a k -wise product of the sum of a set of elements with itself. This is true because each (identical) product term contains each of the k distinct elements, giving us $k \cdot k - 1 \cdot \dots \cdot 1 = k!$ permutations of a distinct k -matching.

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

B.5 Proof of Theorem 3.6

PROOF. For the sake of contradiction, assume we can solve our problem in $o(T_{\text{match}}(k, G))$ time. Given a graph G by Lemma B.1 we can compute the PDB encoding in $O(m)$ time. Then after we run our algorithm on $\tilde{\Phi}_G^k$, we get $\tilde{\Phi}_G^k(p_i, \dots, p_i)$ for every $0 \leq i \leq 2k$ in additional $O(k) \cdot o(T_{\text{match}}(k, G))$ time. Lemma B.2 then computes the number of k -matchings in G in $O(k^3)$ time. Adding the runtime of all of these steps, we have an algorithm for computing the number of k -matchings that runs in time

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

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

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

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

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

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

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

$$\#(G, \mathfrak{g}) = m, \tag{15}$$

¹⁵Since $\tilde{\Phi}_G^k(\mathbf{X})$ does not have any monomial with degree < 2 , it is the case that $c_0 = c_1 = 0$ but for the sake of simplicity we will ignore this observation.

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

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

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

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

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

B.7 Proofs of Eq. (15)-Eq. (20)

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

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

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

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

B.8 Tools to prove Theorem 3.7

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

LEMMA B.3. *For any p , we have:*

$$\begin{aligned} \tilde{\Phi}_G^3(p, \dots, p) = & \#(G, \mathfrak{B}) p^2 + 6\#(G, \mathfrak{A}) p^3 + 6\#(G, \mathfrak{C}) p^4 + 6\#(G, \mathfrak{G}) p^3 \\ & + 6\#(G, \mathfrak{D}) p^4 + 6\#(G, \mathfrak{E}) p^4 + 6\#(G, \mathfrak{F}) p^5 + 6\#(G, \mathfrak{H}) p^6. \end{aligned} \quad (21)$$

B.8.1 Proof for Lemma B.3.

PROOF. By definition we have that

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

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

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

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

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

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

DEFINITION B.4. For $\ell \geq 1$, let graph $G^{(\ell)}$ be a graph generated from an arbitrary graph G , by replacing every edge e of G with an ℓ -path, such that all inner vertices of an ℓ -path replacement edge have degree 2.¹⁶

We will prove Theorem 3.7 by the following reduction:

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

For clarity, we repeat the notion of $\#(G, H)$ to mean the count of subgraphs in G isomorphic to H . The following lemmas relate these counts in $G^{(2)}$ to $G^{(1)}$ and are useful in proving Lemma B.8.

LEMMA B.6. The 3-matchings in graph $G^{(2)}$ satisfy the identity:

$$\begin{aligned} \#(G^{(2)}, \mathfrak{I}) &= 8 \cdot \#(G^{(1)}, \mathfrak{I}) + 6 \cdot \#(G^{(1)}, \mathfrak{A}) \\ &\quad + 4 \cdot \#(G^{(1)}, \mathfrak{B}) + 4 \cdot \#(G^{(1)}, \mathfrak{C}) + 2 \cdot \#(G^{(1)}, \mathfrak{D}). \end{aligned}$$

LEMMA B.7. For $\ell > 1$ and any graph $G^{(\ell)}$, $\#(G^{(\ell)}, \mathfrak{A}) = 0$.

Finally, the following result immediately implies Theorem B.5:

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

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

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

B.9 Proofs for Lemma B.6, Lemma B.7, and Lemma B.8

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

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

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

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

DEFINITION B.11 ($\binom{E}{t}$ AND $\binom{E}{\leq t}$). Let $\binom{E}{t}$ denote the set of subsets in E with exactly t edges. In a similar manner, $\binom{E}{\leq t}$ is used to mean the subsets of E with t or fewer edges.

The following function f_t is a mapping from every 3-edge shape in $G^{(\ell)}$ to its 'projection' in $G^{(1)}$.

¹⁶Note that $G \equiv G^{(1)}$.

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

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

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

Note, importantly, that when we discuss f_ℓ^{-1} , that each edge present in $E_S^{(1)}$ must have an edge in each $s \in f_\ell^{-1}(E_S^{(1)})$ that projects down to it. In particular, if $|E_S^{(1)}| = 3$, then it must be the case that each $s \in f_\ell^{-1}(E_S^{(1)})$ consists of the following set of edges: $\{(e_i, b_1), (e_j, b_2), (e_m, b_3)\}$, where i, j and m are distinct.

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

B.9.1 Proof of Lemma B.6.

PROOF. For each subset $E_S^{(1)} \in \binom{E_1}{\leq 3}$, we count the number of 3-matchings in the 3-edge subgraphs of $G^{(2)}$ in $f_2^{-1}(E_S^{(1)})$. Denote the set of 3-matchings in $f_2^{-1}(E_S^{(1)})$ as $f_{2, \mathfrak{B}}^{-1}(E_S^{(1)})$ and its size as $|f_{2, \mathfrak{B}}^{-1}(E_S^{(1)})|$. We first consider the case of $E_S^{(1)} \in \binom{E_1}{3}$, where $E_S^{(1)}$ is composed of the edges e_1, e_2, e_3 and $f_2^{-1}(E_S^{(1)})$ is the set s of all 3-edge subsets of edges $e_i \in \{(e_1, 0), (e_1, 1), (e_2, 0), (e_2, 1), (e_3, 0), (e_3, 1)\}$ such that $f_\ell(s) = \{e_1, e_2, e_3\}$. For the case that $f_2^{-1}(E_S^{(1)}) = f_{2, \mathfrak{B}}^{-1}(E_S^{(1)})$ we have that $|f_{2, \mathfrak{B}}^{-1}(E_S^{(1)})| = 8$.

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

- 3-matching (\mathfrak{B})

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

- Disjoint Two-Path (\mathfrak{A})

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

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

Since these two selections can be made independently, $|f_{2, \mathfrak{A}}^{-1}(E_S^{(1)})| = 2 \cdot 3 = 6$.

- 3-star (\mathfrak{A})

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

- 3-path (\mathfrak{B})

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

- Triangle (\mathfrak{A})

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

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

- 2-matching (\mathfrak{B}), 2-path (\mathfrak{A}), 1 edge (\mathfrak{B})

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

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

B.9.2 Proof of Lemma B.7.

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

B.9.3 Proof of Lemma B.8.

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

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

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

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

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

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

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

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

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

$$\begin{aligned} & 0 - (8\#(G, \mathfrak{D}) + 6\#(G, \mathfrak{E}) + 4\#(G, \mathfrak{F}) + 4\#(G, \mathfrak{G}) + 2\#(G, \mathfrak{H})) (3p^2-p^3) = \\ & \frac{\tilde{\Phi}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{I})}{6p} - \#(G^{(2)}, \mathfrak{E}) - \#(G^{(2)}, \mathfrak{F}) p - \#(G^{(2)}, \mathfrak{G}) p \\ & - [\#(G^{(2)}, \mathfrak{B}) p^2 + 3\#(G^{(2)}, \mathfrak{D}) p^2] - [\#(G^{(2)}, \mathfrak{H}) p + 3\#(G^{(2)}, \mathfrak{A}) p] \end{aligned} \quad (25)$$

$$\begin{aligned} & (10\#(G, \mathfrak{A}) + 10\#(G, \mathfrak{D})) (3p^2-p^3) = \\ & \frac{\tilde{\Phi}_{G^{(2)}}^3(p, \dots, p)}{6p^3} - \frac{\#(G^{(2)}, \mathfrak{I})}{6p} - \#(G^{(2)}, \mathfrak{E}) - \#(G^{(2)}, \mathfrak{F}) p - \#(G^{(2)}, \mathfrak{G}) p \\ & - [\#(G^{(2)}, \mathfrak{H}) p + 3\#(G^{(2)}, \mathfrak{A}) p] - [\#(G^{(2)}, \mathfrak{B}) p^2 - 3\#(G^{(2)}, \mathfrak{D}) p^2] \\ & + (4\#(G, \mathfrak{F}) + [6\#(G, \mathfrak{E}) + 18\#(G, \mathfrak{D})] + [4\#(G, \mathfrak{H}) + 12\#(G, \mathfrak{A})]) (3p^2-p^3) \end{aligned} \quad (26)$$

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

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

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

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

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

B.10 Proof of Theorem B.5

PROOF. We can compute $G^{(2)}$ from $G^{(1)}$ in $O(m)$ time. Additionally, if in time $O(T(m))$, we have $\widetilde{\Phi}_{G^{(\ell)}}^3(p, \dots, p)$ for $\ell \in [2]$, then the theorem follows by Lemma B.8. \square

B.11 Proof of Theorem 3.7

PROOF. For the sake of contradiction, assume that for any G , we can compute $\widetilde{\Phi}_G^3(p, \dots, p)$ in $o(m^{1+\epsilon_0})$ time. Let G be the input graph. Then by Theorem B.5 we can compute $\#(G, \mathfrak{L})$ in further time $o(m^{1+\epsilon_0}) + O(m)$. Thus, the overall, reduction takes $o(m^{1+\epsilon_0}) + O(m) = o(m^{1+\epsilon_0})$ time, which violates Conjecture 3.3. \square

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

C MISSING DETAILS FROM SECTION 4

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

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

DEFINITION C.1 (PURE EXPANSION). *The pure expansion of a polynomial Φ is formed by computing all product of sums occurring in Φ , without combining like monomials. The pure expansion of Φ generalizes Definition 2.1 by allowing monomials $m_i = m_j$ for $i \neq j$.*

Note that similar in spirit to Definition 2.5, $E(C)$ Definition 4.1 reduces all variable exponents $e > 1$ to $e = 1$. Further, it is true that $E(C)$ encodes the pure expansion of C .

EXAMPLE C.2 (EXAMPLE OF PURE EXPANSION). *Consider the factorized representation $(X + 2Y)(2X - Y)$ of the polynomial in Eq. (28). Its circuit C is illustrated in Fig. 4. The pure expansion of the product is $2X^2 - XY + 4XY - 2Y^2$. As an additional example of Definition 4.1, $E(C) = [(X, 2), (XY, -1), (XY, 4), (Y, -2)]$.*

$E(C)$ effectively¹⁷ encodes the *reduced* form of $\text{POLY}(C)$, decoupling each monomial into a set of variables v and a real coefficient c . However, unlike the constraint on the input Φ to compute $\widetilde{\Phi}$, the input circuit C does not need to be in SMB/SOP form.

EXAMPLE C.3 (EXAMPLE FOR DEFINITION 4.2). *Using the same factorization from Example C.2, $\text{POLY}(|C|) = (X + 2Y)(2X + Y) = 2X^2 + XY + 4XY + 2Y^2 = 2X^2 + 5XY + 2Y^2$. Note that this is not the same as the polynomial from Eq. (28). As an example of the slight abuse of notation we alluded to in Sec. 4, $|C|(1, \dots, 1) = 2(1)^2 + 5(1)(1) + 2(1)^2 = 9$.*

DEFINITION C.4 (SUBCIRCUIT). *A subcircuit of a circuit C is a circuit S such that S is a DAG subgraph of the DAG representing C . The sink of S has exactly one gate g .*

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

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

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

such that

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

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

C.1 Proof of Theorem C.5

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

¹⁷The minor difference here is that $E(C)$ encodes the *reduced* form of the pure expansion of the compressed representation, as opposed to the SMB representation

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

```

1:  $\text{acc} \leftarrow 0$ 
2:  $N \leftarrow \left\lceil \frac{2 \log \frac{2}{\delta}}{\epsilon^2} \right\rceil$ 
3:  $(C_{\text{mod}}, \text{size}) \leftarrow \text{ONEPASS}(C)$  ▷ ONEPASS is Algorithm 2
4: for  $i \in 1$  to  $N$  do ▷ Perform the required number of samples
5:    $(M, \text{sgn}_i) \leftarrow \text{SAMPLEMONOMIAL}(C_{\text{mod}})$  ▷ 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 ▷  $M$  is the sampled monomial's set of variables (ref. appendix C.9)
7:      $Y_i \leftarrow \prod_{X_j \in M} p_j$ 
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  $\text{acc}$ 

```

LEMMA C.6. *The ONEPASS function completes in time:*

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

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

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

Aaron says: I don't think the word *only* is needed.

LEMMA C.7. *The function SAMPLEMONOMIAL completes in time*

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

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

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

THEOREM C.8. *For any C with $\text{DEG}(\text{POLY}(|C|)) = k$, algorithm 1 outputs an estimate acc of $\tilde{\Phi}(p_1, \dots, p_n)$ such that*

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

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

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

PROOF OF THEOREM C.5. Set $\mathcal{E} = \text{APPROXIMATE}\tilde{\Phi}(C, (p_1, \dots, p_n), \delta, \epsilon')$, where

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

which achieves the claimed error bound on \mathcal{E} (acc) trivially due to the assignment to ϵ' and theorem C.8, since $\epsilon' \cdot |C|(1, \dots, 1) = \epsilon \cdot \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|C|(1, \dots, 1)} \cdot |C|(1, \dots, 1) = \epsilon \cdot \tilde{\Phi}(p_1, \dots, p_n)$.

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

$$\begin{aligned} \frac{1}{(\epsilon')^2} \cdot \log\left(\frac{1}{\delta}\right) &= \frac{\log \frac{1}{\delta}}{\epsilon^2 \left(\frac{\tilde{\Phi}(p_1, \dots, p_N)}{|C|(1, \dots, 1)}\right)^2} \\ &= \frac{\log \frac{1}{\delta} \cdot |C|^2(1, \dots, 1)}{\epsilon^2 \cdot \tilde{\Phi}^2(p_1, \dots, p_n)}. \end{aligned}$$

□

Let us now prove Theorem C.8:

C.2 Proof of Theorem C.8

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

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

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

$$\begin{aligned} \mathbb{E}[Y_i] &= \sum_{(v,c) \in E(C)} \frac{\mathbb{1}_{(\text{ISIND}(v_m))} \cdot c \cdot \prod_{X_i \in \text{VAR}(v)} p_i}{|C|(1, \dots, 1)} \\ &= \frac{\tilde{\Phi}(p_1, \dots, p_n)}{|C|(1, \dots, 1)}, \end{aligned}$$

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

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

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

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

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

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

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

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

Aaron says: Why does the \geq sign change to $>$?

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

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

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

C.3 Proof of Theorem 4.7

PROOF. The result follows by first noting that by definition of γ , we have

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

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

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

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

C.4 Proof of Lemma 4.8

PROOF. The circuit \mathcal{C}' is built from \mathcal{C} in the following manner. For each input gate g_i with $g_i.\text{val} = X_t$, replace g_i with the circuit \mathcal{S} encoding the sum $\sum_{j=1}^c j \cdot X_{t,j}$. We argue that \mathcal{C}' is a valid circuit by the following facts. Let $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ be the original c -TIDB \mathcal{C} was generated from. Then, by Proposition 2.4 there exists a Binary-BIDB $\mathcal{D}' = (\times_{t \in D'} \{0, c_t\}, \mathcal{P}')$, with $D' = \{\langle t, j \rangle \mid t \in D, j \in [c]\}$, from which the conversion from \mathcal{C} to \mathcal{C}' follows. Both $\text{POLY}(\mathcal{C})$ and $\text{POLY}(\mathcal{C}')$ have the same expected multiplicity since (by Proposition 2.4) the distributions \mathcal{P} and \mathcal{P}' are equivalent and $\sum_{j=1}^c j \cdot \mathbf{W}'_{t,j} = \mathbf{W}_t$ for $\mathbf{W}' \in \{0, 1\}^{cn}$ and $\mathbf{W} \in \{0, \dots, c\}^D$ such that $\mathbf{W}_t \equiv \mathbf{W}'_t$. Finally, note that because there exists a (sub) circuit encoding $\sum_{j=1}^c j \cdot X_{t,j}$ that is a *balanced* binary tree, the above conversion implies the claimed size and depth bounds of the lemma.

Next we argue the claim on $\gamma(\mathcal{C}')$. Consider the list of expanded monomials $\mathbb{E}(\mathcal{C})$ for c -TIDB circuit \mathcal{C} . Let $v_m = X_{t_1}^{d_1} \cdots X_{t_\ell}^{d_\ell}$ be an arbitrary monomial with ℓ variables and let (abusing notation) $v'_m = \left(\sum_{j=1}^c j \cdot X_{t_1,j}\right)^{d_1} \cdots \left(\sum_{j=1}^c j \cdot X_{t_\ell,j}\right)^{d_\ell}$. Then, for $f_\ell = \sum_{i=1}^\ell d_i$, v_m induces the set of monomials $\left\{\prod_{i=1}^{f_\ell} j_i \cdot X_{t_i,j_i}^{d_i}\right\}_{j_i \in [c]}$ in the pure expansion of v'_m . Recall that a cancellation occurs in v'_m when there exists $t_{i,j} \neq t_{i,j'}$ in the same block B where variables $X_{t_i,j}, X_{t_i,j'}$ are in the set of variables v'_i of v_m . Observe that cancellations can only occur for each $X_{t_i}^{d_i} \in v_m$, where the expansion $\left(\sum_{j=1}^c j \cdot X_{t_i,j}\right)^{d_i}$ represents the monomial $X_{t_i}^{d_i}$ in D' . Consider the number of cancellations for $\left(\sum_{j=1}^c j \cdot X_{t_i,j}\right)^{d_i}$. Then $\gamma \leq 1 - c^{-(d_i-1)}$, since for each element in the set of cross products $\left\{\times_{i \in [d_i], j_i \in [c]} X_{t_i,j_i}\right\}$ there are *exactly* c surviving elements with $j_1 = \dots = j_{d_i} = j$, i.e. $X_{t_i,j}^{d_i}$ for each $j \in [c]$. The rest of the $c^{d_i} - c$ cross terms cancel. Regarding all of v'_m , it is the case that the proportion of non-cancellations for each $\left(\sum_{j=1}^c j \cdot X_{t_i,j}\right)^{d_i} \in v'_m$ multiply because non-cancelling terms for $\left(\sum_{j=1}^c j \cdot X_{t_i,j}\right)^{d_i}$ can only be joined with non-cancelling terms of $\left(\sum_{j=1}^c X_{t_{i'},j}\right)^{d_{i'}} \in v'_m$ for $t \neq t'$. This then yields the fraction of cancelled monomials $\gamma \leq 1 - \prod_{i=1}^\ell c^{-(d_i-1)} \leq 1 - c^{-(k-1)}$ where the inequalities take into account the fact that $f_\ell \leq k$.

Since this is true for arbitrary v , the bound follows for $\text{POLY}(\mathcal{C}')$. \square

C.5 Proof of Lemma 4.9

We will prove Lemma 4.9 by considering the two cases separately. We start by considering the case when \mathcal{C} is a tree:

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

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

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

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

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

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

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

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

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

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

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

The upper bound in Lemma 4.9 for the general case is a simple variant of the above proof (but we present a proof sketch of the bound below for completeness):

LEMMA C.10. *Let C be a (general) circuit. Then we have*

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

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

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

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

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

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

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

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

C.6 ONEPASS Remarks

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

The evaluation of $|C|(1, \dots, 1)$ can be defined recursively, as follows (where C_L and C_R are the 'left' and 'right' inputs of C if they exist):

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

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

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

$$C.\text{Rweight} \leftarrow \frac{|C_R|(1, \dots, 1)}{|C_L|(1, \dots, 1) + |C_R|(1, \dots, 1)} \quad (34)$$

C.7 ONEPASS Example

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

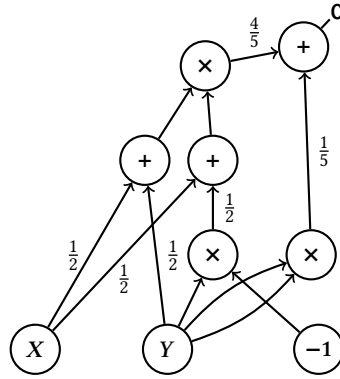


Figure 5: Weights computed by ONEPASS in Example C.11.

Algorithm 2 ONEPASS (C)**Input:** C: Circuit**Output:** C: Annotated Circuit**Output:** $\text{sum} \in \mathbb{N}$

```

1: for g in TOPORD (C) do
2:   if g.type = VAR then
3:     g.partial ← 1
4:   else if g.type = NUM then
5:     g.partial ← |g.val|
6:   else if g.type = × then
7:     g.partial ← gL.partial × gR.partial
8:   else
9:     g.partial ← gL.partial + gR.partial
10:    g.Lweight ←  $\frac{g_L.\text{partial}}{g.\text{partial}}$ 
11:    g.Rweight ←  $\frac{g_R.\text{partial}}{g.\text{partial}}$ 
12:   end if
13:   sum ← g.partial
14: end for
15: return (sum, C)

```

▷ TOPORD (·) is the topological order of C

C.8 Proof of ONEPASS (Lemma C.6)

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

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

For the inductive hypothesis, assume that ONEPASS correctly computes S.partial , S.Lweight , and S.Rweight for all gates g in C with $k \geq 0$ iterations over TOPORD . We now prove for $k + 1$ iterations that ONEPASS correctly computes the partial , Lweight , and Rweight values for each gate g_i in C for $i \in [k + 1]$. The g_{k+1} must be in the last ordering of all gates g_i . When $\text{SIZE}(C) > 1$, if g_{k+1} is a leaf node, we are back to the base case. Otherwise g_{k+1} is an internal node which requires binary input.

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

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

Algorithm 3 SAMPLEMONOMIAL (C)**Input:** C: Circuit**Output:** vars: TreeSet**Output:** sgn $\in \{-1, 1\}$

```

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

```

Runtime Analysis. It is known that $\text{TOPORD}(G)$ is computable in linear time. There are $\text{SIZE}(C)$ iterations. Each iteration has runtime $O\left(\overline{M}(\log(|C(1 \dots, 1)|), \log(\text{SIZE}(C)))\right)$ time. This can be seen since each of all the numbers which the algorithm computes is at most $|C(1 \dots, 1)|$. Hence, by definition each such operation takes $\overline{M}(\log(|C(1 \dots, 1)|), \log \text{SIZE}(C))$ time, which proves the claimed runtime. \square

C.9 SAMPLEMONOMIAL Remarks

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

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

C.10 Proof of SAMPLEMONOMIAL (Lemma C.7)

PROOF. We first need to show that SAMPLEMONOMIAL samples a valid monomial v_m by sampling and returning a set of variables v , such that (v, c) is in $E(C)$ and v_m is indeed a monomial of the $\tilde{\Phi}(X)$ encoded in C. We show this via induction over the depth of C. For the base case, let the depth d of C be 1. We have that the single gate is either a constant c for which by line 15 we return $\{\}$, or we have that C.type = VAR and C.val = x , and by line 17 we return $\{x\}$. By definition 4.1, both cases return a valid v for some (v, c) from $E(C)$, and the base case is proven.

Aaron says: I think it is slightly confusing to say that depth = 0 in view of the definition of depth in S.4. To say $k = 0$ is also strange, since, for a single join, we have that $k = 2$.

For the inductive hypothesis, assume that for $d \leq k$ for some $k \geq 1$, that it is indeed the case that SAMPLEMONOMIAL returns a valid monomial.

For the inductive step, let us take a circuit C with $d = k + 1$. Note that each input has depth $d - 1 \leq k$, and by inductive hypothesis both of them sample a valid monomial. Then the sink can be either a + or \times gate. For the case when C.type = +, line 3 of SAMPLEMONOMIAL will choose one of the inputs of the source. By inductive hypothesis it is the case that some valid monomial is being randomly sampled from each of the inputs. Then it follows when C.type = + that a valid monomial is sampled by SAMPLEMONOMIAL. When the C.type = \times , line 10

computes the set union of the monomials returned by the two inputs of the sink, and it is trivial to see by definition 4.1 that v_m is a valid monomial encoded by some (v, c) of $E(C)$.

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

For the base case $d = 1$, by definition 2.9 we know that the $SIZE(C) = 1$ and $C.type = NUM$ or VAR . For either case, the probability of the value returned is 1 since there is only one value to sample from. When $C.val = x$, the algorithm always return the variable set $\{x\}$. When $C.type = NUM$, $SAMPLEMONOMIAL$ will always return the variable set \emptyset .

Aaron says: I don't think this is technically right, since $SAMPLEMONOMIAL$ returns a tuple of two values.

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

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

Consider the case when $C.type = \times$. For the term (v, c) from $E(C)$ that is being sampled it is the case that $v = v_L \cup v_R$, where v_L is coming from C_L and v_R from C_R . The probability that $SAMPLEMONOMIAL(C_L)$ 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 v_R are sampled with independent randomness, the final probability for sample v is then $\frac{|c_{v_L}| \cdot |c_{v_R}|}{|C_L|(1, \dots, 1) \cdot |C_R|(1, \dots, 1)}$. For (v, c) in $E(C)$, by definition 4.1 it is indeed the case that $|c| = |c_{v_L}| \cdot |c_{v_R}|$ and that (as shown in eq. (32)) $|C|(1, \dots, 1) = |C_L|(1, \dots, 1) \cdot |C_R|(1, \dots, 1)$, and therefore v is sampled with correct probability $\frac{|c|}{|C|(1, \dots, 1)}$.

For the case when $C.type = +$, $SAMPLEMONOMIAL$ will sample v from one of its inputs. By inductive hypothesis we know that any v_L in $E(C_L)$ and any v_R in $E(C_R)$ will 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 equal v , depending on whether C_L or C_R is sampled. Assume that v is sampled from C_L , and note that a symmetric argument holds for the case when v is sampled from C_R . Notice also 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 $ONEPASS$. Then, since $SAMPLEMONOMIAL$ goes top-down, and each sampling choice is independent (which follows from the randomness in the root of C being independent from the randomness used in its subtrees), the probability for v to be sampled from C is equal to the product of the probability that C_L is sampled from C and v is sampled in C_L , and

$$\begin{aligned} Pr(SAMPLEMONOMIAL(C) = v) &= \\ Pr(SAMPLEMONOMIAL(C_L) = v) \cdot Pr(SampledChild(C) = C_L) &= \\ = \frac{|c_v|}{|C_L|(1, \dots, 1)} \cdot \frac{|C_L|(1, \dots, 1)}{|C_L|(1, \dots, 1) + |C_R|(1, \dots, 1)} &= \\ = \frac{|c_v|}{|C|(1, \dots, 1)}, & \end{aligned}$$

and we obtain the desired result.

Lastly, we show by simple induction of the depth d of C that $SAMPLEMONOMIAL$ indeed returns the correct sign value of c in (v, c) .

In the base case, $C.type = NUM$ or VAR . For the former by Line 15, $SAMPLEMONOMIAL$ correctly returns the sign value of the gate. For the latter by Line 17, $SAMPLEMONOMIAL$ returns the correct sign of 1, since a variable is a neutral element, and 1 is the multiplicative identity, whose product with another sign element will not change that sign element.

For the inductive hypothesis, we assume for a circuit of depth $d \leq k$ and $k \geq 1$ that the algorithm correctly returns the sign value of c .

Similar to before, for a depth

Aaron says: Why do we use $d = k + 1$ for the inductive cases above?

$d \leq k + 1$, it is true that C_L and C_R both return the correct sign of c . For the case that $C.type = \times$, the sign value of both inputs are multiplied, which is the correct behavior by definition 4.1. When $C.type = +$, only one input of C is sampled, and the algorithm returns the correct sign value of c by inductive hypothesis.

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

Aaron says: We don't really justify why we can bound the number of recursive calls as we claim in what follows.

Since there can be at most $k = \text{DEG}(C)$ nodes visited at every level of the circuit, and each of the first $d - 1$ levels (going from the sink to the source nodes) will contain at least one recursive call, we can upperbound the number of recursive calls in `SAMPLEMONOMIAL` by $O((\text{DEG}(C) + 1) \cdot \text{DEPTH}(C))$, which by the above will prove the claimed runtime of Lemma C.7.

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

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

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

We prove the following inequality holds.

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

Note that eq. (36) implies the claimed runtime.

Aaron says: If the claimed runtime is from the first paragraph, then I don't follow.

We prove eq. (36) for the number of gates traversed in `SAMPLEMONOMIAL` using induction over $\text{DEPTH}(C)$. Recall how degree is defined in definition 4.4.

Aaron says: In the following, by Definition 4.3, we would have that $\text{DEPTH}(C) = 1$ *technically*.

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

For the inductive hypothesis, we assume the bound holds for any circuit where $\ell \geq \text{DEPTH}(C) \geq 0$. Now consider the case when `SAMPLEMONOMIAL` has an arbitrary circuit C input with $\text{DEPTH}(C) = \ell + 1$. By definition $C.\text{type} \in \{+, \times\}$. Note that since $\text{DEPTH}(C) \geq 2$, C must have input(s). Further we know that by the inductive hypothesis the inputs C_i for $i \in \{L, R\}$ of the sink gate C uphold the bound

$$2(\text{DEG}(C_i) + 1) \cdot \text{DEPTH}(C_i) + 1 \geq \text{COST}(C_i). \quad (37)$$

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

If $C.\text{type} = +$, then $\text{DEG}(C) = \max(\text{DEG}(C_L), \text{DEG}(C_R))$. Otherwise $C.\text{type} = \times$ and $\text{DEG}(C) = \text{DEG}(C_L) + \text{DEG}(C_R) + 1$. In either case it is true that $\text{DEPTH}(C) = \max(\text{DEPTH}(C_L), \text{DEPTH}(C_R)) + 1$.

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

$$2(\text{DEG}(C_L) + \text{DEG}(C_R) + 2) \cdot (\max(\text{DEPTH}(C_L), \text{DEPTH}(C_R)) + 1) + 1 \quad (38)$$

$$\geq 2(\text{DEG}(C_L) + 1) \cdot \text{DEPTH}(C_L) + 2(\text{DEG}(C_R) + 1) \cdot \text{DEPTH}(C_R) + 3 \quad (39)$$

$$\geq 1 + \text{COST}(C_L) + \text{COST}(C_R) = \text{COST}(C). \quad (40)$$

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

$$2\text{DEG}(C_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(C_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(C_L) + 2\text{DEG}(C_R) + 4 + 1 \quad (41)$$

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

$$2\text{DEG}(C_L) \cdot \text{DEPTH}(C_L) + 2\text{DEPTH}(C_L) + 2\text{DEG}(C_R) \cdot \text{DEPTH}(C_R) + 2\text{DEPTH}(C_R) + 3 \quad (42)$$

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

$$\begin{aligned} & 2\text{DEG}(C_L) \cdot \text{DEPTH}_{\max} + 2\text{DEG}(C_R) \cdot \text{DEPTH}_{\max} + 4\text{DEPTH}_{\max} + 2\text{DEG}(C_L) + 2\text{DEG}(C_R) + 5 \\ & \geq 2\text{DEG}(C_L) \cdot \text{DEPTH}(C_L) + 2\text{DEG}(C_R) \cdot \text{DEPTH}(C_R) + 2\text{DEPTH}(C_L) + 2\text{DEPTH}(C_R) + 3. \end{aligned} \quad (43)$$

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

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

$$2(\max(\text{DEG}(C_L), \text{DEG}(C_R)) + 1) \cdot (\max(\text{DEPTH}(C_L), \text{DEPTH}(C_R)) + 1) + 1 \quad (44)$$

$$\geq \max(2(\text{DEG}(C_L) + 1) \cdot \text{DEPTH}(C_L) + 1, 2(\text{DEG}(C_R) + 1) \cdot \text{DEPTH}(C_R) + 1) + 1 \quad (45)$$

$$\geq 1 + \max(\text{COST}(C_L), \text{COST}(C_R)) = \text{COST}(C) \quad (46)$$

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

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

Since $\text{DEG}_{\max} \cdot \text{DEPTH}_{\max} \geq \text{DEG}(C_i) \cdot \text{DEPTH}(C_i)$, the following upperbounds the expansion of eq. (45):

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

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

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

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

Similar to the case of $C.\text{type} = \times$, (46) follows by equations (35) and (37).

This proves (36) as desired. \square

C.11 Experimental Results

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

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

For the data we used the MayBMS data generator [41] tool to randomly generate uncertain versions of TPC-H tables. The queries computed over the database instance are Q_1 , Q_2 , and Q_3 from [4], all of which are modified versions of TPC-H queries Q_3 , Q_6 , and Q_7 where all aggregations have been dropped.

As written, the queries disallow BIDD cross terms. We first ran all queries, noting the result size for each. Next the queries were rewritten so as not to filter out the cross terms. The comparison of the sizes of both result sets should then suggest in one way or another whether or not there exist many cross terms in practice. As seen, the experimental query results contain little to no cancelling terms. Fig. 6 shows the result sizes of the queries, where column CF is the result size when all cross terms are filtered out, column CI shows the number of output tuples when the cancelled tuples are included in the result, and the last column is the value of γ . The experiments show γ to be in a range between $[0, 0.1]\%$, indicating that only a negligible or constant (compare the result sizes of $Q_1 < Q_2$ and their respective γ values) amount of tuples are cancelled in practice when running queries over a typical BIDD instance. Interestingly, only one of the three queries had tuples that violated the BIDD constraint.

To conclude, the results in Fig. 6 show experimentally that γ is negligible in practice for BIDD queries. We also observe that (i) tuple presence is independent across blocks, so the corresponding probabilities (and hence p_0) are independent of the number of blocks, and (ii) BIDDs model uncertain attributes, so block size (and hence γ) is a function of the ‘‘messiness’’ of a dataset, rather than its size. Thus, we expect Theorem 4.7 to hold in general.

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

Figure 6: Number of Cancellations for Queries Over BIDD.

D CIRCUITS

D.1 Representing Polynomials with Circuits

D.1.1 Circuits for query plans. We now formalize circuits and the construction of circuits for \mathcal{RA}^+ queries. As mentioned earlier, we represent lineage polynomials as arithmetic circuits over \mathbb{N} -valued variables with $+$, \times . A circuit for query Q and $\mathbb{N}[X]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[X]}$ ¹⁸ is a directed acyclic graph $\langle V_{Q, \mathcal{D}_{\mathbb{N}[X]}} \cup E_{Q, \mathcal{D}_{\mathbb{N}[X]}}, \phi_{Q, \mathcal{D}_{\mathbb{N}[X]}}, \ell_{Q, \mathcal{D}_{\mathbb{N}[X]}} \rangle$ with vertices $V_{Q, \mathcal{D}_{\mathbb{N}[X]}}$ and directed edges $E_{Q, \mathcal{D}_{\mathbb{N}[X]}} \subset V_{Q, \mathcal{D}_{\mathbb{N}[X]}}^2$. The sink function $\phi_{Q, \mathcal{D}_{\mathbb{N}[X]}} : \mathcal{U}^n \rightarrow V_{Q, \mathcal{D}_{\mathbb{N}[X]}}$ is a partial function that maps the tuples of the n -ary relation $Q(\mathcal{D}_{\mathbb{N}[X]})$ to vertices. We require that $\phi_{Q, \mathcal{D}_{\mathbb{N}[X]}}$'s range be limited to sink vertices (i.e., vertices with out-degree 0). A function $\ell_{Q, \mathcal{D}_{\mathbb{N}[X]}} : V_{Q, \mathcal{D}_{\mathbb{N}[X]}} \rightarrow \{+, \times\} \cup \mathbb{N} \cup X$ assigns a label to each node: Source nodes (i.e., vertices with in-degree 0) are labeled with constants or variables (i.e., $\mathbb{N} \cup X$), while the remaining nodes are labeled with the symbol $+$ or \times . We require that vertices have an in-degree of at most two. Note that we can construct circuits for

¹⁸For background on $\mathbb{N}[X]$ -encoded PDBs, see Appendix A.1

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

D.2 Modeling Circuit Construction

We now connect the size of a circuit (where the size of a circuit is the number of vertices in the corresponding DAG) for a given \mathcal{RA}^+ query Q and $\mathbb{N}[X]$ -encoded PDB $\mathcal{D}_{\mathbb{N}[X]}$ to the runtime $T_{det}(Q, D, c)$ of the PDB's deterministic bounding database D .

Aaron says: @atri: do we use D or \bar{D} here?

We do this formally by showing that the size of the circuit is asymptotically no worse than the corresponding runtime of a large class of deterministic query processing algorithms.

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

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

encodes a polynomial, realized as

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

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

Algorithm 4 defines how the circuit for a query result is constructed. Denote the set of active output tuples as $\text{DOM}(\phi)$. We quickly review the number of vertices emitted in each case.

Base Relation. This circuit has $|D.R|$ vertices.

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

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

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

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

D.2.1 Bounding circuit depth. We first show that the depth of the circuit (DEPTH ; Definition 4.3) is bounded by the size of the query. Denote by $|Q|$ the number of relational operators in query Q , which recall we assume is a constant.

PROPOSITION D.1 (CIRCUIT DEPTH IS BOUNDED). *Let Q be a relational query and D be a deterministic bounding database with n tuples. There exists a (lineage) circuit C^* encoding the lineage of all tuples $t \in Q(D)$ for which $\text{DEPTH}(C^*) \leq O(k|Q| \log(n))$.*

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

D.2.2 Circuit size vs. runtime.

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

PROOF. We prove by induction that $|V_{Q, \mathcal{D}_{\mathbb{N}[X]}} \setminus \{v_0\}| \leq kT_{det}(Q, D, c)$. For clarity, we implicitly exclude v_0 in the proof below.

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

Algorithm 4 LINEAGECIRCUIT (Q, D, E, V, ℓ)**Input:** Q : query**Input:** D : a deterministic bounding database**Input:** E, V, ℓ : accumulators for the edge list, vertex list, and vertex label list.**Output:** $C = \langle V, E, \phi, \ell \rangle$: a circuit encoding the lineage of each tuple in $Q(D)$

```

1: if  $Q$  is  $R$  then ▷ Case 1:  $Q$  is a relation atom
2:   for  $t \in D.R$  do
3:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, R(t))\}$  ▷ Allocate a fresh node  $v_t$ ; note that when  $R(t) = c \cdot X_t$  for  $c > 1$ , we assume the algorithm
      generates a 3 node circuit encoding the multiplication of  $c \cdot X_t$ , adding the new vertices, edges, and vertice/label pairs to their respective
      sets.
4:      $\phi(t) \leftarrow v_t$ 
5:   end for
6:   return  $\langle V, E, \phi, \ell \rangle$ 
7: else if  $Q$  is  $\sigma_\theta(Q_1)$  then ▷ Case 2:  $Q$  is a Selection
8:    $\langle V, E, \phi', \ell \rangle \leftarrow \text{LINEAGECIRCUIT}(Q_1, D, V, E, \ell)$ 
9:   for  $t \in \text{DOM}(\phi')$  do
10:    if  $\theta(t)$  then  $\phi(t) \leftarrow \phi'(t)$  else  $\phi(t) \leftarrow v_0$ 
11:    end for
12:   return  $\langle V, E, \phi, \ell \rangle$ 
13: else if  $Q$  is  $\pi_A(Q_1)$  then ▷ Case 3:  $Q$  is a Projection
14:    $\langle V, E, \phi', \ell \rangle \leftarrow \text{LINEAGECIRCUIT}(Q_1, D, V, E, \ell)$ 
15:   for  $t \in \pi_A(Q_1(D))$  do
16:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, +)\}$  ▷ Allocate a fresh node  $v_t$ 
17:      $\phi(t) \leftarrow v_t$ 
18:   end for
19:   for  $t \in Q_1(D)$  do
20:      $E \leftarrow E \cup \{(\phi'(t), \phi(\pi_A t))\}$ 
21:   end for
22:   Correct nodes with in-degrees  $> 2$  by appending an equivalent fan-in two tree instead
23:   return  $\langle V, E, \phi, \ell \rangle$ 
24: else if  $Q$  is  $Q_1 \cup Q_2$  then ▷ Case 4:  $Q$  is a Bag Union
25:    $\langle V', E', \phi_1, \ell' \rangle \leftarrow \text{LINEAGECIRCUIT}(Q_1, D, V, E, \ell)$ 
26:    $\langle V, E, \phi_2, \ell \rangle \leftarrow \text{LINEAGECIRCUIT}(Q_2, D, V', E', \ell')$ 
27:    $\phi \leftarrow \phi_1 \cup \phi_2$ 
28:   for  $t \in \text{DOM}(\phi_1) \cap \text{DOM}(\phi_2)$  do
29:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, +)\}$  ▷ Allocate a fresh node  $v_t$ 
30:      $\phi(t) \leftarrow v_t$ 
31:      $E \leftarrow E \cup \{(\phi_1(t), v_t), (\phi_2(t), v_t)\}$ 
32:   end for
33:   return  $\langle V, E, \phi, \ell \rangle$ 
34: else if  $Q$  is  $Q_1 \bowtie \dots \bowtie Q_m$  then ▷ Case 5:  $Q$  is a  $m$ -ary Join
35:   for  $i \in [m]$  do
36:      $\langle V, E, \phi_i, \ell \rangle \leftarrow \text{LINEAGECIRCUIT}(Q_i, D, V, E, \ell)$ 
37:   end for
38:   for  $t \in \text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_m)$  do
39:      $V \leftarrow V \cup \{v_t\}; \ell \leftarrow \ell \cup \{(v_t, \times)\}$  ▷ Allocate a fresh node  $v_t$ 
40:      $\phi(t) \leftarrow v_t$ 
41:      $E \leftarrow E \cup \{(\phi_i(\pi_{\text{sch}}(Q_i(D))(t)), v_t) \mid i \in [m]\}$ 
42:   end for
43:   Correct nodes with in-degrees  $> 2$  by appending an equivalent fan-in two tree instead
44:   return  $\langle V, E, \phi, \ell \rangle$ 
45: end if

```

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

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

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

(From the inductive assumption)

$$\leq kT_{det}(Q_1, D, c) + |Q_1|$$

(By definition of $T_{det}(Q, D, c)$)

$$\leq kT_{det}(Q, D, c).$$

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

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

(From the inductive assumption)

$$\leq k(T_{det}(Q_1, D, c) + T_{det}(Q_2, D, c)) + (|Q_1| + |Q_2|)$$

(By definition of $T_{det}(Q, D, c)$)

$$\leq k(T_{det}(Q, D, c)).$$

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

$$|V_{Q, \mathcal{D}_{\mathbb{N}[X]}}| = |V_{Q_1, \mathcal{D}_{\mathbb{N}[X]}}| + \dots + |V_{Q_m, \mathcal{D}_{\mathbb{N}[X]}}| + (m-1)|Q_1 \bowtie \dots \bowtie Q_m|$$

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

$$\begin{aligned} &\leq kT_{det}(Q_1, D, c) + \dots + kT_{det}(Q_m, D, c) + \\ &\quad (m-1)|Q_1 \bowtie \dots \bowtie Q_m| \\ &\leq k(T_{det}(Q_1, D, c) + \dots + T_{det}(Q_m, D, c) + \\ &\quad |Q_1 \bowtie \dots \bowtie Q_m|) \end{aligned}$$

(By definition of $T_{det}(Q, D, c)$ and assumption on $T_{join}(\cdot)$)

$$\leq kT_{det}(Q, D, c).$$

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

D.2.3 Runtime of LINEAGECIRCUIT. We next need to show that we can construct the circuit in time linear in the deterministic runtime.

LEMMA D.3. *Given a query Q over a deterministic bounding database D and the C^* output by Algorithm 4, the runtime $T_{LINEAGECIRCUIT}(Q, D, C^*) \leq O(T_{det}(Q, D, c))$.*

PROOF. By analysis of Algorithm 4, invoked as $C^* \leftarrow \text{LINEAGECIRCUIT}(Q, D, \emptyset, \{v_0\}, \{(v_0, 0)\})$.

We assume that the vertex list V , edge list E , and vertex label list ℓ are mutable accumulators with $O(1)$ amortized append. We assume that the tuple to sink mapping ϕ is a linked hashmap, with $O(1)$ insertions and retrievals, and $O(n)$ iteration over the domain of keys. We assume that the n -ary join $\text{DOM}(\phi_1) \bowtie \dots \bowtie \text{DOM}(\phi_n)$ can be computed in time $T_{join}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_n))$ (Definition 2.13) and that an intersection $\text{DOM}(\phi_1) \cap \text{DOM}(\phi_2)$ can be computed in time $O(|\text{DOM}(\phi_1)| + |\text{DOM}(\phi_2)|)$ (e.g., with a hash table).

Before proving our runtime bound, we first observe that $T_{det}(Q, D, c) \geq \Omega(|Q(D)|)$. This is true by construction for the relation, projection, and union cases, by Definition 2.13 for joins, and by the observation that $|\sigma(R)| \leq |R|$.

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

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

Projection. Projection requires a recursive call to Algorithm 4, which by the recursive assumption is bounded by $O(T_{det}(Q', D, c))$, which in turn is a term in $T_{det}(\pi_A Q', D, c)$. What remains is an iteration over $\pi_A(Q(D))$ (lines 13–16), an iteration over $Q'(D)$ (lines 17–19), and the construction of a fan-in tree (line 20). The first iteration is $O(|Q(D)|) \leq O(T_{det}(Q, D, c))$. The second iteration and the construction of the bounded fan-in tree are both $O(|Q'(D)|) \leq O(T_{det}(Q', D, c)) \leq O(T_{det}(Q, D, c))$, by the the observation above that $T_{det}(Q, D, c) \geq \Omega(|Q(D)|)$.

Bag Union. As above, the recursive calls explicitly correspond to terms in the expansion of $T_{det}(Q_1 \cup Q_2, D, c)$. Initializing ϕ (Line 27) can be accomplished in $O(\text{DOM}(\phi_1) + \text{DOM}(\phi_2)) = O(|Q_1(D)| + |Q_2(D)|) \leq O(T_{det}(Q_1, D) + T_{det}(Q_2, D, c))$. The remainder requires computing $Q_1 \cap Q_2$ (Line 28) and iterating over it (Lines 29 to 31), which is $O(|Q_1| + |Q_2|)$ as noted above – this directly corresponds to terms in $T_{det}(Q_1 \cup Q_2, D, c)$.

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

(i) $T_{join}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_m))$

(ii) $O(|Q_1(D) \bowtie \dots \bowtie Q_m(D)|) \leq O(T_{join}(\text{DOM}(\phi_1), \dots, \text{DOM}(\phi_m)))$ (Definition 2.13)

(iii) $O(m|Q_1(D) \bowtie \dots \bowtie Q_m(D)|)$ (as (ii), noting that $m \leq k = O(1)$) □

E HIGHER MOMENTS

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

F THE KARP-LUBY ESTIMATOR

Computing the marginal probability of a tuple in the output of a set-probabilistic database query has been studied extensively. To the best of our knowledge, the current state of the art approximation algorithm for this problem is the Karp-Luby estimator [32], which first appeared in MayBMS/Sprout [39], and more recently as part of an online “anytime” approximation algorithm [16, 21].

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

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

Aaron says: Why are we computing the cardinality of variables that equal 1?

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

Aaron says: A bit confused on the above sentence:

- i) what is meant by conditioning \mathcal{P} on a W_i ,
- ii) why is each W_i a monomial?

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

The Karp-Luby estimator is employed on the SMB representation¹⁹ of C (to solve the set-PDB version of Problem 1.6), where each W_i represents the event that one monomial is true. By simple inspection, if there are ℓ monomials, this estimator has runtime $\Omega(\ell)$. Further, a minimum of $\left\lceil \frac{3 \cdot \ell \cdot \log(\frac{2}{\delta})}{\epsilon^2} \right\rceil$ invocations of the estimator are required to achieve $1 \pm \epsilon$ approximation with probability at least $1 - \delta$ [39], entailing a runtime at least quadratic in ℓ . As an arbitrary lineage circuit C may encode $\Omega(|C|^k)$ monomials, the worst case runtime is at least $\Omega(|C|^{2k})$ (where k is the ‘degree’ of lineage polynomial encoded by C). By contrast note that by the discussion after Lemma 4.9 we can solve Problem 1.6 in time $O(|C|^2)$ for all BIDD circuits *independent* of the degree k .

G PARAMETERIZED COMPLEXITY

In Sec. 3, we utilized common conjectures from fine-grained complexity theory. The notion of $\#W[1]$ – *hard* is a standard notion in *parameterized complexity*, which by now is a standard complexity tool in providing data complexity bounds on query processing results [24]. E.g. the fact that k -matching is $\#W[1]$ – *hard* implies that we cannot have an $n^{\Omega(1)}$ runtime. However, these results do not carefully track the exponent in the hardness result. E.g. $\#W[1]$ – *hard* for the general k -matching problem does not imply anything specific for the 3-matching problem. Similar questions have led to intense research into the new sub-field of *fine-grained complexity* (see [50]), where we care about

¹⁹Note that since we are in the set semantics, in the lineage polynomial/formula, addition is logical OR and multiplication is logical AND.

the exponent in our hardness assumptions as well– e.g. Conjecture 3.3 is based on the popular *Triangle detection hypothesis* in this area (cf. [35]).