

You should NEVER use a term before defining it.

23:2 Bag PDB Queries

44 considering bag query semantics. We denote by $Q(\mathbf{W})(t)$ the multiplicity of t in query Q
 45 over possible world $\mathbf{W} \in \{0, \dots, c\}^D$.

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

48 ► **Problem 1.1.** Given a c -TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, \mathcal{RA}^+ query Q ¹, and result tuple
 49 t , compute the expected multiplicity of t : $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}} [Q(\mathbf{W})(t)]$.

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

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

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

67 Let $T_{det}(\text{OPT}(Q), \bar{D}, c)$ (see Sec. 2.4 for further details) denote the runtime for query
 68 $\text{OPT}(Q)$, deterministic database \bar{D} , and multiplicity bound c . Being we consider \mathcal{RA}^+
 69 queries in which order of operators can impact runtime, we denote the optimal query as
 70 $\text{OPT}(Q) = \min_{Q' \in \mathcal{RA}^+, Q' \equiv Q} T_{det}(Q', \bar{D}, c)$.

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

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

71 **Our lower bound results.** Our question is whether or not it is always true that $T^*(Q, \mathcal{D}) \leq$
 72 $T_{det}(\text{OPT}(Q), D, c)$. Unfortunately this is not the case. Table 1 shows our results.

¹ A query Q is an \mathcal{RA}^+ query if it is composed entirely of one or more of the positive relational operators $\{\sigma, \pi, \bowtie, \cup\}$.

210 ▶ **Problem 1.6.** Given one circuit \mathcal{C} that encodes $\Phi[Q, D, t]$ for all result tuples t (one sink
 211 per t) for bag-PDB \mathcal{D} and \mathcal{RA}^+ query Q , does there exist an algorithm that computes a
 212 $(1 \pm \epsilon)$ -approximation of $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}} [Q(\mathbf{W})(t)]$ (for all result tuples t) in $O(|\mathcal{C}|)$ time?

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

$$217 \quad \Phi_1^2(\mathbf{p}) = p_{AP}^2 p_{XP}^2 p_B^2 + p_{BP}^2 p_{YP}^2 p_E^2 + p_{BP}^2 p_{ZP}^2 p_C^2 + 2p_{AP} p_X p_{BP} p_{YP} p_E + 2p_{AP} p_X p_{BP} p_{ZP} p_C + 2p_{BP}^2 p_{YP} p_E p_{ZP} p_C$$

$$218 \quad \leq p_{AP} p_X p_B + p_{BP} p_Y p_E + p_{BP} p_Z p_C + 2p_{AP} p_X p_{BP} p_{YP} p_E + 2p_{AP} p_X p_{BP} p_{ZP} p_C + 2p_{BP} p_{YP} p_E p_{ZP} p_C = \tilde{\Phi}_1^2(\mathbf{p})$$

220 If we assume that all seven probability values are at least $p_0 > 0$, we get that $\Phi_1^2(\mathbf{p})$ is in
 221 the range $[(p_0)^3 \cdot \tilde{\Phi}_1^2(\mathbf{p}), \tilde{\Phi}_1^2(\mathbf{p})]$. In sec. 4 we demonstrate that a $(1 \pm \epsilon)$ (multiplicative)
 222 approximation with competitive performance is achievable. To get an $(1 \pm \epsilon)$ -multiplicative
 223 approximation and solve Problem 1.6, using \mathcal{C} we uniformly sample monomials from the
 224 equivalent SMB representation of Φ (without materializing the SMB representation) and
 225 ‘adjust’ their contribution to $\tilde{\Phi}(\cdot)$.

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

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

239 **2 Background and Notation**

240 **2.1 Polynomial Definition and Terminology**

241 A polynomial over a set of variables \mathbf{S} with $|S| = m$ and individual degree $B < \infty$ is formally
 242 defined as (where $c_{\mathbf{d}} \in \mathbb{N}$):

$$243 \quad \Phi(S_1, \dots, S_m) = \sum_{\mathbf{d} \in \{0, \dots, B\}^D} c_{\mathbf{d}} \cdot \prod_{i \in [m]} S_i^{d_i}. \tag{1}$$

244 **Definition 2.1** (Standard Monomial Basis). The term $\prod_{i \in [m]} S_i^{d_i}$ in Eq. (1) is a monomial.
 245 A polynomial $\Phi(\mathbf{X})$ is in standard monomial basis (SMB) when we keep only the terms with
 246 $c_{\mathbf{d}} \neq 0$ from Eq. (1).
 247

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

250 **Definition 2.2** (Degree). The degree of polynomial $\Phi(\mathbf{X})$ is the largest $d_{\mathbf{d}} = \sum_{i \in [m]} d_i$ such
 251 that $c_{(d_1, \dots, d_n)} \neq 0$. We denote the degree of Φ as $\text{deg}(\Phi)$.

If the idea is you want to preserve X for Φ use X or 2. Why are you using S instead of X?

Perhaps better to use some other notation than Φ

don't match Φ for a general poly & use Φ only for lineage polys.

23:8 Bag PDB Queries

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

257 2.2 1-BIDB

258 A block independent database BIDB \mathcal{D}' can be viewed as a 1-TIDB \mathcal{D} with the added flexibility
259 that each $t \in \mathcal{D}$ has multiple disjoint alternatives, i.e., all $t \in \mathcal{D}'$ are partitioned into m
260 independent blocks with the condition that tuples $t \in b_i$ for $i \in [m]$ are disjoint events. We
261 define next a specific construction of BIDB that is useful for our work.

262 **Definition 2.3** (1-BIDB). Define a 1-BIDB to be the pair $\mathcal{D}' = (\prod_{t \in \mathcal{D}'} \{0, c_t\}, \mathcal{P}')$, where
263 \mathcal{D}' is the set of possible tuples such that each $t \in \mathcal{D}'$ has a multiplicity domain of $\{0, c_t\}$,
264 with $c_t \in \mathbb{N}$. The operation $\prod_{t \in \mathcal{D}'}$ is the direct product of all such multiplicity domain pairs.
265 The tuples $t \in \mathcal{D}'$ are partitioned into m independent blocks b_i , $i \in [m]$, of disjoint tuples. \mathcal{P}'
266 is the probability distribution across all worlds such that, given $\mathbf{W} \in \prod_{t \in \mathcal{D}'} \{0, c_t\}$, $t, t' \in$
267 $b_i : Pr[\mathbf{W}_t, \mathbf{W}_{t'} > 0] = 0$.

268 We now present a reduction that is useful in deriving our results:

269 **Definition 2.4** (*c*-TIDB reduction). Given *c*-TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$, let $\mathcal{D}' =$
270 $(\prod_{t \in \mathcal{D}'} \{0, c_t\}^{\mathcal{D}'}, \mathcal{P}')$ be the 1-BIDB obtained in the following manner: for each $t \in \mathcal{D}$,
271 create block $b_t = \{ \langle t, j \rangle_{j \in [c]} \}$ of disjoint tuples, for all $j \in [c]$. The probability distribution
272 \mathcal{P}' is the one induced by $\mathbf{p} = ((p_{t,j})_{t \in \mathcal{D}, j \in [c]})$ and the BIDB disjoint requirement, where
273 given any $\mathbf{W} \in \prod_{t \in \mathcal{D}'} \{0, c_t\}^{\mathcal{D}'}$, $Pr[\mathbf{W}_{t,j}, \mathbf{W}_{t,j'} > 0] = 0$ for any $j \neq j' \in [c]$, such that
274 for any $W \in \prod_{t \in \mathcal{D}'} \{0, c_t\}^{\mathcal{D}'}$, $Pr[\mathbf{W} = W] = \prod_{t \in \mathcal{D}'} \prod_{j \in [c]} W_{t,j} \cdot j \cdot p_t$ if $\forall t \in \mathcal{D}' \forall j \neq j' \in [c]$
275 $[c], W_{t,j}, W_{t,j'} \geq 1$; otherwise $Pr[\mathbf{W} = W] = 0$.⁴

276 For the *c*-TIDB \mathcal{D} , each $X_t \in [c]$, while in the reduced 1-BIDB \mathcal{D}' , each $X_{t,j} \in$
277 $\{0, 1\}$. Hence, in the setting of 1-BIDB, we have the following semantics for generating
278 lineage polynomials in \mathcal{RA}^+ queries: $\Phi'[\pi_A(Q), \mathcal{D}', t_j] = \sum_{t_{j'} \in \pi_A(Q(\mathcal{D}')): t_{j'} = t_j} \Phi'[Q, \mathcal{D}', t_{j'}]$,
279 $\Phi'[\sigma_\theta(Q), \mathcal{D}', t_j] = \begin{cases} \theta = 1 & \Phi'[Q, \mathcal{D}', t_j] \\ \theta = 0 & 0 \end{cases}$, $\Phi'[Q_1 \bowtie Q_2, \mathcal{D}', t_j] = \Phi'[Q_1, \mathcal{D}', \pi_{attr(Q_1)}(t_j)] \cdot$
280 $\Phi'[Q_2, \mathcal{D}', \pi_{attr(Q_2)}(t_j)]$, $\Phi'[Q_1 \cup Q_2, \mathcal{D}', t_j] = \Phi'[Q_1, \mathcal{D}', t_j] + \Phi'[Q_2, \mathcal{D}', t_j]$, and the base
281 case now becomes $\Phi'[R, \mathcal{D}', t_j] = j \cdot X_{t,j}$ (c.f. Fig. 1). Then given the disjoint requirement and
282 the semantics for constructing the lineage polynomial over a 1-BIDB, $\Phi'[R, \mathcal{D}', t]$ is of the same
283 structure as the reformulated polynomial Φ_R of step i) from Definition 1.3, which then implies
284 that $\tilde{\Phi}'$ is the reduced polynomial that results from step ii) of Definition 1.3, and further
285 that Lemma 1.4 immediately follows for 1-BIDB polynomials: $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}'}[\Phi'(\mathbf{W})] = \tilde{\Phi}'(\mathbf{p})$.

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

287 **Corollary 2.5.** If Φ is a 1-BIDB lineage polynomial already in SMB, then the expectation
288 of Φ , i.e., $\mathbb{E}[\Phi] = \tilde{\Phi}(p_1, \dots, p_n)$ can be computed in $O(|\Phi|)$ time.

⁴ We 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} .

Handwritten notes:
- Can just say that \mathcal{D}' is disjoint union of sets of tuples B_1, \dots, B_n . Also upper case for blocks is better.
- I think this is for confusing. If you want you can say 1-BIDB is a gen of TIDB.
- if clear from the context \rightarrow PROPAGATE
- not defined $\rightarrow \mathcal{D}'$ is partitioned
- but you sentence seems to suggest t is partitioned
- As X 's standard for Cartesian product. part is not specific
- This part should be in Def 2.3. Here you are defining what the ops do for \mathcal{D}' . should be sorry should have said this before
- Actually, this should be stated as lemma.
- Update like this
- W ?
- Put this in its own fig.
- it should have its own defn
- here are for Φ not the PDB
- $\Phi \triangleq \tilde{\Phi}$ for TIDB. Finally state this as a lemma as well.
- First give fig to define $\tilde{\Phi}'$. Then give definition of $\tilde{\Phi}'$. Then compare the two to $\Phi \triangleq \tilde{\Phi}$ for TIDB.

Need a filler sentence
say readers can skip
this section w/o losing anything
S. Feng, B. Glavic, A. Huber, O. Kennedy, A. Rudra

Connect this
defs to the
defs we have
used for C-7IDB
& 1-BIDB
spec what is P
here

289 2.2.1 Possible World Semantics

290 Queries over probabilistic databases are traditionally viewed as being evaluated using the
291 so-called possible world semantics. A general bag-PDB can be defined as the pair $\mathcal{D} = (\Omega, \mathcal{P})$
292 where Ω is the set of possible worlds represented by \mathcal{D} . Under the possible world semantics,
293 the result of a query Q over an incomplete database Ω is the set of query answers produced
294 by evaluating Q over each possible world $\omega \in \Omega$: $\{Q(\omega) : \omega \in \Omega\}$. The result of a query is
295 the pair $(Q(\omega), \mathcal{P}')$ where \mathcal{P}' is a probability distribution that assigns to each possible query
296 result the sum of the probabilities of the worlds that produce this answer: $Pr[\omega \in \Omega] =$
297 $\sum_{\omega' \in \Omega, Q(\omega')=Q(\omega)} Pr[\omega']$.

Aaron says: I am not sure the following paragraph is needed, since the reduction definition says pretty much the same thing. Unless that definition changes, we can get rid of this paragraph.

298
299 Suppose that \mathcal{D} is a 1-BIDB. Instead of looking only at the possible worlds of \mathcal{D} , one
300 can consider all worlds, including those that cannot exist due to, e.g., disjointness. The
301 all worlds set can be modeled by $\mathbf{W} \in \{0, 1\}^{cn}$, such that $\mathbf{W}_{t,j} \in \mathbf{W}$ represents whether
302 or not the multiplicity of t is j (here and later, especially in Sec. 4, we will rename the
303 variables as X_1, \dots, X_n , where $n = \sum_{t \in D} |b_t|$). We can denote a probability distribution
304 over all $\mathbf{W} \in \{0, 1\}^{cn}$ as \mathcal{P}' . When \mathcal{P}' is the one induced from each $p_{t,j}$ while assigning
305 $Pr[\mathbf{W}] = 0$ for any \mathbf{W} with $\mathbf{W}_{t,j}, \mathbf{W}_{t,j'} \geq 1$ for $j \neq j'$, we end up with a bijective mapping
306 from \mathcal{P} to \mathcal{P}' , such that each mapping is equivalent, implying the distributions are equivalent.
307 Appendix B.2 has more details.

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

311 ► **Proposition 2.6** (Expectation of polynomials). Given a bag-PDB $\mathcal{D} = (\Omega, \mathcal{P})$, $\mathcal{R}A^+$ query
312 Q , and lineage polynomial $\Phi[Q, D, t]$ for arbitrary result tuple t , we have (denoting \mathbf{D} as the
313 random variable over Ω): $\mathbb{E}_{\mathbf{D} \sim \mathcal{P}}[Q(\mathbf{D})(t)] = \mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$.

314 A formal proof of Proposition 2.6 is given in Appendix B.3.⁵ We focus on the problem of
315 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
316 $\Phi[Q, D, t]$ (i.e., $\Phi(\mathbf{X})$ will denote a polynomial).

317 2.3 Formalizing Problem 1.6

318 Problem 1.6 asks if there exists a linear time approximation algorithm in the size of a given
319 circuit C which encodes $\Phi(\mathbf{X})$. In this work we represent lineage polynomials via arithmetic
320 circuits [9], a standard way to represent polynomials over fields (particularly in the field of
321 algebraic complexity) that we use for polynomials over \mathbb{N} in the obvious way. Since we are
322 particularly using circuits to model lineage polynomials, we can refer to these circuits as
323 lineage circuits. However, when the meaning is clear, we will drop the term lineage and only
324 refer to them as circuits.

325 ► **Definition 2.7** (Circuit). A circuit C is a Directed Acyclic Graph (DAG) whose source
326 gates (in degree of 0) consist of elements in either \mathbb{N} or $\mathbf{X} = (X_1, \dots, X_n)$. For each result

⁵ Although Proposition 2.6 follows, e.g., as an obvious consequence of [28]’s Theorem 7.1, we are unaware of any formal proof for bag-probabilistic databases.

Technically
only for
C-7IDB
& not
bag
PDBs in
genl.

This part
can do this
instead

Specific

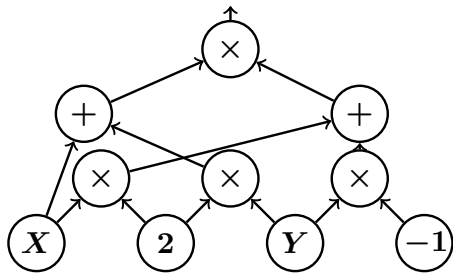
327 tuple there exists one sink gate. The internal gates have binary input and are either sum (+)
 328 or product (\times) gates. Each gate has the following members: *type*, *partial*, *input*, *degree*,
 329 *Lweight*, and *Rweight*, where *type* is the value type $\{+, \times, \text{VAR}, \text{NUM}\}$ and *input* the list of
 330 inputs. Source gates have an extra member *val* storing the value. C_L (C_R) denotes the left
 331 (right) input of C .

Aaron says: Does the following matter, i.e., does it point anything out special for our research? **EDIT:** Lemma 4.8 does use this (when C is a tree) to answer Problem 1.6 with a yes.

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

336 The circuits in Fig. 2 encode their respective polynomials in column Φ . Note that each
 337 circuit C encodes a tree, with edges pointing towards the root.

the last one is not a tree



346 **Figure 3** Circuit encoding of $(X + 2Y)(2X - Y)$
 347

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.

maps

Definition 2.8 ($\text{POLY}(\cdot)$). ~~Denote $\text{POLY}(C)$ to be the function from the sink of circuit C to its corresponding polynomial (in SMB).~~ $\text{POLY}(\cdot)$ is recursively defined on C as follows, with addition and multiplication following the standard interpretation for polynomials:

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

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

352 **Definition 2.9** (Circuit Set). $\text{CSet}(\Phi(\mathbf{X}))$ is the set of all possible circuits C such that
 353 $\text{POLY}(C) = \Phi(\mathbf{X})$.

354 The circuit of Fig. 3 is an element of $\text{CSet}(2X^2 + 3XY - 2Y^2)$. One can think of
 355 $\text{CSet}(\Phi(\mathbf{X}))$ as the infinite set of circuits where for each element C , $\text{POLY}(C) = \Phi(\mathbf{X})$.

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

357 **Definition 2.10** (The Expected Result Multiplicity Problem). Let \mathcal{D} be an arbitrary BIDB-
 358 PDB and \mathbf{X} be the set of variables annotating tuples in D_Ω . Fix an \mathcal{RA}^+ query Q and a
 359 result tuple t . The EXPECTED RESULT MULTIPLICITY PROBLEM is defined as follows:

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

362 2.4 Relationship to Deterministic Query Runtimes

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

365 ▶ **Definition 2.11** (Join Cost). Denote by $T_{join}(R_1, \dots, R_m)$ the runtime of an algorithm
 366 for computing the m -ary join $R_1 \bowtie \dots \bowtie R_m$. We require only that the algorithm must
 367 enumerate its output, i.e., that $T_{join}(R_1, \dots, R_m) \geq |R_1 \bowtie \dots \bowtie R_m|$.

needs to be updated. OPT(Q) needs to be brought in here.

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

$$T_{det}(R, \bar{D}, c) = |\bar{D}.R| \quad T_{det}(\sigma Q, \bar{D}, c) = T_{det}(Q, \bar{D}) \quad T_{det}(\pi Q, \bar{D}, c) = T_{det}(Q, \bar{D}, c) + |Q(\bar{D})|$$

$$T_{det}(Q \cup Q', \bar{D}, c) = T_{det}(Q, \bar{D}, c) + T_{det}(Q', \bar{D}, c) + |Q(\bar{D})| + |Q'(\bar{D})|$$

$$T_{det}(Q_1 \bowtie \dots \bowtie Q_m, \bar{D}, c) = T_{det}(Q_1, \bar{D}, c) + \dots + T_{det}(Q_m, \bar{D}, c) + T_{join}(Q_1(\bar{D}), \dots, Q_m(\bar{D}))$$

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

375 Finally, Lemma E.2 and Lemma E.3 show that for any \mathcal{RA}^+ query Q and D , there exists
 376 a circuit \mathbf{C}^* such that $T_{LC}(Q, D, \mathbf{C}^*)$ and $|\mathbf{C}^*|$ are both $O(T_{det}(Q, D, c))$. Recall we assumed
 377 these two bounds when we moved from Problem 1.5 to Problem 1.6.

OPT(Q)

3 Hardness of Exact Computation

381 In this section, we will prove the hardness results claimed in Table 1 for a specific (family) of
 382 hard instance (Q, \mathcal{D}) for Problem 1.2 where \mathcal{D} is a 1-TIDB. Note that this implies hardness
 383 for c -TIDBs ($c \geq 1$), BIDBs and general bag-PDB, showing Problem 1.2 cannot be done in
 384 $O(T_{det}(\text{OPT}(Q), D, c))$ runtime.

Problem 1.2 only talks about c-TIDBs. You should definitely mention the results also work for 1-BIDB & other more general bag PDBs.

3.1 Preliminaries

385 Our hardness results are based on (exactly) counting the number of (not necessarily induced)
 386 subgraphs in G isomorphic to H . Let $\#(G, H)$ denote this quantity. We can think of H
 387 as being of constant size and G as growing. In particular, we will consider the problems of
 388 computing the following counts (given G in its adjacency list representation): $\#(G, \triangle)$ (the
 389 number of triangles), $\#(G, \text{3-matchings})$ (the number of 3-matchings), and the latter's generalization
 390 $\#(G, \text{3-matchings})$ (the number of 3-matchings), and the latter's generalization
 391 $\#(G, \text{3-matchings})$ (the number of k -matchings). We use $T_{match}(k, G)$ to denote the optimal
 392 runtime of computing $\#(G, \text{3-matchings})$ exactly. Our hardness results in Sec. 3.2 are based on
 393 the following hardness results/conjectures:

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

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

399 We note that the above conjecture is somewhat non-standard. In particular, the best known
 400 algorithm to compute $\#(G, \text{3-matchings})$ takes time $\Omega(|V|^{k/2})$ (i.e. if this is the best algorithm
 401 then $c_0 = \frac{1}{4}$) [11]. What the above conjecture is saying is that one can only hope for a
 402 polynomial improvement over the state of the art algorithm to compute $\#(G, \text{3-matchings})$.

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

Also why do you have different p_i for each i ? We only use same p for all i .

404 ► **Conjecture 3.3.** *There exists a constant $\epsilon_0 > 0$ such that given an undirected graph*
 405 $G = (V, E)$, *computing $\#(G, \mathfrak{A})$ exactly cannot be done in time $o(|E|^{1+\epsilon_0})$.*

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

408 All of our hardness results rely on a simple lineage polynomial encoding of the edges
 409 of a graph. To prove our hardness result, consider a graph $G = (V, E)$, where $|E| = m$,
 410 $V = [n]$. Our lineage polynomial has a variable X_i for every i in $[n]$. Consider the polynomial
 411 $\Phi_G(\mathbf{X}) = \sum_{(i,j) \in E} X_i \cdot X_j$. The hard polynomial for our problem will be a suitable power $k \geq 3$
 412 of the polynomial above:

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

414
$$\Phi_G^k(X_1, \dots, X_n) = \left(\sum_{(i,j) \in E} X_i \cdot X_j \right)^k$$

spell it out. could be a simple as sum apply Fig. 1 to the Q2 we get the data

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

```
417 SELECT 1 FROM T t1, R r, T t2
418 WHERE t1.city = r.city1 AND t2.city = r.city2
```

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

```
422 SELECT COUNT(*) FROM R1 JOIN R2 JOIN ... JOIN Rk
```

make sure notation is consistent w/ def / notation from Sec 1.

425 Consider again the c -TIDB instance \mathcal{D} of Fig. 2 and, for our hard instance, let $c = 1$.
 426 \mathcal{D} generalizes to one compatible to Definition 3.4 as follows. Relation T has n tuples
 427 corresponding to each vertex for i in $[n]$, each with probability p_i and R has tuples
 428 corresponding to the edges E (each with probability of 1).⁶ In other words, for this instance
 429 \mathcal{D} contains the set of n unary tuples in T (which corresponds to V) and m binary tuples
 430 in R (which corresponds to E). Note that this implies that Φ_G^k is indeed a c -TIDB-lineage
 431 polynomial.

Aaron says: @atri, we discussed this last meeting, but I am not sure if we really pinpointed how we want to treat (in a consistent manner) the runtime of Lemma 3.5 since k is a constant and m is growing. Would it be a good idea to be consistent with the O_ϵ notation of Problem 1.5 and say $O_k(m)$

Score!

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

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

436 3.2 Multiple Distinct p Values

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

⁶ Technically, $\Phi_G^k(\mathbf{X})$ should have variables corresponding to tuples in R as well, but since they always are present with probability 1, we drop those. Our argument also works when all the tuples in R also are present with probability p but to simplify notation we assign probability 1 to edges.

438 ▶ **Theorem 3.6.** Let p_0, \dots, p_{2k} be $2k + 1$ distinct values in $(0, 1]$. Then computing
 439 $\tilde{\Phi}_G^k(p_i, \dots, p_i)$ (over all $i \in [2k + 1]$ for arbitrary $G = (V, E)$) needs time $\Omega(T_{match}(k, G))$,
 440 assuming $T_{match}(k, G) \geq \omega(|E|)$.

441 Note that the second row of Table 1 follows from Proposition 2.6, Theorem 3.6, Lemma 3.5,
 442 and Theorem 3.1 while the third row is proved by Proposition 2.6, Theorem 3.6, Lemma 3.5,
 443 and Conjecture 3.2. Since Conjecture 3.2 is non-standard, the latter hardness result should
 444 be interpreted as follows. Any substantial polynomial improvement for Problem 1.2 (over the
 445 trivial algorithm that converts Φ into SMB and then uses Corollary 2.5 for EC) would lead
 446 to an improvement over the state of the art *upper* bounds on $T_{match}(k, G)$. Finally, note
 447 that Theorem 3.6 needs one to be able to compute the expected multiplicities over $(2k + 1)$
 448 distinct values of p_i , each of which corresponds to distinct \mathcal{P} (for the same D), which explain
 449 the ‘Multiple’ entry in the second column in the second and third row in Table 1. Next, we
 450 argue how to get rid of this latter requirement.

451 3.3 Single p value

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

456 ▶ **Theorem 3.7.** Fix $p \in (0, 1)$. Then assuming Conjecture 3.3 is true, any algorithm that
 457 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
 458 ϵ_0 is as defined in Conjecture 3.3.

459 Note that Proposition 2.6 and Theorem 3.7 above imply the hardness result in the first
 460 row of Table 1. We note that Theorem 3.1 and Conjecture 3.2 (and the lower bounds in the
 461 second and third row of Table 1) need k to be large enough (in particular, we need a family
 462 of hard queries). But the above Theorem 3.7 (and the lower bound in first row of Table 1)
 463 holds for $k = 3$ (and hence for a fixed query).

464 4 $1 \pm \epsilon$ Approximation Algorithm

465 In Sec. 3, we showed that Problem 1.2 cannot be solved in $O(T_{det}(\text{OPT}(Q), D, c))$ runtime.
 466 With this result, we now design an approximation algorithm for our problem that runs in
 467 $O(|C|)$ for a very broad class of circuits, (thus affirming Problem 1.6) see the discussion after
 468 Lemma 4.8 for more). The following approximation algorithm applies to c -TIDB lineage
 469 polynomials and general BIDB (over bag- \mathcal{RA}^+ query semantics) lineage polynomials in
 470 practice, where for the latter we note that a 1-TIDB is equivalently a 1-BIDB (blocks are
 471 size 1) and our experimental results (see Appendix D.10) using queries from the PDBench
 472 benchmark [1] show a low γ (see Definition 4.6) supporting the notion that our bounds hold
 473 for general BIDB in practice. Corresponding proofs and pseudocode for all formal statements
 474 and algorithms can be found in Appendix D.

475 4.1 Preliminaries and some more notation

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

→ First we want approximation ↑ time. THEN can we do that by showing the result via circuits
 This sentence is too long
 ↑ give flow gr
 ↓ Did not read rest of sec 4 super carefully

you should lead with the main result here. \uparrow i.e. we solve P 1.6 from any fixed $\epsilon > 0$

4.2 Our main result

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

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

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

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

► **Theorem 4.7.** Let \mathcal{C} be an arbitrary 1-BIDB circuit, define $\Phi(\mathbf{X}) = \text{POLY}(\mathcal{C})$, let $k = \text{DEG}(\mathcal{C})$, and let $\gamma = \gamma(\mathcal{C})$. 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

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

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

The restriction on γ is satisfied by any 1-TIDB (where $\gamma = 0$ in the equivalent 1-BIDB of Definition 2.4) as well as for all three queries of the PDBench BIDB benchmark (see Appendix D.10 for experimental results).

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

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

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

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 $O_k \left(\frac{1}{(\epsilon')^2} \cdot \text{SIZE}(\mathcal{C})^2 \cdot \log \frac{1}{\delta} \right)$ for general circuits \mathcal{C} . If \mathcal{C} is a tree, then the runtime simplifies to $O_k \left(\frac{1}{(\epsilon')^2} \cdot \text{SIZE}(\mathcal{C}) \cdot \log \frac{1}{\delta} \right)$, which then answers Problem 1.6 with yes for such circuits.

where is the corollary for c -TIDB? That is. the main results of this section

Aaron says: Is it standard to assume that in the asymptotic notation above ϵ and δ are constant? Otherwise this does not uphold Problem 1.6.

542

543

544

545

546

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

547

548

549

550

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

Aaron says: What is $|Q|$? Isn't that just k ?

551

552

553

554

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

555

556

557

558

559

560

561

562

563

564

565

566

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 [46] (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. [14] and Olteanu et al. [21] proved dichotomies for UCQs and two classes of queries with negation, respectively. Amarilli et al. investigated tractable classes of databases for more complex queries [3]. Another line of work studies which structural properties of lineage formulas lead to tractable cases [31, 41, 44]. In this paper we focus on intensional query evaluation with polynomials.

567

568

569

570

571

572

573

574

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

575

576

577

578

579

580

581

582

583

584

585

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

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

Query size.
No, it is not exactly k

why should ϵ & δ be constants?
Results hold for any $\epsilon, \delta > 0$