

General suggestion: Do NOT define B1DB/T1DB in possible world

S. Feng, B. Glavic, A. Huber, O. Kennedy, A. Rudra

23:7

semantics. I think trying to fit stuff into possible world just make notation more complicated

216 equivalent SMB representation of Φ (without materializing the SMB representation) and
 217 'adjust' their contribution to $\tilde{\Phi}(\cdot)$.

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

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

231 **2 Background and Notation**

232 **2.1 Polynomial Definition and Terminology**

233 A polynomial over $\mathbf{X} = (X_1, \dots, X_n)$ with individual degree $B < \infty$ is formally defined as
 234 (where $c_d \in \mathbb{N}$):

235
$$\Phi(X_1, \dots, X_n) = \sum_{d \in \{0, \dots, B\}^D} c_d \cdot \prod_{t \in D} X_t^{d_t} \quad (1)$$

236 **Definition 2.1** (Standard Monomial Basis). The term $\prod_{t \in D} X_t^{d_t}$ in Eq. (1) is a monomial.
 237 A polynomial $\Phi(\mathbf{X})$ is in standard monomial basis (SMB) when we keep only the terms with
 238 $c_d \neq 0$ from Eq. (1).

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

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

243 As an example, the degree of the polynomial $X^2 + 2XY^2 + Y^2$ is 3. Product terms in lineage
 244 arise only from join operations (Fig. 1), so intuitively, the degree of a lineage polynomial
 245 is analogous to the largest number of joins needed to produce a result tuple. We call a
 246 polynomial $\Phi(\mathbf{X})$ a c -T1DB-lineage polynomial (or simply lineage polynomial), if there exists
 247 a \mathcal{RA}^+ query Q , c -T1DB \mathcal{D} , and result tuple t such that $\Phi(\mathbf{X}) = \Phi[Q, \mathcal{D}, t](\mathbf{X})$.

248 **2.1.1 c-T1DBs and 1-B1DBs**

249 An incomplete database Ω is a set of deterministic databases ω called possible worlds.

250 A c -T1DB \mathcal{D} is a pair $(\{0, \dots, c\}^D, \mathcal{P})$ such that $\{0, \dots, c\}^D$ is an incomplete database
 251 whose set of possible worlds is the $c + 1^n$ tuple/multiplicity combinations across all $t \in D$,
 252 where $|D| = n$, $D = \bigcup_{\mathbf{m} \in \{0, \dots, c\}^D, \mathbf{m}_i \geq 1} t$ is the set of possible tuples across possible worlds,
 253 and \mathcal{P} is a probability distribution over $\{0, \dots, c\}^D$.

254 A block independent database (B1DB) is a related probabilistic data model $\mathcal{D} = (\Omega, \mathcal{P})$
 255 such that the base set of tuples $D = \bigcup_{\omega \in \Omega, t \in \omega} t$ is partitioned into a set of n independent

For Sec 2.1 perhaps use a generic set of variables S & define everything in terms of that.

Use def eqs to define B1DB

no need to re-define

Define B1DB "directly" as we did for c-T1DB i.e. of the form $(\{0, \dots, c\}^D, \mathcal{P})$ allow \mathcal{P} to assign certain worlds 0 prob.

spell it out Maybe use D' as set of base tuples for B1DB to distinguish from D of c-T1DB & not for sub-subset Why not use ω (not world)?

→ Current defn for 1-BIDB seems to indicate $\{0,1\}^{\text{count}}$ which is **NOT** enough for the **redix**

no need to define this since we do not need to make it clear that we're allowing for mult. here

256 blocks $\{(b_t)_{t \in [n]}\}$ such that the set of tuples $\{(t_j)_{j \in [c]}\}$ in block b_t are disjoint from one another. This construction produces the set of possible worlds Ω that consists of all unique combinations of tuples in D with the constraint that for any $\omega \in \Omega$, no two tuples $t_j, t_{j'}, j \neq j'$ from the same block b_t exist together. A **c-BIDB** has the further requirement that each block has a multiplicity of at most c . We present a reduction that is useful in producing our results:

261 ► **Definition 2.3** (*c*-TIDB reduction). Given *c*-TIDB $\mathcal{D} = (\{0, \dots, c\}^D, \mathcal{P})$ let $\mathcal{D}' = (\Omega, \mathcal{P}')$ be the 1-BIDB obtained in the following manner: for each $t \in \mathcal{D}$, create block $b_t = \{ \langle t, j X_{t,j} \rangle_{j \in [c]} \}$, such that $X_{t,j} \in \{0, 1\}$. The probability distribution \mathcal{P}' is the one induced by $\mathbf{p} = ((p_{t,j})_{t \in \mathcal{D}, j \in [c]})$ and the BIDB disjoint requirement.

265 For the *c*-TIDB \mathcal{D} , each $X_t \in [c]$, while in the reduced 1-BIDB \mathcal{D}' , each $X_{t,j} \in \{0, 1\}$. Hence, in the setting of 1-BIDB, the base case of Fig. 1 now becomes $\Phi[R, D, t] = \sum_{j \in [c]} j X_{t,j}$. Then given the disjoint requirement and the semantics for constructing the lineage polynomial over a 1-BIDB, $\Phi[R, D', t]$ is of the same structure as the reformulated polynomial Φ_R of step i) from Definition 1.3, which then implies that $\tilde{\Phi}$ is the reduced polynomial that results from step ii) of Definition 1.3, and further that Lemma 1.4 immediately follows for 1-BIDB polynomials: $\mathbb{E}_{\mathbf{W} \sim \mathcal{P}'}[\Phi(\mathbf{W})] = \tilde{\Phi}(\mathbf{p})$.

Aaron says: @atri, not sure if \mathcal{P}' should be \mathcal{P}'' (in the above expectation) as discussed below. Since $\mathcal{P}' \equiv \mathcal{P}''$, then the proof still holds for Lemma 1.4, but maybe it is important to \mathcal{P}'' to drive the point home that we iterate over the all worlds set (as opposed to the set of possible worlds) when computing the expectation of a polynomial. Or maybe it suffices to note that $\mathcal{P}' \equiv \mathcal{P}''$.

273 Instead of looking only at the possible worlds of \mathcal{D} , one can consider all worlds, including those that cannot exist due to disjointness. The all worlds set can be modeled by $\mathbf{M} \in \{0, 1\}^{cn}$,³ such that $\mathbf{M}_{t,j} \in \mathbf{M}$ represents whether or not the multiplicity of t is j . We denote a probability distribution over all $\mathbf{M} \in \{0, 1\}^n$ as \mathcal{P}'' . When \mathcal{P}'' is the one induced from each $p_{t,j}$ while assigning $Pr[\mathbf{M}] = 0$ for any \mathbf{M} with $\mathbf{M}_{t,j} = \mathbf{M}_{t,j'} = 1$ for $j \neq j'$, we end up with a bijective mapping from \mathcal{P}' to \mathcal{P}'' , such that each mapping is equivalent, implying the distributions are equivalent. Appendix B.2 has more details.

280 Let $|\Phi|$ be the number of operators in Φ .
 281 ► **Corollary 2.4.** 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.

283 ~~Queries over probabilistic databases are evaluated using the so-called possible world semantics. 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 $\omega \in \Omega$: $\{Q(\omega) : \omega \in \Omega\}$.~~

287 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: $Pr[\omega \in \Omega] = \sum_{\omega' \in \Omega, Q(\omega') = Q(\omega)} Pr[\omega']$.

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

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

Why does BIDB def have variables in it?

(1) Refstate the full def of Φ (we have space for)

(2) This redix is incorrect. Define in terms of \mathcal{D} (ie base case should be $\Phi[R, D, t] = \sum_j j X_{t,j}$)

In some sense this section is not super relevant to paper. So putting all this in a separate sub-section makes more sense

use Φ for BIDB

need to explicitly state this out

Using Φ for both TIDB & BIDB

If you do not start from possible world semantics this is not an issue. You just start w/ \mathcal{P}''

This is what you can make the connection & possible world semantics

This should be in main text & highlighted

Note: Don't forget to change the opening of S1 to not use the term product distribution, but rather state that \mathcal{P} is a probability distribution.

293 ► **Proposition 2.5** (Expectation of polynomials). Given a bag-PDB $\mathcal{D} = (\Omega, \mathcal{P})$, \mathcal{RA}^+ query
 294 Q , and lineage polynomial $\Phi[Q, D, t]$ for arbitrary result tuple t , we have (denoting \mathbf{D} as the
 295 random variable over Ω): $\mathbb{E}_{\mathbf{D} \sim \mathcal{P}}[Q(\mathbf{D})(t)] = \mathbb{E}_{\mathbf{W} \sim \mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$.

296 A formal proof of Proposition 2.5 is given in Appendix B.3.⁴ We focus on the problem of
 297 computing $\mathbb{E}_{\mathcal{P}}[\Phi[Q, D, t](\mathbf{W})]$ from now on, assume implicit Q, D, t , and drop them from
 298 $\Phi[Q, D, t]$ (i.e., $\Phi(\mathbf{X})$ will denote a polynomial).

299 **2.2 Formalizing Problem 1.6**

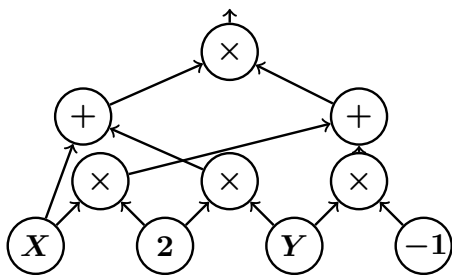
300 We represent lineage polynomials via arithmetic circuits [9], a standard way to represent
 301 polynomials over fields (particularly in the field of algebraic complexity) that we use for
 302 polynomials over \mathbb{N} in the obvious way. Since we are particularly using circuits to model
 303 lineage polynomials, we can refer to these circuits as lineage circuits. However, when the
 304 meaning is clear, we will drop the term lineage and only refer to them as circuits.

305 ► **Definition 2.6** (Circuit). A circuit C is a Directed Acyclic Graph (DAG) whose source
 306 gates (in degree of 0) consist of elements in either \mathbb{N} or \mathbf{X} . For each result tuple there exists
 307 one sink gate. The internal gates have binary input and are either sum (+) or product (\times)
 308 gates. Each gate has the following members: *type*, *partial*, *input*, *degree*, *Lweight*, and
 309 *Rweight*, where *type* is the value type $\{+, \times, \text{VAR}, \text{NUM}\}$ and *input* the list of inputs. Source
 310 gates have an extra member *val* storing the value. C_L (C_R) denotes the left (right) input of C .

Aaron says: Does the following matter, i.e., does it point anything out special for our research?

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

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



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

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 2.7** ($\text{POLY}(\cdot)$). 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:

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

⁴ Although Proposition 2.5 follows, e.g., as an obvious consequence of [28]’s Theorem 7.1, we are unaware of any formal proof for bag-probabilistic databases.

Handwritten notes:
 - "not defined" (crossed out)
 - "move to next section?"
 - "need a filler sentence removal reader about P1.6"
 - "make sure is for bag" (with arrow pointing to Definition 2.6)
 - "for w pr" (with arrow pointing to Definition 2.6)
 - "he" (with arrow pointing to Definition 2.6)

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

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

333 The circuit of Fig. 3 is an element of $\text{CSet}(2X^2 + 3XY - 2Y^2)$. One can think of
 334 $\text{CSet}(\Phi(\mathbf{X}))$ as the infinite set of circuits where for each element \mathcal{C} , $\text{POLY}(\mathcal{C}) = \Phi(\mathbf{X})$.

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

336 **Definition 2.9** (The Expected Result Multiplicity Problem). Let \mathcal{D} be an arbitrary BIDB-
 337 PDB and \mathbf{X} be the set of variables annotating tuples in D_Ω . Fix an \mathcal{RA}^+ query Q and a
 338 result tuple t . The EXPECTED RESULT MULTIPLICITY PROBLEM is defined as follows:

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

341 **2.3 Relationship to Deterministic Query Runtimes**

342 To decouple our results from specific join algorithms, we first abstract the cost of a join.

343 **Definition 2.10** (Join Cost). Denote by $T_{\text{join}}(R_1, \dots, R_m)$ the runtime of an algorithm
 344 for computing the m -ary join $R_1 \bowtie \dots \bowtie R_m$. We require only that the algorithm must
 345 enumerate its output, i.e., that $T_{\text{join}}(R_1, \dots, R_m) \geq |R_1 \bowtie \dots \bowtie R_m|$.

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

350
$$T_{\text{det}}(R, D) = |D \cdot R| \quad T_{\text{det}}(\sigma Q, D) = T_{\text{det}}(Q, D) \quad T_{\text{det}}(\pi Q, D) = T_{\text{det}}(Q, D) + |Q(D)|$$

351
$$T_{\text{det}}(Q \cup Q', D) = T_{\text{det}}(Q, D) + T_{\text{det}}(Q', D) + |Q(D)| + |Q'(D)|$$

352
$$T_{\text{det}}(Q_1 \bowtie \dots \bowtie Q_m, D) = T_{\text{det}}(Q_1, D) + \dots + T_{\text{det}}(Q_m, D) + T_{\text{join}}(Q_1(D), \dots, Q_m(D))$$

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

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

358 **3 Hardness of Exact Computation**

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

363 **3.1 Preliminaries**

364 Our hardness results are based on (exactly) counting the number of (not necessarily induced)
 365 subgraphs in G isomorphic to H . Let $\#(G, H)$ denote this quantity. We can think of H

Filler sentence to connect back to intro

make sure
intro has
a few
pts to
cover

Use D
instead
of D

correct
notation.