# Combining Sampling and Synopses with Worst-Case Optimal Runtime and Quality Guarantees for Graph Pattern Cardinality Estimation

Kyoungmin Kim[†]
George Fletcher[‡]

Hyeonji Kim[†]
Wook-Shin Han[†*]

Pohang University of Science and Technology (POSTECH), Korea[†], Eindhoven University of Technology (TU/e), Netherlands[‡]
{kmkim, hjkim, wshan}@dblab.postech.ac.kr[†], {g.h.l.fletcher}@tue.nl[‡]

## ABSTRACT

Graph pattern cardinality estimation is the problem of estimating the number of embeddings $|\mathcal{M}|$ of a query graph in a data graph. This fundamental problem arises, for example, during query planning in subgraph matching algorithms. There are two major approaches to solving the problem: sampling and synopsis. Synopsis (or summary)-based methods are fast and accurate if synopses capture information of graphs well. However, these methods suffer from large errors due to loss of information during summarization and inherent assumptions. Sampling-based methods are unbiased but suffer from large estimation variance due to large sample space.

To address these limitations, we propose ALLEY, a hybrid method that combines both sampling and synopses. ALLEY employs 1) a novel sampling strategy, *random walk with intersection*, which effectively reduces the sample space, 2) *branching* to further reduce variance, and 3) a novel mining approach that extracts and indexes *tangled* patterns as synopses which are inherently difficult to estimate by sampling. By using them in the online estimation phase, we can effectively reduce the sample space while still ensuring unbiasedness. We establish that ALLEY has worst-case optimal runtime and approximation quality guarantees for any given error bound $\epsilon$ and required confidence $\mu$. In addition to the theoretical aspect of ALLEY, our extensive experiments show that ALLEY outperforms the state-of-the-art methods by up to orders of magnitude higher accuracy with similar efficiency.

## 1 INTRODUCTION

Subgraph matching is one of the most fundamental and heavily researched types of graph querying [39]. Given a query graph $q$ and a data graph $g$, subgraph matching is the problem of finding the set $\mathcal{M}$ of all (isomorphic or homomorphic) embeddings of $q$ in $g$. Figure 1 shows an example of subgraph matching with 20 embeddings $\{u_1 \rightarrow v_{10}, u_2 \rightarrow v_{29}, u_3 \rightarrow v_{30}, u_4 \rightarrow v_{49}, u_5 \rightarrow v_{30}\}$, $\{u_1 \rightarrow v_{10}, u_2 \rightarrow v_{29}, u_3 \rightarrow v_{30}, u_4 \rightarrow v_{49}, u_5 \rightarrow v_{59}\}$, ..., $\{u_1 \rightarrow v_{10}, u_2 \rightarrow v_{29}, u_3 \rightarrow v_{30}, u_4 \rightarrow v_{49}, u_5 \rightarrow v_{77}\}$. Here, $u \rightarrow v$ denotes the mapping from a query vertex $u$ to a data vertex $v$. Subgraph matching has many well-known applications, including chemical compound search [49], protein interaction analysis [50], social analysis [15], and knowledge base queries [21].

As well as subgraph matching, graph pattern cardinality estimation (i.e., estimating the number of embeddings, $|\mathcal{M}|$) is an important problem. For example, estimates are used to determine the cost of query plans in subgraph matching systems [33]. Accurate estimation enhances the quality of plans while efficient estimation minimizes the query optimization overhead. Estimates can also be

used in fast approximate query answering when calculating the exact answer takes too long, and approximate results are sufficient for data analytics. An online aggregation method can even progress the estimation towards a more accurate answer [28].

In this paper, we focus on combining two major approaches – sampling and synopses – to solve the graph pattern cardinality estimation problem. Synopsis-based methods pre-build summary structures in the offline phase and use the structures to perform cardinality estimation in the online phase. Sampling-based methods perform sampling, calculate weights of samples, and aggregate these weights to estimate the cardinality in the online phase.

Synopsis-based methods have been mainly developed for RDF graphs. C-SET [35] summarizes the data graph into a set of star-shaped structures, while SUMRDF summarizes it into a smaller graph. These methods are fast and accurate if synopses capture information of data graphs well. However, it is a well-known problem that the information loss in summarization and the ad-hoc assumptions (e.g., uniformity and independence) in computing the estimates may produce large errors [27, 31, 38]. On various real and synthetic datasets, a very recent work by Park et al. [38] shows that these methods actually suffer from serious under-estimation problems due to these limitations.

Sampling-based methods have been mainly developed for relational data, but show good performance on graph data. Park and colleagues further show that, surprisingly, an online aggregation method designed for relational data, WANDERJOIN, significantly outperforms all techniques designed for graph data [38]. However, sampling-based methods suffer from large *sample space* that leads to large estimation variance and a high probability of sampling failures (i.e., samples have zero-weights). The failures lead to significant under-estimation in sampling-based methods, including WANDERJOIN, which we will show in Section 3. Here, the sample space is the set of all possible random walks following a particular sampling strategy and order. The set of candidates for each vertex/edge in order is called the *local* sample space. We illustrate these problems using an example in Figure 1 and WANDERJOIN.

WANDERJOIN for relational data can be translated into an estimator for graph data which performs *edge-at-a-time* sampling, by regarding query edges as relations and overlapping query vertices as join keys. For example, the query graph in Figure 1a is equivalent to a join graph in Figure 1c. WANDERJOIN first determines the order of query edges, e.g., $\langle (u_1, u_2), (u_2, u_4), (u_1, u_4), (u_1, u_3), (u_3, u_4), (u_4, u_5) \rangle$. Following the order, it breaks cycles and transforms the query graph into a tree with split vertices, as in Figure 1b. For example, $u_4$ is first split into $u_4$ and $u_4'$ since the first three edges form a cycle and $u_4$ is the latest visited. $u_4$ is then split into $u_4$ and $u_4''$ again due to $(u_3, u_4)$. The sampling order is defined on the edges

---
*corresponding author

(a) Query graph $q$.  (b) Query tree of $q$.  (c) Join graph of $q$.  (d) Data graph $g$.  (e) Statistics of edges in $g$.

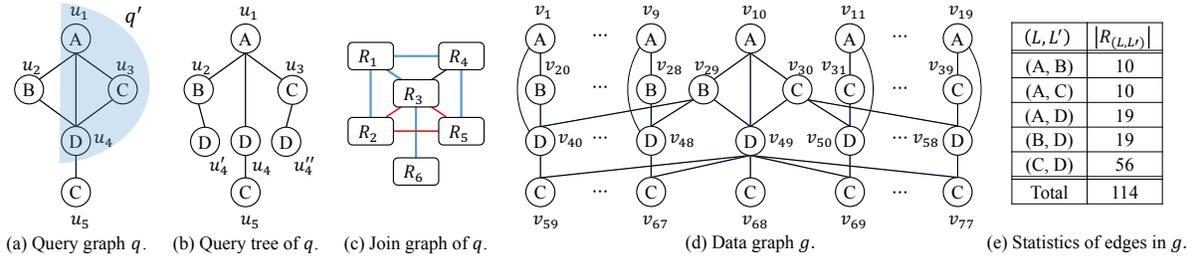| $(L, L')$ | $|R_{(L,L')}|$ |
|---|---|
| (A, B) | 10 |
| (A, C) | 10 |
| (A, D) | 19 |
| (B, D) | 19 |
| (C, D) | 56 |
| Total | 114 |

**Figure 1: A running example.**

of the tree as $\langle (u_1, u_2), (u_2, u'_4), (u_1, u_4), (u_1, u_3), (u_3, u''_4), (u_4, u_5) \rangle$. Blue edges in Figure 1c represent the corresponding join path.

WANDERJOIN then randomly walks on the data graph following the sampling order. For instance, $(v_1, v_{20})$ is first uniformly sampled for $(u_1, u_2)$ out of ten data edges that match $(u_1, u_2)$, i.e., $\{(v_1, v_{20}), (v_2, v_{21}), ..., (v_{10}, v_{29})\}$. For the second edge $(u_2, u'_4)$, the walk continues from $v_{20}$. While the walk continues to sample $(v_{20}, v_{40})$ and $(v_1, v_{40})$ for $(u_2, u'_4)$ and $(u_1, u_4)$, respectively, it cannot continue to sample for $(u_1, u_3)$ since $v_1$ has no candidate that matches $(u_1, u_3)$ in its adjacency list. We say that the random walk (or sampling) failed at $(u_1, u_3)$. If WANDERJOIN starts the walk from $(v_{10}, v_{29})$ and succeeds to sample a data edge for each query edge, it finally checks the join conditions between the split query vertices $u_4$, $u'_4$, and $u''_4$ (equivalent to red edges in Figure 1c). For example, if $(v_{10}, v_{49})$ and $(v_{29}, v_{40})$ are sampled for $(u_1, u_4)$ and $(u_2, u'_4)$, respectively, the join fails between the split vertices $u_4$ and $u'_4$ since $v_{49}$ is sampled for $u_4$ while $v_{40}$ is for $u'_4$. Similarly, the condition between $u_4$ and $u''_4$ is checked. If any of the join conditions fails, we also say that the sampling has failed. Sampling succeeds if and only if the sampled edges form an embedding of $q$ in $g$.

WANDERJOIN suffers from a large sample space for two main reasons. First, in performing random walks, it considers only one query edge at a time, which blindly leads to zero candidates for later edges in the sampling order. Second, it breaks cycles in $q$ and samples multiple data vertices for the same query vertex (e.g., $\{u_4, u'_4, u''_4\}$) which must eventually be joined. This can lead to a significant number of failures for cyclic queries.

Since performing a successful sampling (i.e., sampling an embedding) is important in sampling-based estimators, a naive approach to combine sampling and synopses would pre-compute embeddings by subgraph matching and use them as sample space. However, this is an NP-hard problem, and since we do not know which patterns will be queried online, we might have to calculate embeddings for all possible patterns that appear in the data graph. This is, of course, an infeasible approach.

A second naive approach would to pre-compute *domains* of small patterns (e.g., up to five edges), inspired by work for frequent pattern mining in a single large graph [14]. Here, the domain is defined as the set of data vertices that participate in an embedding of a pattern. For example, the domain of $u_1$ in $q'$ in Figure 1a is $\{v_{10}\}$, since only $v_{10}$ participates in an embedding of $q'$ in $g$. By indexing the domains in the offline phase and using them as the local sample spaces of vertices in $q$ in the online phase, we can safely prune out the candidates that lead to sampling failures.

However, frequent pattern mining methods (aka *fpms*) suffers from scalability for large heterogeneous graphs [2, 22, 44]. Enumerating all small-size patterns and calculating their domains using NP-hard subgraph matching incurs significant overhead. Especially in our scenario, the scalability issue would be even more onerous since 1) we do not differentiate between frequent or infrequent patterns so the number of patterns to index can be much larger, 2) we have to materialize all the domains to use them online, rather than repeatedly discarding the searched ones as in *fpms*, and 3) at least a superset of each domain should be indexed (no false-negative candidates), where *fpms* are optimized to probe a subset of domain and determines whether a pattern is frequent or not.

**Contributions.** To address the above limitations, we present AL-LEY, a hybrid method that carefully combines both sampling and synopses. We first design a novel sampling strategy that reduces the sample space. Instead of using sampling only in the online phase, we interleave sampling and mining domains in the offline phase, and index only hard patterns that are inherently difficult to handle by sampling. We call such patterns *tangled*, if it is hard to reduce sample space, and the rate of sampling failures exceeds some threshold. By indexing only tangled patterns as important synopses, we can greatly improve scalability of the offline phase. Again, in the online phase, we use the domains of tangled subqueries to prune local sample spaces of queries. In short, we improve the efficiency of building synopses by combining with sampling, also the accuracy of sampling-based estimation by combining with synopses.

We briefly illustrate our ideas using examples, starting from a novel sampling strategy, called *random walk with intersection*. Instead of considering only one edge at a time as in WANDERJOIN, ALLEY samples a vertex at a time from a reduced number of candidates by considering all its incident edges. For example, if the sampling starts from $u_1$ in Figure 1a, ALLEY intersects the candidates for $(u_1, u_2)$, $(u_1, u_3)$, and $(u_1, u_4)$ on $u_1$, i.e., $\{v_1, v_2, ..., v_{10}\} \cap \{v_{10}, v_{11}, ..., v_{19}\} \cap \{v_1, v_2, ..., v_{19}\} = \{v_{10}\}$. We readily obtain the candidates that are more likely to lead to a successful sampling. Along with intersections, ALLEY uses *branching* to bound variance and increase efficiency. Branching is to sample multiple data vertices *without replacement* for a query vertex $u$ if the number of candidates for $u$ is large. These enable a significant accuracy increase with a reasonable trade-off of efficiency.

For indexing tangled patterns, we use a novel mining approach, called *walk-fail-then-calculate*. ALLEY performs random walks using the domains of smaller patterns and calculates the domains if the rate of sampling failures exceeds some threshold. For a data graph $g_2$ in Figure 2 and a subquery $q'$ in Figure 1a, the number of embeddings of $q'$ is just one, i.e., $\{u_1 \to v_1, u_3 \to v_{12}, u_4 \to v_{40}\}$. However, the
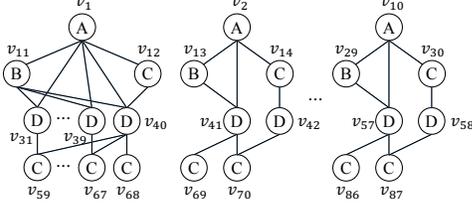
Figure 2: Another data graph $g_2$ with same edge statistics.

number of candidates for any vertex $u$ in $q'$, i.e., the local sample space for $u$, is at least 10, even if we consider the intersections. The domains of smaller subqueries, e.g., $\{(u_1, u_3), (u_1, u_4)\}$, have no effect on pruning local sample spaces of $q'$. Therefore, a random walk will fail with a probability near 0.9. If we set the threshold for sampling failures less than 0.9, say 0.8, such patterns can be considered tangled. Then, we index the domains for $u_1$, $u_3$, and $u_4$ of $q'$, which are $\{v_1\}$, $\{v_{12}\}$, and $\{v_{40}\}$, respectively. Along with our novel mining approach, we propose three optimization techniques based on the fact that we only need to obtain reduced local sample spaces, not necessarily the exact domains of each tangled pattern, and we can determine tangled patterns without performing too many random walks.

In addition to its practicality, we establish that ALLEY has formal guarantees on estimation error, confidence, and time complexity as follows: given an error bound $\epsilon$ and confidence $\mu$, ALLEY guarantees $\Pr(|Z - |\mathcal{M}|| < \epsilon \cdot |\mathcal{M}|) > \mu$ always in $O(AGM(q))$ time, where $Z$ is the random variable for the cardinality estimate and $AGM(q)$ is the worst-case optimal bound of $|\mathcal{M}|$ [36]. Achieving this guarantee is meaningful since it bounds both estimation error and runtime. In addition, an emerging class of subgraph matching algorithms that run in $O(AGM(q))$ time has recently been studied actively [32, 36]. In summary, in this paper

- We show that small sample space leads to small variance and fewer sampling failures, resulting in high accuracy in sampling-based estimators (Section 3).
- We present ALLEY, an accurate and efficient graph pattern cardinality estimator based on a novel sampling strategy (Sections 4-5).
- We propose a novel mining method to increase effectiveness for tangled patterns and to make ALLEY a hybrid method that combines both sampling and synopses (Section 6).
- We prove the probabilistic theoretical guarantees of ALLEY in worst-case optimal time (Section 7).
- With extensive experiments, we establish that ALLEY consistently and significantly outperforms all the state-of-the-art estimators (Section 8). Specifically, ALLEY outperforms WANDERJOIN by up to two orders of magnitude in terms of accuracy with similar efficiency.

## 2 BACKGROUND

### 2.1 Problem Definition

An undirected labeled graph is a triple, $(V, E, L)$, such that $V$ is a set of vertices, $E$ is a set of edges, and $L$ is a label function that maps a vertex to a set of labels. Given two labeled graphs $q$ and $g$, we can define a match or an embedding between $q$ and $g$ as follows [25].

**Definition 1.** A graph $q = (V_q, E_q, L_q)$ is homomorphic to a subgraph of a data graph $g = (V_g, E_g, L_g)$ if there is a mapping (or an embedding) $m : V_q \rightarrow V_g$ such that 1) $\forall u \in V_q, L_q(u) \subseteq L_g(m(u))$, and 2) $\forall (u, u') \in E_q, (m(u), m(u')) \in E_g$.

Here, $(u, u')$ represents an undirected edge between vertex $u$ and $u'$. Given a query graph $q$ and a data graph $g$, let $\mathcal{M}$ denote the set of all embeddings of $q$ in $g$. Finding $\mathcal{M}$ is called the *subgraph matching* problem, and estimating $|\mathcal{M}|$ is called the *graph pattern cardinality estimation* problem.

We assume that $q$ is connected, otherwise we can estimate the cardinalities of disjoint subqueries and multiply them. For ease of explanation, we assume edges are not labeled, and the query and the data graph are not multi-graphs. The actual implementation of ALLEY supports directed, edge-labeled, and multi-graphs as well.

### 2.2 Notation

We next present notation used throughout the paper (Table 1). Set operations are often used, for example, we regard $e = \{u, u'\}$ for an edge $e = (u, u')$. An embedding $m$ is also a set of mappings $u \mapsto v$ where $u \in V_q$ and $v \in V_g$. If $u \mapsto v \in m$, we write $m(u) = v$.

We also use notation from relational algebra. Given a query edge $e$, $R_e$ denotes the set of data edges that match $e$. Hence, $R_e = R_{(u,u')} := \{(v, v') \in E_g \mid L_q(u) \subseteq L_g(v) \land L_q(u') \subseteq L_g(v')\}$. Since $R_{(u,u')}$ is determined by labels of $u$ and $u'$, we also use $R_{(L,L')}$ ($L = L_q(u)$, $L' = L_q(u')$). This corresponds to a binary relation in relational algebra, regarding $u$ and $u'$ as two attributes and each $(v, v') \in R_{(u,u')}$ as a tuple. Note that, a subgraph matching query can be translated into a natural join query $\bowtie_{e \in E_q} R_e$ as in Figure 1c.

We use a set projection of $R_{(u,u')}$ onto $u'$ as $\pi_{u'}(R_{(u,u')}) := \{v' \in V_g \mid \exists v \in V_g : (v, v') \in R_{(u,u')}\}$ which we simplify as $V_u^{u'}$, representing the set of data vertices that match $u'$ and have an incident edge in $R_{(u,u')}$. We also use $adj_u^{u'}(v) := \{v' \in V_u^{u'} \mid (v, v') \in E_g\}$, which is the set of adjacent vertices of $v$ in $V_u^{u'}$. For example, $V_{u_2}^{u_1} = \{v_1, v_2, ..., v_{10}\}$ and $adj_{u_2}^{u_1}(v_{20}) = \{v_1\}$ in Figure 1.

**Table 1: Notation used in the paper.**

| | |
|---|---|
| $\mathcal{M}$ | set of embeddings of $q$ in $g$ |
| $\mathcal{P}, P_u$ | sample space and local sanple space for $u$ |
| $p, w(p)$ | potential embedding in $\mathcal{P}$ and weight of $p$ |
| $\mathbb{I}(p)$ | indicator variable for $p$ whether $p \in \mathcal{M}$ or not |
| $R_e$ | set of data edges that match a query edge $e$ |
| $V_u^{u'}$ | set of data vertices that match $u'$ and have an incident edge in $R_{(u,u')}$ |
| $adj_u^{u'}(v)$ | set of adjacent vertices of $v$ in $V_u^{u'}$ |
| $o$ | sampling order of query vertices (or edges) |
| $D_q(u)$ | domain of a vertex $u$ in a pattern $q$ |

### 2.3 AGM Bound

The AGM bound [8] is defined for a given query graph $q$ and a set of data graphs $G_q$, denoted as $AGM(q)$. Here, $G_q$ consists of all possible data graphs with the same $|R_e|$ for each $e \in E_q$ (e.g., Figures 1d and 2 have the same edge statistics in Figure 1e.)

The AGM bound is worst-case optimal in the sense that it is a *tight* upper bound of $|\mathcal{M}|$, i.e., there exists a data graph $g^* \in G_q$ that satisfies $|\mathcal{M}| = AGM(q)$. If a subgraph matching algorithm

runs in $O(AGM(q))$ time, the algorithm is also called worst-case optimal. Recent systems [3, 6, 23] implement such algorithms. As in typical worst-case optimal algorithms [7, 37], we omit $|V_q|$ and $\log|E_g|$ term in $O$ notation.

## 3  HT ESTIMATOR

Horvits-Thompson (HT) estimators are a class of unbiased sampling-based estimators [40]. In our problem, an HT estimator samples data vertices/edges for query vertices/edges and checks whether the sampled vertices/edges form an embedding of $q$ in $g$. Here, sampled vertices/edges may not form an embedding but any (partial) mapping, so we call each sample a *potential* embedding of $q$ as [7].

After sampling a potential embedding $p$, the estimator assigns a *weight*, $w(p)$, as $w(p) = 1/\Pr(p) \cdot \mathbb{I}(p)$. Here, $\Pr(p)$ represents the probability of sampling $p$ dependent on a sampling method, and $\mathbb{I}(p)$ is an indicator variable for sampling failure of $p$, i.e., $\mathbb{I}(p) = 1$ if $p \in \mathcal{M}$ ($p$ is an embedding of $q$) and $\mathbb{I}(p) = 0$ otherwise. This weight is used as the estimated cardinality. For instance, using WANDERJOIN as an HT estimator in our running example in Figure 1, $\Pr(m) = \frac{1}{10} \cdot \frac{1}{10} \cdot \frac{1}{1} \cdot \frac{1}{1} \cdot \frac{1}{1} \cdot \frac{1}{20}$ for any $m \in \mathcal{M}$ whose weight $w(m)$ is $1/\Pr(m) \cdot \mathbb{I}(p) = 2,000 \cdot 1 = 2,000$. The sampling is repeated $s$ times, and the final estimate is calculated as the average of weights, e.g., if $s = 1,000$ and the sampling failed for 999 times out of 1,000 repetitions, the estimate becomes $\frac{2,000}{1,000} = 2$ in the above example. Here, $s$ is called the *sample size*. The whole process is given in Algorithm 1 with an additional step of determining the *sampling order* of query vertices/edges.

---

**Algorithm 1.** SimpleHTEstimator($q, g, s$)

**Input:** A query graph $q$, a data graph $g$, and a sample size $s$
1: $o \leftarrow$ CHOOSESAMPLINGORDER$(q, g)$
2: $sum \leftarrow 0$
3: **foreach** ($i = 1$ to $s$) **do**
4:      $p \leftarrow$ SAMPLEPOTENTIALEMBEDDING$(o, q, g)$
5:      $sum \leftarrow sum + w(p)$
6: **return** $sum/s$

---

### 3.1  Characteristics of HT Estimators

HT estimators have two important characteristics: they are unbiased and there is a fundamental trade-off between estimation accuracy and efficiency.

The unbiasedness can be explained as follows. Let $Z_i$ denote the random variable for $w(p_i)$ where $p_i$ is the $i$'th potential embedding sampled ($1 \le i \le s$). Then, let $Z = \sum_{i=1}^{s} Z_i/s$ denote the random variable for the estimated cardinality.

$$\mathbb{E}[Z_i] = \sum_p \Pr(p) \cdot \frac{1}{\Pr(p)} \cdot \mathbb{I}(p) = \sum_p \mathbb{I}(p) = \sum_{p \in \mathcal{M}} 1 = |\mathcal{M}| \quad (1)$$

$$\mathbb{E}[Z] = \mathbb{E}\left[\frac{\sum_{i=1}^{s} Z_i}{s}\right] = \sum_{i=1}^{s} \frac{\mathbb{E}[Z_i]}{s} = |\mathcal{M}| \quad (2)$$

(2) holds since the $Z_i$'s are identical and independent. The accuracy-efficiency trade-off can be explained as follows:

$$\mathbb{V}\text{ar}[Z] = \mathbb{V}\text{ar}\left[\frac{\sum_{i=1}^{s} Z_i}{s}\right] = \sum_{i=1}^{s} \frac{\mathbb{V}\text{ar}[Z_i]}{s^2} = \frac{\mathbb{V}\text{ar}[Z_i]}{s} \quad (3)$$

Thus, increasing the sample size $s$ increases the runtime (i.e., decreases the efficiency) of an HT estimator but reduces the estimation variance (i.e., increases the accuracy).

### 3.2  Sample Space

Each HT estimator has a different sampling strategy that determines the set of potential embeddings to be sampled. For example, an HT estimator $T_1$ samples any data edge from $E_g$ for each query edge, while another HT estimator $T_2$ samples a data edge from $R_e$ for each query edge $e$. Obviously, $T_1$ can sample (useless) potential embeddings which $T_2$ does not. If $\mathcal{P}_T$ denotes the set of all potential embeddings that can be sampled by an estimator $T$, then $\mathcal{P}_{T_2} \subset \mathcal{P}_{T_1}$. Here, $\mathcal{P}$ is called the *sample space*. If $\mathcal{M} \subseteq \mathcal{P}_T$, then we say that $T$ is *consistent*. So far, we assumed that the estimator is consistent.

Not surprisingly, the sample space plays a critical role in estimation accuracy. Since $\mathbb{V}\text{ar}[Z_i] = \mathbb{E}[Z_i^2] - \mathbb{E}[Z_i]^2$, the following equation holds for a consistent HT estimator.

$$\mathbb{V}\text{ar}[Z_i] = \left(\sum_p \Pr(p) \cdot \left(\frac{1}{\Pr(p)} \cdot \mathbb{I}(p)\right)^2\right) - \mathbb{E}[Z_i]^2$$
$$= \sum_p \frac{1}{\Pr(p)} \cdot \mathbb{I}(p) - |\mathcal{M}|^2 = \sum_{p \in \mathcal{M}} \frac{1}{\Pr(p)} - |\mathcal{M}|^2. \quad (4)$$

That is, the estimation variance depends on the probabilities of sampling embeddings of $q$. If the probabilities are low, then the variance will be large. To increase such probabilities, we must reduce the sample space size $|\mathcal{P}|$ to decrease the probabilities of sampling $p \notin \mathcal{M}$ and increase the probabilities of sampling $p \in \mathcal{M}$. If we assume uniform sampling, i.e., $\Pr(p) = 1/|\mathcal{P}|$, (4) can be approximated as follows:

$$\mathbb{V}\text{ar}[Z_i] \approx |\mathcal{M}| \cdot |\mathcal{P}| - |\mathcal{M}|^2. \quad (5)$$

Thus, the variance is roughly proportional to the size of the sample space. If $|\mathcal{P}|$ is significantly larger than $|\mathcal{M}|$, it means that most of the potential embeddings sampled are not embeddings. In Line 4 of Algorithm 1, it is highly likely that $p \notin \mathcal{M}$, resulting in sampling failures. Since the weight is zero for such samples, this can cause a serious under-estimation problem as shown in [38]. For example, $|\mathcal{P}_{T_1}|$ is $114^6 \approx 2 \cdot 10^{12}$ since there are 114 data edges and six query edges in Figure 1. For $T_2$, $|\mathcal{P}_{T_2}|$ is $\prod_{e \in E_q} |R_e| \approx 10^8$. Note that both $T_1$ and $T_2$ are consistent, but have significantly large sample spaces compared to $|\mathcal{M}| = 20$. In the following sections, we show how we reduce the sample space.

## 4  OVERVIEW OF SOLUTION

As explained in Section 1, a naive combination approach would mine embeddings or domains of all small-size patterns that appear in the data graph. However, this might require a tremendous amount of computation. Instead, our approach interleaves sampling during the mining to improve mining efficiency, and utilizes the mined domains to improve online estimation accuracy.

Figure 3 shows the overall architecture of ALLEY. In the offline phase, the *tangled pattern index* is built based on the data graph $g$ only. The goal is to automatically find all tangled patterns, which are expected to result in low accuracy in the online phase. The patterns are investigated from smaller to larger ones, that is, each pattern is extended from its sub-patterns. A new pattern $q$ is input to our
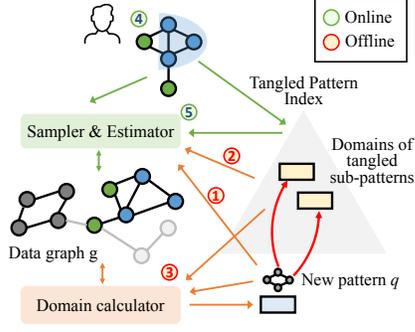
**Figure 3: Architecture of ALLEY.**

sampling-based estimator (①) along with the indexed domains of its sub-patterns (②), which are used to prune the sample space for $q$. If the ratio of sampling failures for $q$ exceeds some threshold, i.e., the pruning effect is poor, the domains of $q$ are calculated (③) and indexed as synopses. Thus, Alley interleaves sampling with mining, that is, sampling is performed to judge the tangledness of a pattern. In the online phase, given a query graph specified by a user (④), we use the built index to prune the sample space in estimating the cardinality of the query (⑤). We first explain online sampling-based estimation and move to building synopses offline.

## 5 SAMPLING-BASED ESTIMATION

This section explains the main algorithm of online estimation (Algorithm 2). The algorithm consists of three main parts, 1) choosing sampling order, 2) searching domains, and 3) random walk with intersection.

---

**Algorithm 2.** ALLEY$(q, g, s, I)$

---
**Input:** A query graph $q$, a data graph $g$, a sample size $s$, and a
      tangled pattern index $I$ of $g$

1: $o \leftarrow$ CHOOSESAMPLINGORDER$(q, g)$
2: $\{D_{q_j}\} \leftarrow$ SEARCHDOMAINSRECURSIVE$(q, o, I)$
3: $sum, count \leftarrow 0, 0$ /* sum of weights and # of calls   */
4: $p \leftarrow \emptyset$ /* initialize potential embedding       */
5: **while** $(s > 0)$ **do**
6:    $w_1, s_1 \leftarrow$ RANDOMWALKWITHINTERSECT$(q, g, \{D_{q_j}\}, o, p, 1)$
7:    $sum \leftarrow sum + w_1$
8:    $s \leftarrow s - s_1; count \leftarrow count + 1$
9: **return** $sum/count$

---

**CHOOSESAMPLINGORDER.** We choose the sampling order $o$ (Line 1) by adopting a simple greedy approach considering the selectivity of each vertex. Specifically, we define the rank for each vertex $u$ as $rank(u) := \min_{(u,u') \in E_q} |V_{u'}^u|$. This means the set of candidates that match $u$ considering the query edge $e$ only. Then, we select the minimum-rank vertex as the starting query vertex. For a tie, we consider two additional ranks, one is $\min_{e \in E_q, e \ni u} |R_e|$ and the other one is the degree of $u$. For remaining vertices, we repeatedly select the vertex who has the largest number of neighboring selected vertices. For a tie, we choose one with a smaller rank. While selecting a good sampling order is an important research problem as noted in typical subgraph matching papers [18, 51], we leave it for future investigation.

**Example 1.** *In Figure 1, $rank(u_1) = \min_{(u_1,u)} |V_u^{u_1}| = \min\{|V_{u_2}^{u_1}|,$ $|V_{u_3}^{u_1}|, |V_{u_4}^{u_1}|\} = \min\{10, 10, 19\} = 10$. Similarly, $rank(u_2) = rank(u_3) = rank(u_4) = 10$, and $rank(u_5) = 29$. We then compare $\min_{e \ni u} |R_e|$ for the four tied vertices except $u_5$. Again, there's a tie between $u_1, u_2,$ and $u_3$. We finally choose $u_1$ as the starting query vertex since it has the largest degree (=3) among $u_1, u_2,$ and $u_3$. In subsequent examples, we use $o = \langle u_1, u_2, u_4, u_3, u_5 \rangle$ as the sampling order of ALLEY. Since using the data graph in Figure 2 results in the same rank values, the same sampling order $o$ is chosen.*

**SEARCHDOMAINSRECURSIVE.** Next, for each vertex $u$, we find a domain in the index by searching for subqueries $q_j$ of $q$ that contain $u$ (Line 2). The searched domain is used for pruning a local sample space for $u$. The detailed algorithm will be explained in Section 6.2.

**Example 2.** *Assume that the index of the graph in Figure 2 stores only $q'$ in Figure 1a and its domains for simplicity. Since $q'$ is a subquery of $q$, we search for $q'$ in the index. Then, we can use the searched domains $D_{q'}$ of $q'$ which are $\{v_1\}, \{v_{12}\},$ and $\{v_{40}\},$ as the candidates of $u_1, u_3,$ and $u_4$, respectively.*

**RANDOMWALKWITHINTERSECT.** We then iteratively sample potential embeddings. As explained in Algorithm 1, we calculate the weight $w_1$ until we run out of sample size (Lines 5-8), and finally return the average of $w_1$ as the estimated cardinality (Line 9). However, the difference with Algorithm 1 is that each (recursive) call to RANDOMWALKWITHINTERSECT can sample $s_i$ ($\geq 1$) potential embeddings (we call this *branching*), and $w_i$ is the *aggregated weight* of $s_i$ potential embeddings at each recursion depth $i$ ($1 \leq i \leq |V_q|$). This variation is essential to guarantee the worst-case optimal time complexity and approximation quality of the algorithm, which will be explained in Section 7. Note that ALLEY is a stack of HT estimators and thus unbiased, also proved in Section 7.

### 5.1 Random Walk with Intersection

We now explain how RANDOMWALKWITHINTERSECT works in Algorithm 3. At recursion depth $i$, the function first calculates the set $P_{o_i}$ of candidates for $o_i$ given the current potential embedding ($\{o_j \rightarrow v_j\}_{j<i}$). It then randomly walks to candidate vertices, extends the potential embedding, and continues to the next recursion. For simplicity, we omit in Algorithm 3 that, if the total number of random walks exceeds $s$, we stop further random walks and calculate the estimate from the collected random walks so far.

**INTERSECT.** First, we calculate $P_{o_i}$ which is the local sample space for $o_i$ using a multi-way intersection (Line 1). Given the current potential embedding $p$, we perform an intersection $\left(\cap_{(o_i,o_j) \in E_q, j<i} adj_{o_j}^{o_i}(p(o_j))\right) \cap \left(\cap_{(o_i,o_j) \in E_q j \geq i} V_{o_j}^{o_i}\right)$. If $D_{q'}(o_i)$ has been searched, since we assure $D_{q'}(o_i) \subseteq \left(\cap_{(o_i,u) \in E_{q'}} V_u^{o_i}\right)$, we can remove such $V_u^{o_i}$ terms (already considered in calculating $D_{q'}(o_i)$) from our initial equation and intersect with $D_{q'}(o_i)$ instead.

If $P_{o_i}$ is empty, it indicates that the sampling has failed at $o_i$. We directly return $w_i = 0$ and $s_i = 1$ (Lines 2-3). If $i = |V_q|$, the sampling has succeeded with $|P_{o_i}|$ embeddings found (Lines 4-5). In Section 7, we prove that for each of $v \in P_{o_{|V_q|}}, p \cup \{o_{|V_q|} \rightarrow v\}$ forms an embedding of $q$.

**Example 3.** *For the graph in Figure 1 and $i = 1$ where $o_1 = u_1$, we compute $P_{o_1}$ by intersecting three sets, i.e., $V_{u_2}^{u_1} \cap V_{u_3}^{u_1} \cap V_{u_4}^{u_1} = \{v_1, v_2, ..., v_{10}\} \cap \{v_{10}, v_{11}, ..., v_{19}\} \cap \{v_1, v_2, ..., v_{19}\} = \{v_{10}\}$. As in*

**Algorithm 3.** RANDOMWALKWITHINTERSECT($q, g, \{D_{q'}\}, o, p, i$)

**Input:** A query graph $q$, a data graph $g$, a sampling order $o$,
current potential embedding $p$, and current recursion depth $i$

1: $P_{o_i} \leftarrow$ INTERSECT($g, p, \{D_{q'}\}, o_i$) /* multi-way intersection */

2: **if** ($P_{o_i} = \emptyset$) **then**

3:      **return** $0, 1$

4: **if** ($i = |V_q|$) **then**

5:      **return** $|P_{o_i}|, 1$

     /* branching, sample without replacement */

6: $C \leftarrow$ CHOOSERANDOMVERTICESWOR($P_{o_i}, \lceil b \cdot |P_{o_i}| \rceil$)

7: $sum, num \leftarrow 0$ /* sum of weights and # of recursion */

8: **foreach** ($v \in C$) **do**

9:      $p \leftarrow p \cup \{o_i \rightarrow v_i\}$ /* $p(o_i) = v_i$ */

10:      $w_{i+1}, s_{i+1} \leftarrow$
       RANDOMWALKWITHINTERSECT($q, g, \{D_{q'}\}, o, p, i+1$)

11:      $num \leftarrow num + s_{i+1}; sum \leftarrow sum + w_{i+1}$

12:      $p \leftarrow p \setminus \{o_i \rightarrow v\}$ /* $p(o_i) = NIL$ */

13: **return** $|P_{o_i}| \cdot sum/|C|, num$

---

this example, we can effectively reduce the local sample space through 3-way intersection.

**Example 4.** *We now consider $g_2$ in Figure 2 and the query graph in Figure 1a. We compute an intersection of three sets for $P_{o_1}$, i.e., $V_{u_2}^{u_1} \cap V_{u_3}^{u_1} \cap V_{u_4}^{u_1} = \{v_1, v_2, ..., v_{10}\} \cap \{v_1, v_2, ..., v_{10}\} \cap \{v_1, v_2, ..., v_{10}\} = \{v_1, v_2, ..., v_{10}\}$. Unfortunately, the intersection does not prune candidates in contrast to Example 3.*

**Example 5.** *Unlike Example 4, if $D_{q'}(u_1) = \{v_{10}\}$ is available, we only need to intersect $D_{q'}(u_1)$ with $V_{u_2}^{u_1}$ to find the local sample space $P_{o_1}$ for $u_1$, since the other two edges $(u_1, u_3)$ and $(u_1, u_4)$ are already considered in $D_{q'}(u_1)$. Then, we calculate $D_{q'}(u_1) \cap V_{u_2}^{u_1} = \{v_{10}\} \cap \{v_1, v_2, ..., v_{10}\} = \{v_{10}\}$. This example shows that using domains for tangled patterns can further reduce local sample space.*

We then sample a vertex from $P_{o_i}$ and walk to the vertex, extending $p$ and calling RANDOMWALKWITHINTERSECT recursively (Lines 8-12). Each call returns $w_{i+1}$ and $s_{i+1}$ where $w_{i+1}$ denotes the aggregated weight of $s_{i+1}$ potential embeddings sampled in the call. This $w_{i+1}$ is averaged and multiplied with $|P_{o_i}|$, then returned with the summation of $s_{i+1}$ (Line 13).

**Edge/Vertex-at-a-Time Sampling.** We define edge-at-a-time sampling as sampling a data edge for each query edge, and vertex-at-a-time sampling as sampling a data vertex for each query vertex considering all its incident edges. To sample an embedding that matches the query, each strategy requires sampling $|E_q|$ and $|V_q|$ times, respectively. In general, $|E_q|$ is larger than $|V_q|$ (e.g., $|E_q| = 6$, $|V_q| = 5$ in Figure 1a). By considering all incident edges of each query vertex, the number of candidates is smaller in vertex-at-a-time sampling, although it's bounded to the number of outgoing edges ($|adj_{o_i}^{o_j}|$) as in edge-at-a-time sampling except for $|P_{o_1}|$, which is bounded to $|V_{o_1}^u|$. In Section 8.3, we report the number of candidates in ALLEY (vertex-at-a-time) compared to WANDERJOIN's (edge-at-a-time).

**Time and Space Complexity.** In terms of efficiency, the intersection is the main bottleneck in Algorithm 3. Therefore, the implementation of an intersection requires special care in order to obtain an efficient estimator. We implement the intersection by referencing the implementation of EMPTYHEADED [3] without using its

data compression techniques for a fair comparison with others. Specifically, for a multi-way intersection, we start from the smallest set and intersect it with the second smallest set and repeat this pairwise intersection. Given $k$ sets of size $N_1, N_2, ..., N_k$ (sorted in ascending order), the intersection takes $O(kN_1 \log N_k)$ time and requires $O(N_1)$ space. Meanwhile, the intersection itself has an early stopping effect; that is, it will catch sampling failures at an early stage and avoid further random walks that would eventually fail. As query graphs become larger, the effectiveness of early stopping becomes be greater. We show this in our experiments.

## 5.2 Branching

Branching is to sample data vertices in $P_{o_i}$ for $o_i$ multiple times (i.e., $\lceil b \cdot |P_{o_i}| \rceil$ times) *without* replacement (Line 6). The generated potential embeddings share $p(o_1), p(o_2), ..., $ and $p(o_{i-1})$ but have different $p(o_i)$. The constant $b \in (0, 1)$ is called the *branching factor* which is a hyperparameter. Note that SSTE [7] uses a different branching technique for cycle (sub)queries only with branching factor $1/\sqrt{|E_g|}$ and sampling *with* replacement.

Making branches bounds the estimation variance of $w_i$ regardless of the local sample space $P_{o_i}$ (Proposition 2 in Section 7). Here we explain the advantage of branching at a high level.

**Example 6.** *Figure 4 shows recursion trees of RANDOMWALKWITH-INTERSECT. Here, we consider the case in Example 4. A path from the root to each node shows a potential embedding $p_i$, and $w_i, s_i$ values to be returned are attached at right or center (aggregated). Succeeded walks are colored in green, while failed walks are colored in red.*

*If we do not make any branch (i.e., $b = \frac{1}{\inf}$), the estimator will call an independent recursion path, such as the one in Figure 4a. This can result in multiple identical random walks.*

*In contrast, with branching, we can force a different recursion path to be performed for each recursion level. Assume that $b$ is $\frac{1}{5}$ (Figure 4b). We now sample twice for $o_1$ from $P_{o_1} = \{v_1, v_2, ..., v_{10}\}$, since $\lceil b \cdot |P_{o_1}| \rceil = 2$. Examples of call paths and return values for sampling $v_1$ and $v_2$ for $o_1$, are given in Figure 4b. Here, sampling twice for $o_1$ (and $o_3$) increases the probability of sampling $v_1$ (and $v_{40}$) which leads to an embedding. Note that in Example 3, we always obtain the exact cardinality $|\mathcal{M}|$ with one possible random walk since $|P_{o_1}| = |P_{o_2}| = |P_{o_3}| = |P_{o_4}| = 1$. Thus, the variance of $w_1$ to $w_5$ is zero. For a larger $P_{o_i}$, we must sample multiple times in order to bound each $w_i$, which is the key to proving Proposition 2.*



(a) No branching.
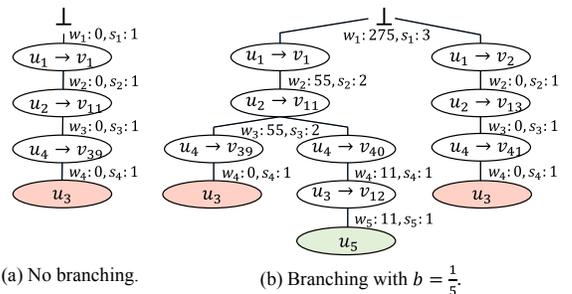(b) Branching with $b = \frac{1}{5}$.

**Figure 4: Recursive calls using $g_2$ in Figure 2.**

Branching can also mitigate the overhead of intersections. For example in Figure 4b, we can reuse the intersection result of $P_{o_3}$

for sampling $v_{39}$ and $v_{40}$ given $p = \{u_1 \to v_1, u_2 \to v_{11}\}$ (the left subtree). Given the same sample size $s$, this has the effect of reducing the number of intersections. This gain results in bounded time complexity below. Note that, in general, a larger branching factor results in a smaller bound of variance and intersection overhead, but it would require a sufficiently large sample size.

**Time Complexity.** If we 1) replace Lines 8-12 in Algorithm 3 by an exhaustive search (i.e., iterate over all $v \in P_{o_i}$ instead of sampling) and 2) replace Lines 4-5 by an exhaustive search for all $v \in P_{o_n}$ and adding $\{o_n \to v\}$ to $p$, then RANDOMWALKWITHINTERSECT starting from $i = 1$ becomes a specialization of Generic-Join [36]; each potential embedding that can be obtained at $i = n$ with $P_{o_n} \neq \emptyset$ is an embedding in $\mathcal{M}$. Then it is trivial that sampling without replacement makes Algorithm 3 a *part of* a worst-case optimal algorithm, without searching for the same embedding multiple times. In other words, sampling *with* replacement cannot guarantee this; a data vertex that leads to many embeddings in the exhaustive search can be sampled multiple times, resulting in more than $O(AGM(q))$ potential embeddings sampled.

## 5.3 Implementation details

In order to increase efficiency, we applied the following optimization techniques. First, $P_{o_1}$ is the same for each call to RANDOMWALK-WITHINTERSECT at Line 6 of Algorithm 2. Thus we compute $P_{o_1}$ only once. Second, we remove the sets in intersections that do not affect the result. In Figure 2, if we start a random walk from $u_2$ (assuming that a data graph different from $g_1$ in Figure 1 is given), we have to intersect three sets, i.e., $V_{u_1}^{u_2}$, $V_{u_3}^{u_2}$, and $V_{u_4}^{u_2}$. However, $L_q(u_3) = L_q(u_4)$ so $V_{u_3}^{u_2} = V_{u_4}^{u_2}$ by the definition of $V_u^{u'}$. Thus, we only intersect the first two sets and remove the last to avoid the redundant computation. Even when we start from $u_3$ and walk to $u_2$, calculating the candidates for $u_2$, we can ignore intersecting with $V_{u_4}^{u_2} = V_{u_3}^{u_2}$ since it is a superset of $adj_{u_3}^{u_2}(p(u_3))$ for any $p$.

To further increase efficiency, one can apply the following techniques which would be an interesting future work: 1) base and state, or a hybrid representation to compress the sets into bitmaps and apply bit intersections [3, 19], 2) SIMD instructions to reduce the number of operations [3, 19] and 3) vertex signatures and multi-core GPUs to highly parallelize the intersections [48].

**Branching.** We apply the following techniques to increase the effectiveness of branching. First, we adaptively choose $b$ depending on $|P_{o_i}|$. If $|P_{o_i}|$ is too small so is $\lceil b \cdot |P_{o_i}| \rceil$, the actual confidence of the estimated cardinality might be too low. Therefore, we use larger $b$ for small $|P_{o_i}|$ as shown in Figure 5. Second, we skip unnecessary branches for a query vertex if all its adjacent vertices are sampled. For example, if the sampling order is given as $\langle u_1, u_2, u_4, u_3, u_5 \rangle$ in Figure 1, it is unnecessary to make branches for $u_3$ and $u_5$ for any data graph, since sampling any candidate from their local sample spaces does not affect the final estimate. Otherwise, if a query vertex $u$ has an unsampled adjacent vertex, different candidates can result in different estimates. Such classification of query vertices has been studied in a subgraph enumeration paper [24].

## 6 TANGLED PATTERN INDEX

We now explain how to build and search in a *tangled pattern index*. We have shown in Section 1 that tangled patterns in queries can
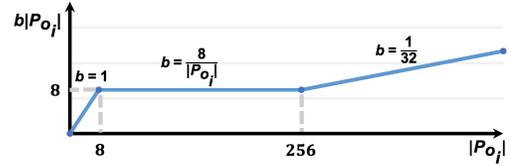


**Figure 5: Adaptive branch size.**

increase the difficulty of sampling and are, in general, inherently difficult to handle by any sampling-based methods. To resolve the problem, we propose an index for such tangled patterns.

**Definition 2.** *Given a data graph $g$, a tangled pattern index for $g$ is a DAG $:= (N, E)$ where 1) each node $n_q \in N$ corresponds to a (tangled or untangled) pattern $q$ appearing in $g$, and 2) each edge $(n_q, n_{q'}, [u, u']) \in E$ with label $[u, u']$ represents a PointingTo relationship between $u \in V_q$ and $u' \in V_{q'}$.*

Here, if $D_q(u)$ denotes the domain of $u$ in $q$, then *PointingTo* relationship from $u$ to $u'$ represents that i) $D_{q'}(u') \supseteq D_q(u)$ from the anti-monotone property (Lemma 1) [14], ii) $D_{q'}(u')$ is used to prune the local sample space of $u$ in $q$ (i.e., $P_u$) during the random walk, and iii) the rate of sampling failures of $q$ does not exceed a predefined threshold (i.e., $q$ is not tangled). If a node $n_q$ has no outgoing *PointingTo* edge, $q$ is *tangled* and the domain of each vertex in $q$ is materialized. Figure 6 shows an example index, where $q_4$ is not tangled and $q_7$ is tangled.

**Lemma 1.** *Given a data graph $g$ and a pattern $q'$ that matches a subgraph of $q$ by a mapping $m$, $D_{q'}(u') \supseteq D_q(m(u'))$ for $u' \in V_{q'}$.*

## 6.1 Building Index

In the offline phase, we start from mining since we do not know which patterns are tangled. However, existing frequent pattern mining work [2, 14] cannot be directly applied to our problem. Their purpose is to mine patterns whose support values are larger than a given threshold, so they adopt optimization techniques to prune infrequent patterns earlier. [2] estimates the frequency of a pattern and prunes if the estimated frequency is low. [14] models frequency evaluation as a constraint satisfaction problem (CSP), and computes the minimal set of domains that are sufficient for CSP to determine whether a pattern is frequent. This avoids finding the whole domains. In contrast, we need to enumerate tangled patterns whose frequencies would be small. Moreover, we must calculate the whole domains in order to prevent false-negatives in the sampling. Those incur extra computational overheads to pattern mining which already has significant time complexity.

To tackle this challenging problem, we present a novel, efficient mining approach, "walk-fail-then-calculate," that interleaves sampling with mining, which allows ALLEY to autonomously determine and index the tangled patterns, avoiding enormous computation for frequent patterns. Like standard frequent subgraph mining algorithms, ALLEY mines from smaller to larger patterns. However, instead of mining all patterns, ALLEY first performs random walks (Algorithm 2) regarding each pattern as a query $q$ and the domains of sub-patterns as $D_{q_j}$. If the random walks fail with a high chance (above a threshold), ALLEY determines that this pattern $q$ is hard to estimate its cardinality from the current index thus extends the index by mining $q$.

Algorithm 4 shows this procedure. It grows patterns from size one to size $maxN$ where size is the number of edges (Lines 4-5). Before calculating the domains for a pattern $q$ of size $N$, we first search for its subgraphs $\{q_i\}$ where each $q_i$ has size $N-1$. If any of the $q_i$ has an empty domain (i.e., $D_{q_i}(u) = \emptyset$ for any $u \in V_{q_i}$), we can skip processing $q$ according to the anti-monotone property; $q$ also has at least one empty domain (Lines 7-8). Otherwise, we start "walk" to determine whether $q$ is tangled (Lines 9-10). Specifically, for each $u \in V_q$, we first select a set of candidates from the set of domains $\{D_{q_i}(u') \mid u' \in V_{q_i} \wedge \exists m : m(u') = u\}$, which will be used in random walks. Among such domains, we select the minimum-size domain, $\hat{D}_q(u) := D_{\hat{q}}(\hat{u})$ where $(\hat{q}, \hat{u}) = \arg\min_{q_i, u' \in V_{q_i} \wedge \exists m: m(u') = u} |D_{q_i}(u')|$, as the set of candidates since it has the highest pruning power. Then, we can get $\hat{D}_q = \{\hat{D}_q(u)\}_{u \in V_q}$ (Line 9). Note that we might refer to different $\hat{q}$ for each $u \in V_q$. According to the anti-monotone property, $\hat{D}_q(u) \supseteq D_q(u)$. Then, we perform random walks by running ALLEY (Line 10). Specifically, we regard $q$ as a query graph and run Algorithm 2. We directly use $\hat{D}_q$ as $D_{q_j}$ instead of running SEARCHDOMAINSRECURSIVE in Line 2 of Algorithm 2.

If ALLEY "fails" (i.e., its failure rate $r$ exceeds a given threshold $\zeta$), we finally "calculate" the domains of $q$ (Line 12). Here, $r$ is the ratio of the number of calls to RANDOMWALKWITHINTERSECT that returns $w_1 = 0$ in Line 6 of Algorithm 2. Again, when calculating the domains of $q$, we use $\hat{D}_q(u)$ as the candidates for $D_q(u)$ according to the property. If ALLEY does not fail, we simply maintain the *PointingTo* edges $(n_q, n_{\hat{q}}, [u, \hat{u}])$ $(\forall u \in V_q)$ and regard $D_q$ as $\hat{D}_q$ in Lines 16-17.

**Example 7.** *Figure 6 shows part of a tangled pattern index for $g_2$ in Figure 2. The 2-edge patterns $q_4$, $q_5$, and $q_6$, are extended from 1-edge patterns, $q_1$, $q_2$, and $q_3$. The $q_7$ is extended from those three 2-edge patterns. Each vertex of $q_4$ points to another vertex in a 1-edge pattern, which represents that $D_q(u)$ has been set to point $\hat{D}_q(u)$ for $u \in V_{q_4}$. This is due to the low failure rate of random walking for $q_4$ using the domains of $q_1$ and $q_2$. Note that, since $D_{q_1}(\text{Ⓒ})$ is smaller than $D_{q_2}(\text{Ⓒ})$, Ⓒ in $q_4$ points to Ⓒ in $q_1$. We omit the pointers of $q_5$ and $q_6$. In contrast, random walks for $q_7$ would result in a large failure rate since the domains of $q_7$ are much smaller than those of its subgraphs. Thus, $q_7$ is determined as tangled and $D_{q_7}$ is calculated.*
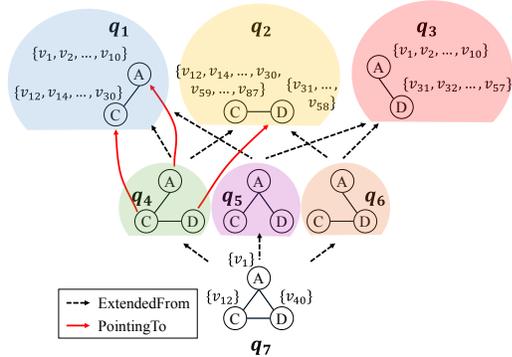


**Figure 6: An example index for $g_2$ in Figure 2.**

In conclusion, our approach has five strong points. First, it reduces the number of expensive domain calculations by placing cheap random walks as filters. Second, for tangled patterns $\{q_t\}$,

---

**Algorithm 4.** TANGLEDPATTERNMINING$(g, maxN, s, \zeta)$

**Input:** A data graph $g$, maximum number of edges of a pattern $maxN$, a sample size $s$, and a threshold $\zeta$
/* $I_N$ stores N-edge patterns and their domains    */
1: $I_N \leftarrow \emptyset$ for $1 \le N \le maxN$
2: **populate** $I_1$ with all 1-edge patterns using labels in $g$
3: **foreach** $(1 \le N \le maxN - 1)$ **do**
4:     **foreach** $(q \in I_N)$ **do**
           /* extend $q$ with an edge, generate multiple (N+1)-size patterns using labels in $g$    */
5:         $I_{N+1} \leftarrow I_{N+1} \cup$ EXTENDPATTERN$(q, L_g)$
6:     **foreach** $(q \in I_{N+1})$ **do**
7:         **if** (HASSUBGRAPHWITHEMPTYDOMAIN$(q, I_N)$) **then**
8:             **remove** $q$ from $I_{N+1}$; **continue**
           /* $\hat{D}_q(u)$ for each $u \in V_q$ points to a precomputed domain $D_{\hat{q}}(\hat{u})$ $(\hat{q} \in \{q_1, q_2, ...\})$    */
9:         $\hat{D}_q \leftarrow$ GETMINIMUMDOMAINSFROMSUBGRAPHS$(q, I_N)$
10:        $r \leftarrow$ GETFAILURERATE(ALLEY, $q, g, s, \hat{D}_q$)
11:        **if** $(r > \zeta)$ **then**
12:            $D_q \leftarrow$ CALCULATEDOMAINS$(q, g)$
13:            **if** $(\exists u \in V_q : D_q(u) = \emptyset)$ **then**
14:                **remove** $q$ from $I_{N+1}$
15:                **continue**
16:        **else**
17:            $D_q \leftarrow \hat{D}_q$ /* just maintain the pointers    */
18: **return** $\{I_2, ..., I_{maxN}\}$

---

$D_{q_t}(u)$ would be much smaller than $\hat{D}_{q_t}(u)$ since the sampling has failed using $\hat{D}_{q_t}(u)$. This has the effect of bounding the index size similar to using a discriminative ratio as in [46]. Third, in order to increase the online effectiveness of ALLEY, we allow ALLEY itself determine hard cases offline. Fourth, we can control the offline overhead and online effectiveness of our index by changing the threshold $\zeta$. A large $\zeta$ let ALLEY reduce the number of domain calculations offline but causes it to use large domains online. Fifth, our index can be generalized to increase the performance of other sampling-based methods, which would be interesting future work.

**Maintenance.** For updates, the index is built by using the RUNSTAT command, as in existing DBMSs. In order to minimize maintenance overhead, we do not update the index whenever there is an update to the graph database. When a considerable amount of updates have occurred, say 5-10% of data edges, a DBA needs to rebuild the index by issuing RUNSTAT again. We confirmed that the estimation accuracy does not deteriorate much even if 5-10% of data edges are inserted through experiments (Figure 23 in Section 8.4). We can further optimize performance by incrementally rebuilding the index.

**Optimizations.** Even though our filtering approach reduces the overhead of domain computation, problems can still occur if $g$ contains many labels and the number of possible patterns explodes. We first reduce the overhead of each call to CALCULATEDOMAINS by implementing a specialized algorithm for computing domains instead of enumerating all embeddings. We utilize dynamic programming to avoid redundant computation, which runs in linear time (i.e., $O(|V_q| \sum_{u \in V_q} \hat{D}_q(u))$) for tree queries. To take advantage of dynamic programming, we extend the optimization to cyclic queries allowing non-linear time and approximate computation;

domains can contain false positive vertices. We also apply the techniques to cache intersection results [32], skip particular query and data vertices [47], and exploit neighborhood equivalence class [20] to reduce duplicated computation. We stop extending $q$ if the size of $\hat{D}_q(u)$ is already too small. To avoid intractable computation of some patterns, we set a threshold for the size of Cartesian product.

Then, we reduce the number of patterns by pruning and grouping labels as follows: First, if $g$ contains edge labels, we use edge labels only without vertex labels. Otherwise, we use vertex labels. If the number of (vertex or edge) labels is large, we group labels and limit the number of groups. This still might be too large, however, but further decreasing the number of groups may result in too many false positives. In such cases, we use hierarchical grouping, e.g., group labels for path patterns, then make supergroups for general tree patterns and finally, make supergroups of the supergroups for general graph patterns. The motivation is that, while labels tend to determine the tangledness of simple queries (e.g., paths), topologies tend to determine the tangledness of complex queries (e.g., cliques).

Finally, instead of using a predefined sample size, we continuously update the failure rate after calculating each $w_1$ in Algorithm 2 and stop if the rate converges. This prevents an unnecessarily large sample size. Also, we only focus on whether $w_1$ is zero or not, so we stop each call to RandomWalkWithIntersect if we sample an embedding and $w_1$ is guaranteed to be nonzero.

**Complexity.** Given a data graph $G = (V_g, E_g, L_g)$, there can be $\varepsilon \times |L_g|^{(maxN+1)}$ tangled patterns where $\varepsilon$ is the ratio of tangled patterns over all patterns up to $maxN$ edges. The factor $\varepsilon$ can significantly reduce time/space complexity. CalculateDomains dominates the other functions which is bounded by $O(|E_g|^{maxN})$ time complexity. Therefore, the time complexity of Algorithm 4 is $O(\varepsilon \times |L_g|^{(maxN+1)} \times |E_g|^{maxN})$, while the space complexity is $O(\varepsilon \times |L_g|^{(maxN+1)} \times |V_g| \times (maxN + 1))$. The complexities of maintenance are similar to above, but we only need to calculate the domains of affected tangled patterns. Note that these are loose upper bounds since we do not consider 1) the impact of search space reduction by using pointed tangled patterns recursively, 2) the important characteristic of tangled patterns that their domain size is much smaller than $|V_g|$, and 3) the optimization techniques. For example, in YAGO dataset, $\varepsilon$ is $10^{-5}$ and the actual average space per pattern is 4.4K whereas $|V_g| \times (maxN + 1)$ is 76.8M.

### 6.2 Searching in Index

In the online phase, we use the built tangled pattern index to reduce the sample space of a given query $q$ (Line 2 of Algorithm 2). The idea is to search for domains of subqueries $\{q_1, q_2, ...\}$ of $q$, and use each domain $D_{q_j}(m^{-1}(o_i))$ to prune candidates $P_{o_i}$ in Line 1 of Algorithm 3, where $m$ is the mapping from $V_{q_j}$ to $V_q$. Due to the anti-monotone property, $D_q(o_i) \subseteq D_{q_j}(m^{-1}(o_i))$. Note that $D_{q_j}$ might not be the actual domains of $q_j$, but it points to one of its subgraphs, as described in Section 6.1.

Among multiple $D_{q_j}(m^{-1}(o_i))$ for different subqueries $q_j$, we search for the minimum domain $D_{q_j}(m^{-1}(o_i))$ (Lines 3-4 in Algorithm 5). Note that the initial SearchDomainsRecursive is called right after ChooseSamplingOrder. Since enumerating all subqueries of $q$ can be costly for large $q$, we greedily search for maximal subqueries (with size $\leq maxN$) that contain the first $k$ vertices

in the given sampling order (Lines 1-4). After searching for such maximal subqueries, we remove the $k$ vertices from $q$ (Line 5) and continue to search from the new first vertex by calling PriorityFirstSearchFrom (Line 7). This is repeated until $q$ is empty.

---

**Algorithm 5.** SearchDomainsRecursive($q, o, I$)

**Input:** A query graph $q$, a sampling order $o$, and a pre-built index $I$ with maximum size $maxN$

1: **foreach** (maximal connected subquery $\tilde{q}$ of $q$ such that $|E_{\tilde{q}}| \leq maxN$ and $\tilde{q}$ contains the first $k$ vertices in $o$) **do**
2: $\quad D_{\tilde{q}} \leftarrow$ SearchDomains($\tilde{q}, I$)
3: $\quad$ **foreach** ($u \in V_{\tilde{q}}$) **do**
4: $\quad\quad D_{q_{min}}(m(u)) \leftarrow \arg\min (|D_{q_{min}}(m(u))|, |D_{\tilde{q}}(u)|)$
5: **remove** the $k$ vertices from $q$
6: **if** ($V_q \neq \emptyset$) **then**
    /* traverse from $o_{k+1}$ with $o$ as priority ($o_i < o_{i+1}$) */
7: $\quad o' \leftarrow$ PriorityFirstSearchFrom($o_{k+1}, o$)
8: $\quad$ SearchDomainsRecursive($q, o', I$)
9: **return** $D_{q_{min}}$

---

**Time complexity.** Given a query $q$, the number of recursive calls is $\lceil |V_q|/2 \rceil$ in the worst case. For each recursive call, PriorityFirstSearchFrom takes $O(|V_q|)$ time (Line 7), and the total number of loops (Lines 1-4) is $O(2^{maxN})$. Calculating the BFS code of $\tilde{q}$ dominates SearchDomains and takes $O(2^{|V_{\tilde{q}}|})$ where $|V_{\tilde{q}}|$ is always less than $maxN$ (Line 2). Therefore, the time complexity of Algorithm 5 is $O(|V_q|^2 + 4^{maxN}|V_q|)$, which is quadratic to the query size.

## 7 WORST-CASE OPTIMAL RUNTIME AND APPROXIMATION QUALITY GUARANTEES

In this section, we establish the runtime and approximation quality guarantees of Alley. First, we formally define the problem of graph pattern cardinality estimation with probabilistic guarantees in worst-case optimal time, which we have discussed in Section 1.

**Definition 3.** *For a given error bound $\epsilon$ and a confidence $\mu \in (0, 1)$, if the random variable for the estimated cardinality $Z$ satisfies $\Pr(|Z - |\mathcal{M}|| < \epsilon \cdot |\mathcal{M}|) > \mu$, then the estimation is a $(1 \pm \epsilon)$-approximation of $|\mathcal{M}|$.*

**Theorem 1.** *Alley performs $(1 \pm \epsilon)$-approximation of $|\mathcal{M}|$ in $O(AGM(q))$ time.*

While the accuracy of sampling-based estimators increases as the sample size increases, the theoretical guarantees of Alley in Theorem 1 indicate that it is enough to set the sample size to $O(AGM(q))$ in order to perform $(1 \pm \epsilon)$-approximation of $|\mathcal{M}|$, for any $\epsilon$ and $\mu$. Still, smaller $\epsilon$ and larger $\mu$ will increase the constant factor in $O(\cdot)$. The proof of Theorem 1 is based on the following lemmas and propositions, which is analogous to SSTE [7].

**Lemma 2.** *Alley is consistent, i.e., $\mathcal{M} \subseteq \mathcal{P}$.*

Proof. We use proof by contradiction. Assume that there is an embedding $m \in \mathcal{M}$ that cannot be sampled by Alley. If $m = \{o_1 \rightarrow v_1, o_2 \rightarrow v_2, ..., o_n \rightarrow v_n\}$ ($n = |V_q|$), then there exists the smallest $i \leq n$ such that $v_i \notin P_{o_i}$. Here, $P_{o_i}$ is calculated from $p = m_i$ where $m_i$ is the restriction of $m$ on $o_1, o_2, ..., o_{i-1}$. Since $v \notin P_{o_i}$, there must be an edge $e^* = (o_k, o_i) \in E_q$ such that $v_i \notin adj_{o_k}^{o_i}(p(o_k))$ if $i > k$ and $v_i \notin V_{o_k}^{o_i}$ if $i < k$.

i) If $i > k$, then $adj_{o_k}^{o_i}(p(o_k)) = adj_{o_k}^{o_i}(v_k)$ and $v_i$ must be in this adjacency list, otherwise there is no edge between $v_k$ and $v_i$ that matches $(o_k, o_i)$. Thus, $m$ cannot be an embedding of $q$, which is a contradiction.

ii) If $k > i$, $v_i \notin V_{o_k}^{o_i}$ means that $v_i$ has no incident edge $(v, v_i)$ that matches $e^*$. Thus, $m$ cannot be an embedding of $q$, which is a contradiction. $\qquad\square$

**Lemma 3.** *ALLEY is a stack of simple HT estimators.*

While a simple HT estimator performs $s$ independent sampling and aggregates their weights at once (Algorithm 1), ALLEY interleaves sampling a data vertex and aggregating weights (Lines 8-12 of Algorithm 3). In that sense, ALLEY can be regarded as a stack of simple HT estimators which returns $w_i$ as results. Then, a natural question arises. What values does $w_i$ estimate for? That is, what is the expectation $\mathbb{E}[w_i]$ of $w_i$ returned by each RANDOMWALKWITH-INTERSECT?

To answer the question, we define the following random variables. Let $a_i$ be the random variable for $|P_{o_i}|$. Let $t_i$ be the random variable for branch size, i.e., $t_i = \lceil b \cdot |P_{o_i}| \rceil$. We then recursively define $Y_i|p_i$ as the random variable for $w_i$; $i$ is the recursion depth, and $p_i$ is the current potential mapping that maps $o_1, ..., o_{i-1}$ to data vertices.

$$Y_i|p_i = \begin{cases} a_i & i = n \\ a_i \cdot \frac{1}{t_i} \sum_{k=1}^{t_i} Y_{i+1}^k|p_{i+1} & i < n \end{cases} \quad (6)$$

Here, superscript $k$ distinguishes $t_i$ identically distributed random variables for $Y_{i+1}|p_{i+1}$. Note that these are not independent due to our sampling without replacement policy.

**Lemma 4.** $\mathbb{E}[Y_i|p_i] = c(q, g|p_i)$. *Here, $c(q, g|p_i)$ denotes the number of embeddings of $q$ in $g$ that contain the vertices and edges specified by $p_i$.*

PROOF. We use proof by induction.

For base case $i = n$, $\mathbb{E}[Y_n|p_n] = a_n = |P_{o_n}|$. Since each $v \in P_{o_n}$ has an incident edge that matches $e$ for every $e \in E_q, e \ni o_n$, $p_n \cup \{o_n \to v\}$ is an embedding of $q$ in $g$, i.e., $a_n \leq c(q, g|p_n)$. Since ALLEY is consistent by Lemma 2, $a_n \geq c(q, g|p_n)$. Combining these two gives $a_n = c(q, g|p_n)$.

For inductive case $i < n$,

$$\mathbb{E}[Y_i | p_i]$$
$$= \sum_{v \in P_{o_i}} \Pr(v | p_i) \cdot a_i \cdot \frac{1}{t_i} \cdot \sum_{k=1}^{t_i} \mathbb{E}[Y_{i+1}^k | p_i \cup \{o_i \to v\}]$$
$$\overset{(1)}{=} \sum_{v \in P_{o_i}} \frac{1}{t_i} \cdot \sum_{k=1}^{t_i} \mathbb{E}[Y_{i+1}^k | p_i \cup \{o_i \to v\}]$$
$$\overset{(2)}{=} \sum_{v \in P_{o_i}} \mathbb{E}[Y_{i+1} | p_i \cup \{o_i \to v\}]$$
$$\overset{(3)}{=} \sum_{v \in P_{o_i}} c(q, g | p_i \cup \{o_i \to v\})$$
$$\overset{(4)}{=} c(q, g | p_i).$$

(1) holds since $\Pr(v|p_i) = 1/a_i$. (2) holds since all $Y_{i+1}^k$'s are identical having the same expectation. (3) holds from the induction

hypothesis. (4) holds since summing up $c(q, g|p_i \cup \{o_i \to v\})$ for every possible $v \in P_{o_i}$ results in counting all embeddings specified by $p_i$.

$\qquad\square$

**Proposition 1.** $\mathbb{E}[Y_1] = |\mathcal{M}|$.

Proposition 1 can be readily obtained by Lemma 4 as $\mathbb{E}[Y_1] = \mathbb{E}[Y_1|p_1] = c(q, g|p_1) = c(q, g) = |\mathcal{M}|$. Therefore, ALLEY preserves the unbiasedness of the HT estimator. Furthermore, we bound the estimation variance of ALLEY.

**Proposition 2.** $\mathbb{V}\text{ar}[Y_1] \leq \frac{b}{1-b} \cdot (\frac{1}{b^n} - 1) \cdot |\mathcal{M}|^2$.

PROOF. Again, we use proof by induction and show that $\mathbb{V}\text{ar}[Y_i|p_i] \leq \frac{b}{1-b} \cdot (\frac{1}{b^{n-i+1}} - 1) \cdot c(q, g|p_i)^2$. If this holds, $\mathbb{V}\text{ar}[Y_1] = \mathbb{V}\text{ar}[Y_1|p_1] \leq \frac{b}{1-b} \cdot (\frac{1}{b^n} - 1) \cdot c(q, g|p_1)^2 = \frac{b}{1-b} \cdot (\frac{1}{b^n} - 1) \cdot |\mathcal{M}|^2$, completing the proof of the proposition.

For base case ($i = n$), $\mathbb{V}\text{ar}[Y_n|p_n] = 0$ as $Y_n|p_n$ is constant as $a_n$. Therefore, $\mathbb{V}\text{ar}[Y_n|p_n] \leq \frac{b}{1-b} \cdot (\frac{1}{b} - 1) \cdot c(q, g|p_n)$.

For inductive case ($i < n$), $\mathbb{V}\text{ar}[Y_i|p_i] = \mathbb{E}[\mathbb{V}\text{ar}[Y_i|p_i, v]] + \mathbb{V}\text{ar}[\mathbb{E}[Y_i|p_i, v]]$ ($v \in P_{o_i}$) by the law of total variance [34]. We bound the second term first.

$$\mathbb{V}\text{ar}[\mathbb{E}[Y_i | p_i, v]]$$
$$= \mathbb{V}\text{ar}[\mathbb{E}[Y_i | p_i \cup \{o_i \to v\}]]$$
$$\overset{(1)}{=} \mathbb{V}\text{ar}[c(q, g | p_i \cup \{o_i \to v\})]$$
$$\leq \mathbb{E}[c(q, g | p_i \cup \{o_i \to v\})^2]$$
$$= \sum_{v \in P_{o_i}} \Pr(v | p_i) \cdot c(q, g | p_i \cup \{o_i \to v\})^2$$
$$\overset{(2)}{\leq} \sum_{v \in P_{o_i}} c(q, g | p_i \cup \{o_i \to v\})^2$$
$$\leq \left( \sum_{v \in P_{o_i}} c(q | p_i \cup \{o_i \to v\}) \right)^2$$
$$= c(q, g | p_i)^2.$$

(1) holds from the inductive proof of Proposition 1. (2) holds since $\Pr(v|p_i) \leq 1$. Now we bound the first term.

$$\mathbb{E}[\mathbb{V}\text{ar}[Y_i | p_i, v]]$$
$$= \sum_{v \in P_{o_i}} \Pr(v | p_i) \cdot \mathbb{V}\text{ar}[Y_i | p_i, v]$$
$$= \sum_{v \in P_{o_i}} \Pr(v | p_i) \cdot \mathbb{V}\text{ar}[Y_i | p_i \cup \{o_i \to v\}]$$
$$= \sum_{v \in P_{o_i}} \Pr(v | p_i) \cdot a_i^2 \cdot \frac{1}{t_i^2} \cdot \mathbb{V}\text{ar}\left[ \sum_{k=1}^{t_i} Y_{i+1}^k | p_i \cup \{o_i \to v\} \right]$$
$$= \sum_{v \in P_{o_i}} \Pr(v | p_i) \cdot a_i^2 \cdot \frac{1}{t_i^2} \cdot \left( \sum_{k=1}^{t_i} \mathbb{V}\text{ar}[Y_{i+1}^k | p_i \cup \{o_i \to v\}] \right.$$

$$+ \sum_{k=1}^{t_i} \sum_{l=1, l \neq k}^{t_i} \mathbb{C}\mathrm{ov}\left(Y_{i+1}^k \mid p_i \cup \{o_i \rightarrow v\}, Y_{i+1}^l \mid p_i \cup \{o_i \rightarrow v\}\right)\Big)$$

$$\overset{(1)}{\leq} \sum_{v \in P_{o_i}} \Pr(v \mid p_i) \cdot a_i^2 \cdot \frac{1}{t_i^2} \cdot \sum_{k=1}^{t_i} \mathbb{V}\mathrm{ar}[Y_{i+1}^k \mid p_i \cup \{o_i \rightarrow v\}]$$

$$\overset{(2)}{=} \sum_{v \in P_{o_i}} a_i \cdot \frac{1}{t_i} \cdot \mathbb{V}\mathrm{ar}[Y_{i+1} \mid p_i \cup \{o_i \rightarrow v\}]$$

$$\overset{(3)}{\leq} \sum_{v \in P_{o_i}} \frac{1}{b} \cdot \mathbb{V}\mathrm{ar}[Y_{i+1} \mid p_i \cup \{o_i \rightarrow v\}]$$

$$\overset{(4)}{\leq} \sum_{v \in P_{o_i}} \frac{1}{b} \cdot \frac{b}{1-b} \cdot \left(\frac{1}{b^{n-i}} - 1\right) \cdot c(q, g \mid p_i \cup \{o_i \rightarrow v\})^2$$

$$= \frac{1}{b} \cdot \frac{b}{1-b} \cdot \left(\frac{1}{b^{n-i}} - 1\right) \cdot \sum_{v \in P_{o_i}} c(q, g \mid p_i \cup \{o_i \rightarrow v\})^2$$

$$\leq \frac{1}{b} \cdot \frac{b}{1-b} \cdot \left(\frac{1}{b^{n-i}} - 1\right) \cdot \left(\sum_{v \in P_{o_i}} c(q, g \mid p_i \cup \{o_i \rightarrow v\})\right)^2$$

$$= \frac{1}{1-b} \cdot \left(\frac{1}{b^{n-i}} - 1\right) \cdot c(q, g \mid p_i)^2.$$

(1) holds since the covariance between $Y_{i+1}^k$ and $Y_{i+1}^l$ is negative due to sampling without replacement [9]. (2) holds since $\Pr(v|p_i) = 1/a_i$, and all $Y_{i+1}^k$'s are identical having the same variance. (3) holds since $t_i = \lceil b \cdot a_i \rceil \geq b \cdot a_i$. (4) holds from the induction hypothesis. Finally, we add the two terms and complete the proof.

$$\mathbb{V}\mathrm{ar}[Y_i \mid p_i]$$
$$= \mathbb{E}[\mathbb{V}\mathrm{ar}[Y_i \mid p_i, v]] + \mathbb{V}\mathrm{ar}[\mathbb{E}[Y_i \mid p_i, v]]$$
$$\leq \frac{1}{1-b} \cdot \left(\frac{1}{b^{n-i}} - 1\right) \cdot c(q, g \mid p_i)^2 + c(q, g \mid p_i)^2$$
$$= \left(\frac{1}{1-b} \cdot \left(\frac{1}{b^{n-i}} - 1\right) + 1\right) \cdot c(q, g \mid p_i)^2$$
$$= \frac{1}{1-b} \cdot \left(\frac{1}{b^{n-i}} - b\right) \cdot c(q, g \mid p_i)^2$$
$$= \frac{b}{1-b} \cdot \left(\frac{1}{b^{n-i+1}} - 1\right) \cdot c(q, g \mid p_i)^2.$$

$\square$

Propositions 1 and 2 explain the accuracy part of Theorem 1. Proposition 3 explains the remaining efficiency part. The proof can be found at the end of Section 5.2.

**Proposition 3.** *Realizing $Y_1$ can be done in $O(AGM(q))$ time.*

**Proof of Theorem 1.** From the propositions, if we repeat Lines 5-8 in Algorithm 2 by $h$ times and take the average of $Y_1$'s as $Z$ (the final estimate), the following statements hold ($h = \frac{\frac{b}{1-b} \cdot (\frac{1}{b^n} - 1)}{\epsilon^2 \cdot (1-\mu)}$).

1) $\mathbb{E}[Z] = \mathbb{E}[Y_1] = |\mathcal{M}|$ since the $Y_1$'s are independent of each other and $\mathbb{E}[Y_1] = |\mathcal{M}|$ from Proposition 1.

2) $\mathbb{V}\mathrm{ar}[Z] = \frac{\mathbb{V}\mathrm{ar}[Y_1]}{h} \leq \frac{1}{h} \cdot \frac{b}{1-b} \cdot (\frac{1}{b^n} - 1) \cdot |\mathcal{M}|^2$ from Proposition 2, and the rightmost term is equal to $\epsilon^2 \cdot (1-\mu) \cdot |\mathcal{M}|^2$ from the definition of $h$.

3) From the Chebyshev inequality, $\Pr(Z - \mathbb{E}[Z] \geq \epsilon \cdot \mathbb{E}[Z]) \leq \frac{\mathbb{V}\mathrm{ar}[Z]}{\epsilon^2 \cdot \mathbb{E}[Z]^2}$ where we know $\frac{\mathbb{V}\mathrm{ar}[Z]}{\epsilon^2 \cdot \mathbb{E}[Z]^2} \leq \frac{\epsilon^2 \cdot (1-\mu) \cdot |\mathcal{M}|^2}{\epsilon^2 \cdot |\mathcal{M}|^2} = 1 - \mu$ from the above statements.

Regarding $b$ and $n$ as constants, $h$ is also a constant. Acquiring $Y_1$ by $h$ times can be done in $O(AGM(q))$ time, and we can obtain $Z$ in $O(AGM(q))$ time that satisfies $\Pr(Z - \mathbb{E}[Z] < \epsilon \cdot \mathbb{E}[Z]) > \mu$. $\square$

## 8 EXPERIMENTS

We now evaluate the performance of Alley to answer the following research questions.

- **Q1.** Compared to existing estimators, how well does Alley perform cardinality estimation on various datasets and queries, in terms of accuracy and efficiency? (Section 8.2)
- **Q2.** How much does the tangled pattern index improve the accuracy of Alley? (Section 8.2)
- **Q3.** Does Alley effectively reduce sampling failures, even for small sampling ratios? (Section 8.3)
- **Q4.** How does our novel mining approach (i.e., walk-fail-then-calculate) improve indexing performance compared to a naive pattern mining approach? (Section 8.4)

### 8.1 Experimental Setup

**Datasets and query sets.** We use six real-world and synthetic datasets and the corresponding query sets shown in Table 2. LUBM [17], YAGO [42], AIDS [1], and Human are used in [38], while HPRD and Youtube are used in [43]. We use the same queries used in [38] and [43] except for Youtube. For Youtube, we use 763 queries with less than 10B embeddings out of 1,800 queries in the original query set, since calculating the true cardinality $|\mathcal{M}|$ requires tremendous computation. For LUBM, we use scale factor 80 by default as in [38], and additionally use larger datasets with scale factors 160, 320, 480, 640, and 800; $|E_g|$ and $|V_g|$ increase linearly to the scale factor. For HPRD and Youtube, query size denotes $|V_q|$ as in [43], while it denotes the number of RDF triples for other datasets as in [38].

**Table 2: Statistics of datasets.**

| Dataset | LUBM | YAGO | AIDS | Human | HPRD | Youtube |
|---|---|---|---|---|---|---|
| # vertices | 2.6M | 12.8M | 254K | 4.7K | 9.5K | 1.1M |
| # edges | 12.3M | 15.8M | 548K | 86K | 70K | 6.0M |
| # vertex labels | 35 | 188K | 50 | 89 | 307 | 25 |
| # edge labels | 35 | 91 | 4 | 0 | 0 | 0 |
| # queries used | 6 | 1,366 | 780 | 49 | 1,800 | 763 |
| Query size | 4 to 6 | 3 to 12 | 3 to 12 | 3 | 4 to 32 | 4 to 32 |
| Cardinality | 15 to 22K | 1 to 28.7M | 1 to 952K | 1 to 9.6K | 1 to 3.2B | 3 to 10B |

**Measure.** We measure accuracy and efficiency using $q$-error [33] and elapsed time, respectively. The $q$-error quantifies the ratio between the actual and the estimated cardinality. The $q$-error is always greater than one, and the smaller the $q$-error, the more accurate the estimation is. Formally, $q$-error $= \max\left(\frac{\max(1,Z)}{\max(1,|\mathcal{M}|)}, \frac{\max(1,|\mathcal{M}|)}{\max(1,Z)}\right)$.

For sampling-based methods, we run 30 times for each query while we run one time for deterministic summary-based methods, as in [38]. For LUBM, we report the mean and the standard deviation of the $q$-error for each query. For the other datasets having numerous queries, we grouped queries by topology, query size, and true cardinality and then, we report the quartiles (i.e., the 25%, 50%, and 75%-tiles) and the whiskers of the $q$-errors for each group.

**Running Environment.** We conducted experiments on a Linux machine with 16 Intel Xeon E5-2450 2.10 GHz CPUs and 32 GB RAM by default, and a machine with 512 GB RAM for scalability test, using a single thread for all experiments. We set one minute as the timeout threshold. For a fair comparison, if a method raises

a timeout at least once in a particular query group, we exclude the results of the method for that group. Following [38], we use sample size $s = N \cdot r$ where $N$ is the number of triples in $g$ that match a triple in $q$, and $r \in (0, 1)$ is the given sampling ratio. We set $r = 0.1\%$ as default and use 1%, 0.1%, and 0.01% for sensitivity analysis. Note that during query optimization a cardinality estimator might be invoked more than a thousand times for a single query. Thus, the estimator should be able to estimate in a few milliseconds. In this respect, we choose the default sampling ratio where sampling-based methods can complete estimation in about one millisecond on most datasets. We report the performance of both using only the sampling technique without any synopsis (denoted by Alley in this section) and with the tangled pattern index (denoted by Alley+TPI). We use $b = 1/32$ as default. When building the tangled pattern index, we use $\zeta = 0.9$, $maxN = 5$ when using edge labels and $maxN = 4$ when using vertex labels. We set the maximum number of label groups in the index to 32.

**Competitors.** We include all seven methods in [38] as our competitors from both graph and relational domains. We additionally considered two more recent competitors, IBJS [27] and subgraph catalog, a summary-based method used in [32]. However, subgraph catalog was originally proposed for graphs with few or zero labels and does not scale well for heterogeneous graphs. It took more than a day to build a summary for almost all datasets we used. Therefore, we exclude subgraph catalog from our experiments. For methods evaluated in [38], C-SET [35] and SumRDF [41] are summary-based methods for graphs, while WanderJoin [28] (aka WJ) and CorrelatedSampling [45] (aka CS) are sampling-based methods for relations. IMPR [11] is a sampling-based method that estimates the cardinality of small graphs with three to five vertices. BSK [10] is a summary-based method for relations that estimates the upper bound of $|\mathcal{M}|$. JSUB in [38] also estimates the upper bound of $|\mathcal{M}|$ using the sampling strategy of [51]. We use the public implementation of G-CARE * for the seven estimators. We additionally implement Alley and IBJS on top of G-CARE. Since IBJS was originally proposed to obtain good-quality samples that aid cost-based query optimizers, it does not choose any sampling order but obtains samples for all subqueries and injects the samples into query optimizers. This is inappropriate for our experiments with fixed sampling ratio; thus, we applied the sampling order of WJ, which is known to be the most accurate and efficient, to IBJS. IBJS works similarly to WJ, but it samples a batch of edges for each walk step instead of sampling an edge as in WJ. For WJ, we additionally implemented two optimization techniques as described in Section 3.6 of [29] on top of the G-CARE implementation.

## 8.2 Overall Performance

**Using Small Queries.** For LUBM and Human, Alley performs a near perfect estimation, achieving almost one $q$-error (Figure 7). This is due to our novel sampling strategy, random walk with intersection that reduces the sample space. However, other sampling-based methods, i.e., WJ, IBJS, IMPR, CS, and JSUB, result in large $q$-error (up to four orders of magnitude higher than Alley's) due to sampling failures even for these small queries. SumRDF is relatively accurate on LUBM but significantly over-estimates on Human,

since SumRDF was originally proposed for RDF graphs, such as LUBM, while Human is a non-RDF graph. These results indicate that summary-based methods do not guarantee consistent performance over various graphs.
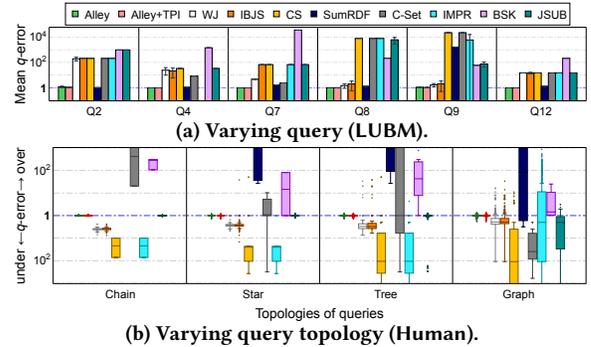


(a) Varying query (LUBM).



(b) Varying query topology (Human).

**Figure 7: Accuracy using small queries.**

**Using Medium-size Queries.** Figure 8 shows the results on YAGO for various query topologies. SumRDF is excluded in this experiment since at least one query of any topology timed out. Overall, more complex queries involving long chains and cycles lead to larger $q$-error. Nevertheless, Alley consistently and significantly outperforms the others for all topologies. For stars, obviously, Alley shows extremely high accuracy since they can be covered by a 1-hop intersection. For trees, Alley still achieves superior accuracy by effectively reducing the sample space. For long chains and cyclic queries (i.e., Cycle and Graph), all sampling-based methods often under-estimate due to the highly selective structures in queries. Alley clearly outperforms all the others for these queries as well.
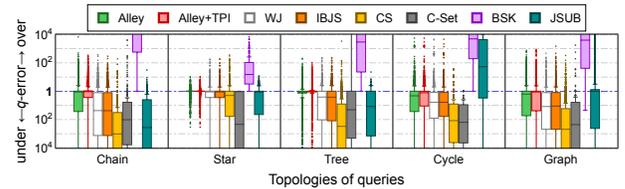


**Figure 8: Varying query topology (YAGO).**

**Using Large Queries.** For larger queries with up to 32 vertices on HPRD, Alley significantly outperforms the others (Figure 9). All competitors suffer from significant under-estimation or time-out (SumRDF and BSK). Specifically, all the sampling-based methods except Alley have similar $q$-error since, for nearly all queries, they fail to sample any embedding and report zero. Alley, however, shows robust performance with $q$-error less than 10 for more than half of the trials. The results show that reducing the sample space is a "must" for large and complex queries. These experiments show the great advantage of interleaving intersections with random walks. For other datasets, we observed similar trends (Figures 10-11).

**Varying True Cardinality.** We now show the trend by varying $|\mathcal{M}|$ on Youtube (Figure 12), where Youtube queries have larger $|\mathcal{M}|$ than HPRD queries. The $q$-error of Alley tends to increase as $|\mathcal{M}|$ increases, but suffers less from severe under-estimation, which
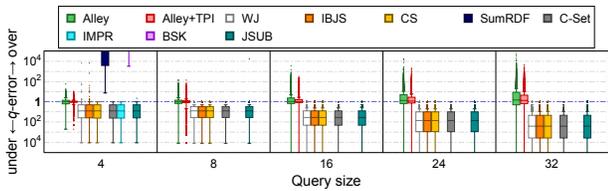
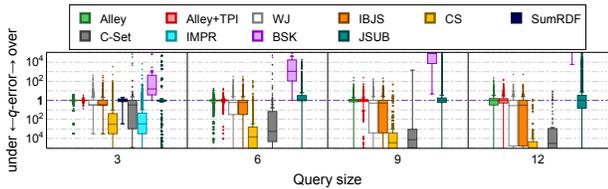Figure 9: Varying query size (HPRD).
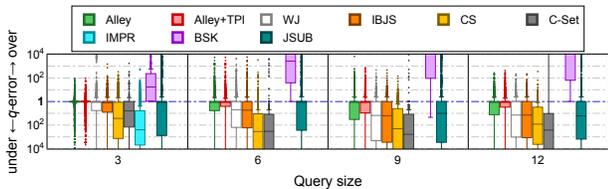


Figure 10: Varying query size (AIDS).



Figure 11: Varying query size (YAGO).

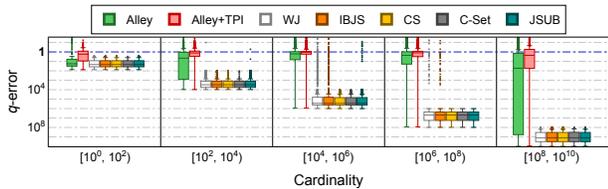occurs in all of the other estimators. This trend is also shown in Figures 13-15.



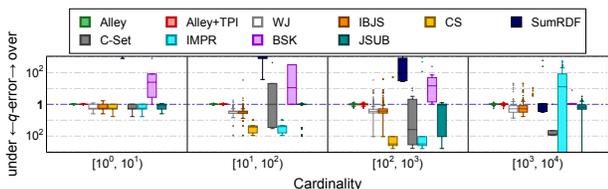Figure 12: Varying true cardinality (Youtube).
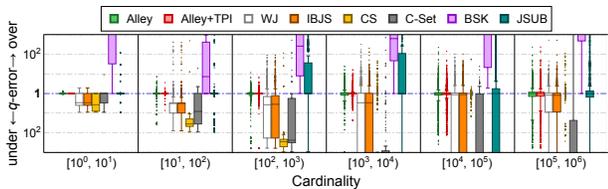


Figure 13: Varying true cardinality (Human).



Figure 14: Varying true cardinality (AIDS).

**Measuring Efficiency.** We measure the elapsed time on AIDS and YAGO (Figure 16). When query graphs are small, ALLEY is slightly
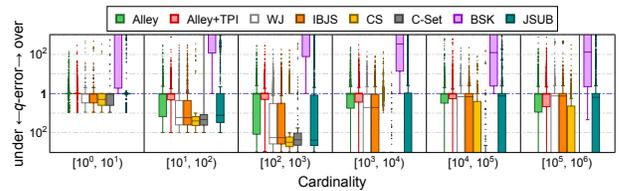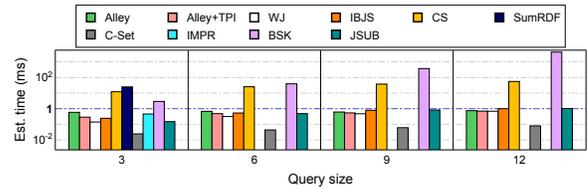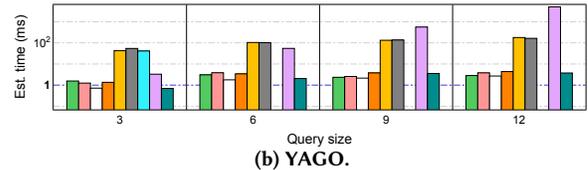


Figure 15: Varying true cardinality (YAGO).

slower than WJ or IBJS, which are the most efficient ones among sampling-based methods. However, as query size grows, ALLEY achieves efficiency similar to theirs. This is due to the early stopping effect of intersections, explained in Section 5.1. Note that, in addition to reasonable accuracy, ALLEY achieves high efficiency with less than a millisecond latency. While C-SET is faster than ALLEY on AIDS, it does not scale well on YAGO since YAGO has many labels. C-SET generates and maintains a large number of entries in its summary, thereby increasing the search time. We observed the same phenomena in the other datasets (Figures 17-18).



(a) AIDS.



(b) YAGO.
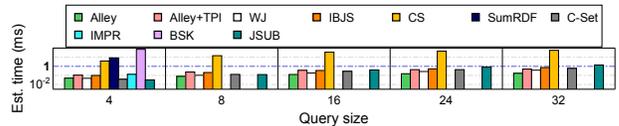
Figure 16: Efficiency by varying query size.



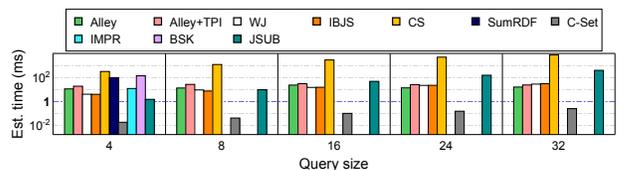Figure 17: Efficiency by varying query size (HPRD).



Figure 18: Efficiency by varying query size (Youtube).

**Combining Sampling with Synopsis.** Compared to ALLEY, AL-LEY+TPI always has better accuracy while having almost the same (sometimes even better) efficiency. This improvement becomes more apparent for large and complex queries, as we claimed earlier. The efficiency of ALLEY+TPI comes from reducing sample space and intersection overhead by using the index, thereby reducing the

estimation time. Due to this advantage along with the superiority of our mining technique that effectively reduces index size, the computational overhead of the index search is almost hidden.

## 8.3 Sampling Failures

In this section, we examine the failures of sampling-based methods. We count extreme failure cases where each method outputs zero due to no sampling success. We compare Alley and Alley+TPI with the two best-performing sampling-based competitors, WJ and IBJS. Table 3 shows the ratio of extreme failure cases to total trials for all queries. Alley and Alley+TPI significantly reduce the failure rates compared to the other methods for all datasets. In particular, when decreasing the sampling ratio on LUBM, the failure rates of the other methods increase greatly by up to 70.6%, while Alley fails only 1.1% of the total trials. This is due to the small sample space of Alley, for instance, the average size of $|P_{o_i}|$ for Alley is 107 on LUBM ($r = 0.01\%$), which is less than 0.05% of WJ's. This results in extremely high accuracy and robustness, as shown in Figure 7a.

**Table 3: The ratio of zero-estimation cases.**

| Dataset | LUBM | | | YAGO | AIDS | HPRD | Youtube |
|---|---|---|---|---|---|---|---|
| **Sampling Ratio ($r$)** | 1% | 0.1% | 0.01% | 0.1% | 0.1% | 0.1% | 0.1% |
| **Alley** | **0.0%** | **0.0%** | **1.1%** | 23.8% | 5.2% | 15.9% | 17.4% |
| **Alley+TPI** | **0.0%** | **0.0%** | **1.1%** | **16.6%** | **1.4%** | **8.2%** | **6.3%** |
| **WJ** | 20.6% | 43.3% | 70.6% | 49.9% | 23.8% | 100% | 99.3% |
| **IBJS** | 21.7% | 41.1% | 67.2% | 49.8% | 23.8% | 100% | 99.3% |

**Varying Sampling Ratio.** In order to further investigate the robustness of Alley, we vary the sampling ratio $r$ from 1% (easy case) to 0.01% (extremely hard case) on AIDS, HPRD, and YAGO (Figures 19-20). We plot notches for better visibility. As the sampling ratio increases, the estimation variance of Alley decreases and Alley becomes more robust, as explained in Section 3 and 7. Again, due to the superiority of our sampling strategy, Alley and Alley+TPI outperform the others by orders of magnitudes for all ratios. Moreover, with $r = 0.01\%$ accuracy, Alley is comparable or better than the other methods with only $r = 1\%$ accuracy.

As shown in Figure 19b, for HPRD, WJ and IBJS fail 100% even with $r = 1\%$, while Alley and Alley+TPI show reasonable accuracy with $r = 1\%$ and 0.1%. However, when $r = 0.01\%$ and the query size is larger than 10, Alley also fails often. Alley+TPI is better than Alley but still has a large variance. Depending on the complexity of the target data and query, one needs to adjust or adaptively determine the sampling ratio in order to perform an accurate estimation. Note that among all methods compared in the experiments, only Alley is able to achieve this functionality in a reasonable time. Others suffer from 1) irrecoverable under-estimation (for sampling-based), 2) irrecoverable and large errors due to loss of information (for summary-based), or even 3) large estimation time.

## 8.4 Mining Performance

We now show the impact of our novel mining approach on the performance of indexing and online estimation. We first evaluate how much our "walk-fail-then-calculate" approach improves the performance of building an index in terms of build time and index size, compared to a naive mechanism that stores domains for all small patterns. We also show the scalability of indexing over varying data sizes. Then, we evaluate whether the tangled pattern index
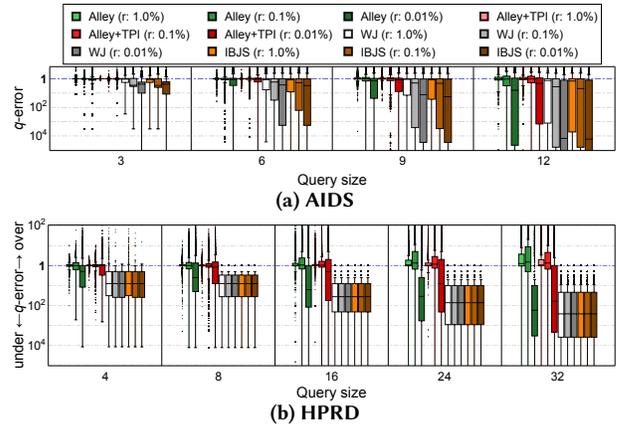


**(a) AIDS**



**(b) HPRD**
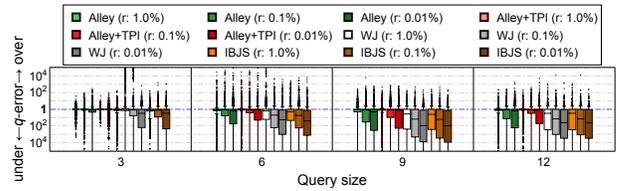
**Figure 19: Varying sampling ratio $r$.**



**Figure 20: Varying sampling ratio (YAGO).**

really retains important information. In other words, we evaluate whether Alley+TPI shows comparable performance with Alley using the naive index (denoted as Alley+Naive). We build the naive index using the same code as Alley+TPI with the only difference being setting $\zeta$ to zero. Note that we use only one thread for building those indices, as mentioned in Section 8.1, although it can easily be parallelized by calculating multiple patterns simultaneously.

The result in Table 4 shows that our "walk-fail-then-calculate" approach effectively reduces both index size and time by up to 56 and 9.7 times, respectively. For $maxN = 5$ in YAGO, even the naive indexing approach fails since materializing the domains of all patterns requires a tremendous amount of memory. While increasing $maxN$ exponentially increases the size of the naive index and time as in typical frequent pattern mining methods, our approach can mitigate the overhead, resulting in a much lighter and more efficient index than using a naive approach.

**Table 4: Ablation results of mining approaches.**

| | Max pattern size: 5 | | Max pattern size: 4 | |
|---|---|---|---|---|
| | Walk-fail-then-calculate | Calculate all ($\zeta = 0$) | Walk-fail-then-calculate | Calculate all ($\zeta = 0$) |
| **Dataset: AIDS** | | | | |
| **Index size in MB (Relative size to input data)** | **88.3** **(6.17)** | 4,950 (346) | **16.2** **(1.13)** | 640 (44.8) |
| **Index time (sec)** | **39.0** | 316 | **4.0** | 38.8 |
| **Dataset: YAGO** | | | | |
| **Index size in MB (Relative size to input data)** | **2,218** **(3.73)** | (Out-of-memory) | **297** **(0.50)** | 3,244 (5.46) |
| **Index time (sec)** | **1,800** | - | **121** | 358 |

Figure 21 shows the indexing performance over large LUBM datasets with scale factors from 160 to 800. The indexing time increases linearly to the scale, and CalculateDomains dominates other functions, occupying 95% of the total time. GetFailureRate

accounts for only 3%. The numbers of calls for CALCULATEDOMAINS and GETFAILURERATE are similar for all data sizes, which are 9K and 150K, respectively; each call to GETFAILURERATE is about 500 times faster than CALCULATEDOMAINS. Therefore, our "walk-fail-then-calculate" approach does not incur a significant overhead yet effectively prunes the patterns to index.
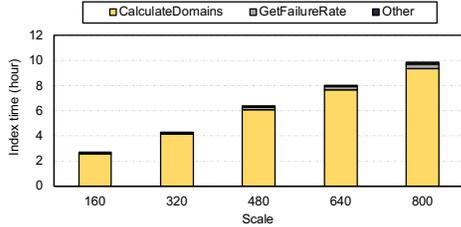


Figure 21: Indexing time by varying data size.

Figure 22 shows that the tangled pattern index is as effective as the naive index. Both ALLEY+TPI and ALLEY+NAIVE show higher accuracy than ALLEY without using an index. These experimental results illustrate that 1) combining with synopses can increase the performance of sampling in cardinality estimation and 2) filtering out domain calculation by random walks during mining can increase scalability while preserving effectiveness.
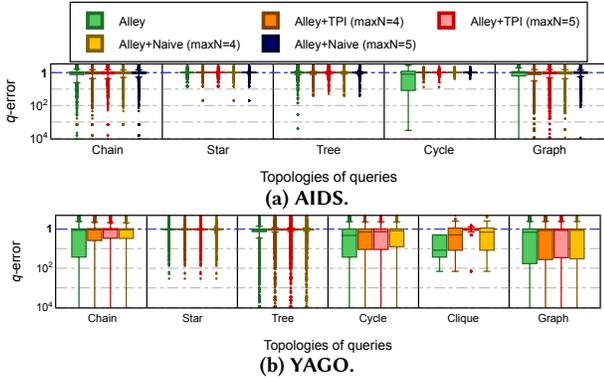


(a) AIDS.



(b) YAGO.

Figure 22: Accuracy by varying indexing approach.

We additionally investigate the robustness of ALLEY+TPI regarding data updates. Edge insertion can introduce additional embeddings, which causes false negatives of the pre-built index. Therefore, ALLEY+TPI would increasingly underestimates as the number of inserted edges increases. In order to simulate evolving graphs, we build an index by using 90% of data edges, and run cardinality estimation on 90%, 92%, 94%, 96%, 98%, and 100% of data edges. We tested with AIDS and YAGO. Figure 23 shows the accuracy of each amount of data edges inserted after the index build. The results show that the estimation accuracy does not deteriorate much even if 5-10% of data edges are inserted. From the result, we can see that we do not need to rebuild a tangled pattern index until there are significant updates to the data graph.

## 9 RELATED WORK

**Graph Pattern Cardinality Estimation.** Table 5 shows a summary of graph pattern cardinality estimators with four comparison aspects. The fourth aspect is based on our empirical study.
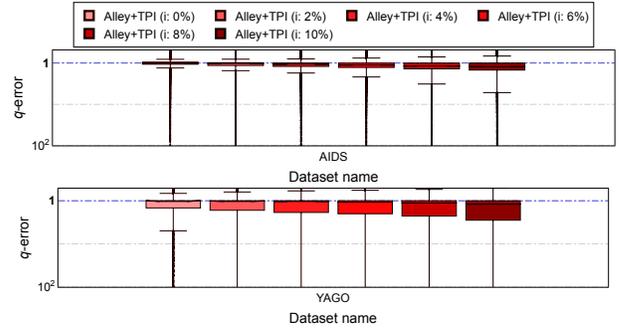


Figure 23: Accuracy of ALLEY+TPI by varying amount of inserted edges.

Table 5: Comparison aspects of estimators. G, R, S, and A refer to Graph, Relation, Synopsis, and sAmpling. ✔, ✘, and ━ represent positive, negative, and omission in our experiments, respectively.

| Estinator | Domain | Method | Unbiased? | Accurate? |
|---|---|---|---|---|
| C-SET [35] | G | S | ✘ | ✘ |
| SUMRDF [41] | G | S | ✘ | ✘ |
| CATALOG [32] | G | S | ✘ | ━ |
| IMPR [11] | G | A | ✔ | ✘ |
| SSTE [7] | G | A | ✔ | ━ |
| WJ [28] | R | A | ✔ | ✔ |
| IBJS [27] | R | A | ✔ | ✔ |
| CS [45] | R | A | ✔ | ✘ |
| JSUB [38, 51] | R | A | ✘ | ✘ |
| BSK [10] | R | S | ✘ | ✘ |
| ALLEY (ours) | G | S, A | ✔ | ✔✔ |

Summary-based approaches have been widely adapted for graph pattern cardinality estimation [4, 30, 35, 41]. Early work [4, 30] stores the exact cardinality of small acyclic patterns (paths or trees) to estimate larger subgraphs in a data graph. More recent work includes C-SET [35], SUMRDF [41], and subgraph catalog [32]. C-SET [35] decomposes $q$ into independent subqueries. SUMRDF [41] merges the vertices and edges of $g$ into a smaller graph $S$, and extends the embeddings calculated in $S$ to the embeddings in $g$ under a uniformity assumption. However, these ad-hoc assumptions have been avoided in sampling-based work as they can significantly degrade estimation accuracy [38]. Subgraph catalog [32] precomputes the average expansion ratio from each $k$-vertex pattern to a $(k+1)$-vertex pattern and multiplies these ratios under the uniformity assumption. However, as mentioned in Section 8.1, storing all small patterns can lead to scalability issues, especially for heterogeneous graphs.

Join cardinality estimation in relational databases, which is closely related to graph pattern cardinality estimation, often adopts sampling approaches [27, 28]. However, we have shown that these sampling-based methods have the serious under-estimation problem due to sampling failures in large sample spaces. While sampling failures have been addressed in [28], it briefly mentions that too many failed random walks will slow down the convergence of estimation, and failures occur more for cyclic queries. However,

resolving these is not the main contribution of [28]. Instead, [28] focuses on selecting a good walk order using a round-robin approach and incorporating trial and failed walks in estimation. In contrast, we propose random walk with intersection, which can reduce sampling failures by cutting down the sample space. We further reduce sampling failures by combining sampling and synopses.

Theoretical bounds on estimation variance or runtime have been established by [5, 7, 12]. For the most part, these works solve very specific problems, e.g., estimating the number of triangles [12], stars [5], or cliques [13]. The most recent work, SSTE [7], generalizes the previous theoretical estimators to arbitrary-shaped queries. However, SSTE lacks practicality since it does not perform random walks but independently samples edges as $T_1$ and $T_2$ in Section 3.
**Frequent Pattern Mining.** Frequent pattern mining finds all patterns (up to a limited size) that appear frequently in graph data. Here, data can be a set of multiple small graphs [46] or a single large graph [2, 14]. Then, a domain of pattern $p$ represents either the set of small graphs that contain $p$ or the set of vertices that participate in an embedding of $p$. Other metrics that satisfy the anti-monotone property are also used [16, 26]. The anti-monotone property permits the mining algorithms to search domains of larger patterns filtered by the domains of smaller sub-patterns. Despite these optimizations, scalability has been a serious concern [2, 22]. Recent methods try to solve this problem by using distributed computing.

Several methods use sampling for frequent pattern mining [2, 52]. These methods employ a sampling-based estimator to detect a set of patterns that are frequent with high probability and prune out infrequent patterns [2], or they sample the data graph to approximate the frequencies of patterns [52]. In contrast, we propose a novel indexing algorithm to store the domains of tangled patterns level-by-level by performing sampling on the previously computed domains of tangled sub-patterns.

## 10 DISCUSSION & CONCLUSION

We presented ALLEY, an accurate and efficient graph pattern cardinality estimator. ALLEY is a hybrid method that combines sampling and synopses, which includes 1) a new sampling strategy based on random walk with intersection and branching, and 2) an efficient mining algorithm and index for tangled patterns. Combining these two novel ideas, ALLEY achieves high accuracy within a reasonable latency. ALLEY also guarantees worst-case optimal time complexity for any given error bound and confidence level.

ALLEY can be used in the following scenarios. Given a pattern matching query $q$, a cost-based optimizer enumerates valid subgraph patterns of $q$ using dynamic programming and computes the cost for each pattern. The cost model calculates the cost based on the estimated cardinality. Here, ALLEY computes such estimated cardinality. Furthermore, from the anti-monotone property, domains of a subquery $q'$ of $q$ can be used as the matching candidates for $q$. Since the tangled patterns have strong structural (e.g., cycles) or label correlations, the search space can be effectively reduced.

While ALLEY can be integrated into any in-memory graph database system in its query optimizer or for answering approximate count queries, it would be beneficial if the system has data structures optimized for intersection, which is the main bottleneck in ALLEY. For example, in order to efficiently process intersections,

the edges in a data graph should be stored in sorted adjacency lists. A recent graph database system, EMPTYHEADED [3], can further benefit from its hybrid representation using bitmaps and SIMD operations, enabling faster intersections. Recent databases with GPU support might further benefit from GPU-based intersection [48].

For exploiting parallelism, index construction can be easily parallelized in a multi-core, distributed setting as in previous pattern mining studies. Our sampling-based estimation can also be easily parallelized in a multi-core setting by calling Line 6 of Algorithm 4 in parallel. However, it is non-trivial to implement an efficient online estimator in a distributed setting due to the communication overhead. We leave it as future work to extend the sampling-based estimation to a distributed setting.

## REFERENCES

[1] 2014. DTP, AIDS antiviral screen. http://dtp.nci.nih.gov/docs/aids/aids_data.html.

[2] Ehab Abdelhamid, Ibrahim Abdelaziz, Panos Kalnis, Zuhair Khayyat, and Fuad Jamour. 2016. Scalemine: Scalable parallel frequent subgraph mining in a single large graph. In *SC'16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*. IEEE, 716–727.

[3] Christopher R Aberger, Susan Tu, Kunle Olukotun, and Christopher Ré. 2016. Emptyheaded: A relational engine for graph processing. In *Proceedings of the 2016 International Conference on Management of Data*. ACM, 431–446.

[4] Ashraf Aboulnaga, Alaa R Alameldeen, and Jeffrey F Naughton. 2001. Estimating the selectivity of XML path expressions for internet scale applications. In *VLDB*, Vol. 1. Citeseer, 591–600.

[5] Maryam Aliakbarpour, Amartya Shankha Biswas, Themis Gouleakis, John Peebles, Ronitt Rubinfeld, and Anak Yodpinyanee. 2018. Sublinear-time algorithms for counting star subgraphs via edge sampling. *Algorithmica* 80, 2 (2018), 668–697.

[6] Molham Aref, Balder ten Cate, Todd J Green, Benny Kimelfeld, Dan Olteanu, Emir Pasalic, Todd L Veldhuizen, and Geoffrey Washburn. 2015. Design and implementation of the LogicBlox system. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. ACM, 1371–1382.

[7] Sepehr Assadi, Michael Kapralov, and Sanjeev Khanna. 2019. A Simple Sublinear-Time Algorithm for Counting Arbitrary Subgraphs via Edge Sampling. In *10th Innovations in Theoretical Computer Science Conference, ITCS 2019, January 10-12, 2019, San Diego, California, USA (LIPIcs, Vol. 124)*, Avrim Blum (Ed.). Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 6:1–6:20. https://doi.org/10.4230/LIPIcs.ITCS.2019.6

[8] Albert Atserias, Martin Grohe, and Dániel Marx. 2008. Size bounds and query plans for relational joins. In *Foundations of Computer Science, 2008. FOCS'08. IEEE 49th Annual IEEE Symposium on*. IEEE, 739–748.

[9] Moulinath Banerjee. 2012. Simple Random Sampling. *Unpublished Manuscript, University of Michigan, Michigan* (2012).

[10] Walter Cai, Magdalena Balazinska, and Dan Suciu. 2019. Pessimistic Cardinality Estimation: Tighter Upper Bounds for Intermediate Join Cardinalities. In *Proceedings of the 2019 International Conference on Management of Data*. ACM, 18–35.

[11] Xiaowei Chen and John C. S. Lui. 2016. Mining Graphlet Counts in Online Social Networks. In *IEEE 16th International Conference on Data Mining, ICDM 2016, December 12-15, 2016, Barcelona, Spain*, Francesco Bonchi, Josep Domingo-Ferrer, Ricardo Baeza-Yates, Zhi-Hua Zhou, and Xindong Wu (Eds.). IEEE Computer Society, 71–80. https://doi.org/10.1109/ICDM.2016.0018

[12] Talya Eden, Amit Levi, Dana Ron, and C Seshadhri. 2017. Approximately counting triangles in sublinear time. *SIAM J. Comput.* 46, 5 (2017), 1603–1646.

[13] Talya Eden, Dana Ron, and C Seshadhri. 2018. On approximating the number of k-cliques in sublinear time. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*. ACM, 722–734.

[14] Mohammed Elseidy, Ehab Abdelhamid, Spiros Skiadopoulos, and Panos Kalnis. 2014. Grami: Frequent subgraph and pattern mining in a single large graph. *Proceedings of the VLDB Endowment* 7, 7 (2014), 517–528.

[15] Wenfei Fan. 2012. Graph pattern matching revised for social network analysis. In *Proceedings of the 15th International Conference on Database Theory*. ACM, 8–21.

[16] Mathias Fiedler and Christian Borgelt. 2007. Subgraph Support in a Single Large Graph. In *Workshops Proceedings of the 7th IEEE International Conference on Data Mining (ICDM 2007), October 28-31, 2007, Omaha, Nebraska, USA*. IEEE Computer Society, 399–404. https://doi.org/10.1109/ICDMW.2007.74

[17] Yuanbo Guo, Zhengxiang Pan, and Jeff Heflin. 2005. LUBM: A benchmark for OWL knowledge base systems. *Web Semantics: Science, Services and Agents on the World Wide Web* 3, 2 (2005), 158–182.

[18] Myoungji Han, Hyunjoon Kim, Geonmo Gu, Kunsoo Park, and Wook-Shin Han. 2019. Efficient Subgraph Matching: Harmonizing Dynamic Programming, Adaptive Matching Order, and Failing Set Together. In *Proceedings of the 2019 International Conference on Management of Data*. ACM, 1429–1446.

[19] Shuo Han, Lei Zou, and Jeffrey Xu Yu. 2018. Speeding Up Set Intersections in Graph Algorithms using SIMD Instructions. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, Gautam Das, Christopher M. Jermaine, and Philip A. Bernstein (Eds.). ACM, 1587–1602. https://doi.org/10.1145/3183713.3196924

[20] Wook-Shin Han, Jinsoo Lee, and Jeong-Hoon Lee. 2013. Turbo iso: towards ultrafast and robust subgraph isomorphism search in large graph databases. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*. ACM, 337–348.

[21] Sen Hu, Lei Zou, Jeffrey Xu Yu, Haixun Wang, and Dongyan Zhao. 2018. Answering natural language questions by subgraph matching over knowledge graphs. *IEEE Transactions on Knowledge and Data Engineering* 30, 5 (2018), 824–837.

[22] Chuntao Jiang, Frans Coenen, and Michele Zito. 2013. A survey of frequent subgraph mining algorithms. *Knowledge Engineering Review* 28, 1 (2013), 75–105.

[23] Chathura Kankanamge, Siddhartha Sahu, Amine Mhedhbi, Jeremy Chen, and Semih Salihoglu. 2017. Graphflow: An Active Graph Database. In *Proceedings of the 2017 ACM International Conference on Management of Data*. ACM, 1695–1698.

[24] Hyeonji Kim, Juneyoung Lee, Sourav S Bhowmick, Wook-Shin Han, JeongHoon Lee, Seongyun Ko, and Moath HA Jarrah. 2016. DUALSIM: Parallel subgraph enumeration in a massive graph on a single machine. In *Proceedings of the 2016 International Conference on Management of Data*. ACM, 1231–1245.

[25] Jinha Kim, Hyungyu Shin, Wook-Shin Han, Sungpack Hong, and Hassan Chafi. 2015. Taming subgraph isomorphism for RDF query processing. *Proceedings of the VLDB Endowment* 8, 11 (2015), 1238–1249.

[26] Michihiro Kuramochi and George Karypis. 2005. Finding Frequent Patterns in a Large Sparse Graph\*. *Data Min. Knowl. Discov.* 11, 3 (2005), 243–271. https://doi.org/10.1007/s10618-005-0003-9

[27] Viktor Leis, Bernhard Radke, Andrey Gubichev, Alfons Kemper, and Thomas Neumann. 2017. Cardinality Estimation Done Right: Index-Based Join Sampling.. In *Cidr*.

[28] Feifei Li, Bin Wu, Ke Yi, and Zhuoyue Zhao. 2016. Wander join: Online aggregation via random walks. In *Proceedings of the 2016 International Conference on Management of Data*. ACM, 615–629.

[29] Feifei Li, Bin Wu, Ke Yi, and Zhuoyue Zhao. 2017. Wander join and XDB: online aggregation via random walks. *ACM SIGMOD Record* 46, 1 (2017), 33–40.

[30] Angela Maduko, Kemafor Anyanwu, Amit Sheth, and Paul Schliekelman. 2008. Graph summaries for subgraph frequency estimation. In *European Semantic Web Conference*. Springer, 508–523.

[31] Volker Markl, Nimrod Megiddo, Marcel Kutsch, Tam Minh Tran, P Haas, and Utkarsh Srivastava. 2005. Consistently estimating the selectivity of conjuncts of predicates. In *Proceedings of the 31st international conference on Very large data bases*. 373–384.

[32] Amine Mhedhbi and Semih Salihoglu. 2019. Optimizing Subgraph Queries by Combining Binary and Worst-Case Optimal Joins. *Proc. VLDB Endow.* 12, 11 (2019), 1692–1704. https://doi.org/10.14778/3342263.3342643

[33] Guido Moerkotte, Thomas Neumann, and Gabriele Steidl. 2009. Preventing bad plans by bounding the impact of cardinality estimation errors. *Proceedings of the VLDB Endowment* 2, 1 (2009), 982–993.

[34] Rajeev Motwani and Prabhakar Raghavan. 1995. *Randomized algorithms*. Cambridge university press.

[35] Thomas Neumann and Guido Moerkotte. 2011. Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins. In *2011 IEEE 27th International Conference on Data Engineering*. IEEE, 984–994.

[36] Hung Q Ngo, Ely Porat, Christopher Ré, and Atri Rudra. 2012. Worst-case optimal join algorithms. In *Proceedings of the 31st ACM SIGMOD-SIGACT-SIGAI symposium on Principles of Database Systems*. ACM, 37–48.

[37] Hung Q Ngo, Christopher Ré, and Atri Rudra. 2014. Skew strikes back: New developments in the theory of join algorithms. *ACM SIGMOD Record* 42, 4 (2014), 5–16.

[38] Yeonsu Park, Seongyun Ko, Sourav S Bhowmick, Kyoungmin Kim, Kijae Hong, and Wook-Shin Han. 2020. G-CARE: A Framework for Performance Benchmarking of Cardinality Estimation Techniques for Subgraph Matching. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 1099–1114.

[39] Siddhartha Sahu, Amine Mhedhbi, Semih Salihoglu, Jimmy Lin, and M Tamer Özsu. 2017. The ubiquity of large graphs and surprising challenges of graph processing. *Proceedings of the VLDB Endowment* 11, 4 (2017), 420–431.

[40] Monroe Sirken and Iris Shimizu. 1999. Population based establishment sample surveys: The Horvitz-Thompson estimator. *Survey Methodology* 25, 2 (1999), 187–192.

[41] Giorgio Stefanoni, Boris Motik, and Egor V Kostylev. 2018. Estimating the cardinality of conjunctive queries over RDF data using graph summarisation. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 1043–1052.

[42] Fabian M Suchanek, Gjergji Kasneci, and Gerhard Weikum. 2008. Yago: A large ontology from wikipedia and wordnet. *Web Semantics: Science, Services and Agents on the World Wide Web* 6, 3 (2008), 203–217.

[43] Shixuan Sun and Qiong Luo. 2020. In-Memory Subgraph Matching: An Indepth Study. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 1083–1098.

[44] Carlos HC Teixeira, Alexandre J Fonseca, Marco Serafini, Georgos Siganos, Mohammed J Zaki, and Ashraf Aboulnaga. 2015. Arabesque: a system for distributed graph mining. In *Proceedings of the 25th Symposium on Operating Systems Principles*. 425–440.

[45] David Vengerov, Andre Cavalheiro Menck, Mohamed Zait, and Sunil P Chakkappen. 2015. Join size estimation subject to filter conditions. *Proceedings of the VLDB Endowment* 8, 12 (2015), 1530–1541.

[46] Xifeng Yan, Philip S. Yu, and Jiawei Han. 2004. Graph Indexing: A Frequent Structure-based Approach. In *Proceedings of the ACM SIGMOD International Conference on Management of Data, Paris, France, June 13-18, 2004*, Gerhard Weikum, Arnd Christian König, and Stefan Deßloch (Eds.). ACM, 335–346. https://doi.org/10.1145/1007568.1007607

[47] Zhengwei Yang, Ada Wai-Chee Fu, and Ruifeng Liu. 2016. Diversified Top-k Subgraph Querying in a Large Graph. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD Conference 2016, San Francisco, CA, USA, June 26 - July 01, 2016*, Fatma Özcan, Georgia Koutrika, and Sam Madden (Eds.). ACM, 1167–1182. https://doi.org/10.1145/2882903.2915216

[48] Li Zeng, Lei Zou, M. Tamer Özsu, Lin Hu, and Fan Zhang. 2020. GSI: GPU-friendly Subgraph Isomorphism. In *36th IEEE International Conference on Data Engineering, ICDE 2020, Dallas, TX, USA, April 20-24, 2020*. IEEE, 1249–1260. https://doi.org/10.1109/ICDE48307.2020.00112

[49] Shijie Zhang, Shirong Li, and Jiong Yang. 2009. GADDI: distance index based subgraph matching in biological networks. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*. ACM, 192–203.

[50] Peixiang Zhao and Jiawei Han. 2010. On graph query optimization in large networks. *Proceedings of the VLDB Endowment* 3, 1-2 (2010), 340–351.

[51] Zhuoyue Zhao, Robert Christensen, Feifei Li, Xiao Hu, and Ke Yi. 2018. Random Sampling over Joins Revisited. In *Proceedings of the 2018 International Conference on Management of Data*. ACM, 1525–1539.

[52] Ruoyu Zou and Lawrence B Holder. 2010. Frequent subgraph mining on a single large graph using sampling techniques. In *Proceedings of the eighth workshop on mining and learning with graphs*. 171–178.