# MSQ-Index: A Succinct Index for Fast Graph Similarity Search

Xiaoyang Chen
Xidian University
Xi'an 710071, China
chenxyu1991@gmail.com

Hongwei Huo
Xidian University
Xi'an 710071, China
hwhuo@mail.xidian.edu.cn

Jun Huan
The University of Kansas
Lawrence, KS 66045, USA
jhuan@ku.edu

Jeffrey Scott Vitter
The University of Mississippi
Oxford, MS 38677-1848, USA
JSV@OleMiss.edu

## ABSTRACT

Graph similarity search has received considerable attention in many applications, such as bioinformatics, data mining, pattern recognition, and social networks. Existing methods for this problem have limited scalability because of the huge amount of memory they consume when handling very large graph databases with millions or billions of graphs.

In this paper, we study the problem of graph similarity search under the graph edit distance constraint. We present a space-efficient index structure based upon the *q-gram* tree that incorporates succinct data structures and hybrid encoding to achieve improved query time performance with minimal space usage. Specifically, the space usage of our index requires only 5%–15% of the previous state-of-the-art indexing size on the tested data while at the same time achieving 2–3 times acceleration in query time with small data sets. We also boost the query performance by augmenting the global filter with range search, which allows us to perform a query in a reduced region. In addition, we propose two effective filters that combine degree structures and label structures. Extensive experiments demonstrate that our proposed approach is superior in space and competitive in filtering to the state-of-the-art approaches. To the best of our knowledge, our index is the first in-memory index for this problem that successfully scales to cope with the large dataset of 25 million chemical structure graphs from the PubChem dataset.

## 1. INTRODUCTION

Graphs are widely used to model complicated data objects in many disciplines, such as bioinformatics, social networks, software and data engineering. Effective analysis and management of graph data become increasingly important. Many queries have been investigated and they can be roughly divided into two broad categories: graph exact search [19] and graph similarity search [6]. Compared with exact search, similarity search can provide a robust solution that permits error-tolerant and supports to search patterns that are not precisely defined.

Similarity computation between two attributed graphs is a core operation of graph similarity search and it has been used in various applications such as pattern recognition, graph classification and chemistry analysis [15]. There are at least four metrics being well investigated: graph edit distance [9, 16, 22, 23, 26], maximal common subgraph distance [2], graph alignment [3] and graph kernel functions [11, 18]. In this paper, we focus on the graph edit distance since it is applicable to virtually all types of data graphs and can also capture precisely structural differences. The graph edit distance $ged(g, h)$ between two graphs $g$ and $h$ is defined as the minimum number of edit operations needed to transform one graph to another.

Given a graph database $G$, a query graph $h$ and an edit distance threshold $\tau$, the graph similarity search problem aims to find all graphs $g$ in $G$ satisfying $ged(g, h) \leq \tau$. Unfortunately, computing the graph edit distance is known to be an NP-hard problem [22]. Therefore, for a large transaction database, such as PubChem, which stores information about roughly 50 million chemical compounds, similarity search is very challenging.

Most of the existing methods adopt the filter-and-verify schema to speed up the search. With such a schema, we first filter data graphs that are not possible results to generate a candidate set, and then validate the candidate graphs with the expensive graph edit distance computations. In general, the existing filters can be divided into four categories: global filter, *q-gram* counting filter, mapping distance-based filter and disjoint partition-based filter. Specifically, number count filter [22] and label count filter [24] are two global filters. The former is derived based upon the differences of the number of vertices and edges of comparing graphs. The later takes labels as well as structures into account, further improving the former. $\kappa$-AT [16] and GSimJoin [24] are two major *q-gram* counting filters. They considered a $\kappa$-*adjacent* subtree and a simple path of length $p$ as a *q-gram*, respectively. C-Star [22] and Mixed [25, 26] are two major mapping distance-based filters. The lower bounds are derived based on the minimum weighted bipartite graphs

between the star and branch structures of comparing graphs, respectively. Pars [23] is a disjoint partition-based filter. It divides each data graph $g$ into several disjoint substructures and prunes $g$ by the subgraph isomorphism.

Even though promising preliminary results have been achieved by existing methods GSimJoin [24], C-Star [22] and Mixed [26], our empirical evaluation of all the methods aforementioned showed that they are not scalable to large graph databases. The critical limitation of existing methods are: (1) existing filters having a weak filter ability produce large candidate sets, resulting in an unacceptable computational cost for verification, (2) the index storage cost of the existing methods is too expensive to run properly. For example, for a database of 10 million graphs, C-Star on average produces $5 \times 10^5$ number of candidates for verification when $\tau = 5$. Both GSimJoin and Mixed produce an index that is too large to fit into the main memory for large input data. The details of the empirical study are presented in Section 7.

To solve the above issues, we propose a space-efficient index structure for graph similarity search which significantly reduces the storage space. Our contributions in this paper are summarized below.

- We propose two effective filters, i.e. degree-based *q-gram* counting filter and degree-sequence filter, by using the degree structures and label structures.

- We create a *q-gram* tree to speed up filtering process. More importantly, we propose the succinct representation of the *q-gram* tree which combines with hybrid coding, significantly reducing the space required for the representation of the *q-gram* tree.

- We convert the number count filter to a two-dimensional orthogonal range searching, which helps us perform a query at a reduced region and hence further improves the filtering performance.

- We have conducted extensive experiments over both real and synthetic datasets to evaluate the index storage space, construction time, filtering capability, and response time. The result is graph similarity search index that we refer to as "MSQ-Index". It confirms the effectiveness and efficiency of our proposed approaches and show that our method can scale well to cope with the large dataset of 25 million chemical compounds from the PubChem dataset.

The rest of this paper is organized as follows: In Section 2, we introduce the problem definition. In Section 3, we present the degree-based *q-gram* counting filter and the degree-sequence filter. In Section 4, we give a method to reduce the query region. In Section 5, we introduce the index structure. In Section 6, we give the query algorithm. In Section 7, we report the experimental results. We investigate the research work related to this paper in Section 8. Finally, we make concluding remarks in Section 9.

## 2. PRELIMINARIES

In this section, we introduce the basic notations and definitions of graph edit distance and graph similarity search.

DEFINITION 1 (ATTRIBUTED GRAPH). *A labeled graph is defined as a six-tuple $g = (V_g, E_g, \mu, \zeta, \Sigma_{V_g}, \Sigma_{E_g})$, where $V_g$ is the set of vertices, $E_g \subseteq V_g \times V_g$ is the set of edges, $\mu : V_g \to \Sigma_{V_g}$ is the vertex labeling function which assigns a label $\mu(v)$ to the vertex $v$, $\zeta : E_g \to \Sigma_{E_g}$ is the edge labeling function which assigns a label $\zeta(e)$ to the edge $e$, $\Sigma_{V_g}$ and $\Sigma_{E_g}$ are the label multisets of $V_g$ and $E_g$, respectively.*

In this paper, we only focus on simple undirected graphs without multi-edge or self-loop. We use $|V_g|$ and $|E_g|$ to denote the number of vertices and edges in $g$, respectively. The graph size refers to $|V_g|$ in this paper. Although in the following discussion we only focus on undirected graphs, our methods can be extended to handle directed graphs.

DEFINITION 2 (GRAPH ISOMORPHISM [19]). *Given two graphs $g$ and $h$, an isomorphism of graphs $g$ and $h$ is a bijection $f : V_g \to V_h$, such that (1) for all $v \in V_g$, $f(v) \in V_h$ and $\mu(v) = \mu(f(v))$. (2) for all $e(u, v) \in E_g$, $e(f(u), f(v)) \in E_h$ and $\zeta(e(u, v)) = \zeta(e(f(u), f(v)))$. If $g$ is isomorphic to $h$, we denote $g \cong h$.*

There are six primitive edit operations that can transform one graph to another [1]. These edit operations are inserting/deleting an isolated vertex, inserting/deleting an edge between two vertices and substituting the label of a vertex or an edge. We denote the substitution of two vertices $u$ and $v$ by $(u \to v)$, the deletion of vertex $u$ by $(u \to \epsilon)$, and the insertion of vertex $v$ by $(\epsilon \to v)$. For edges, we use a similar notation. Given two graphs $g$ and $h$, an edit path $P = \langle p_1, p_2, \ldots, p_k \rangle$ is a sequence of edit operations that transforms $h$ to $g$, such as $h = h^0 \xrightarrow{p_1} h^1 \xrightarrow{p_2} \ldots \xrightarrow{p_k} h^k \cong g$. In Figure 1, we give an example of an edit path $P$ between $g$ and $h$, where the vertex labels are represented by different symbols. The length of $P$ is 6, which consists of two edge deletions, one vertex deletion, one vertex insertion and two edge insertions. In the following sections, we use $|P|$ to denote the length of $P$.

DEFINITION 3 (OPTIMAL EDIT PATH). *Given two graphs $g$ and $h$, an edit path $P$ between $g$ and $h$ is an optimal edit path if and only if there does not exist another edit path $P'$ such that $|P'| < |P|$. The graph edit distance between them, denoted by $ged(g, h)$, is the length of the optimal edit path.*

**Problem statement**: Given a graph database $G = \{g_1, g_2, \ldots, g_{|G|}\}$, a query graph $h$, and an edit distance threshold $\tau$, the problem is to find all the graphs $g$ in $G$ such that $ged(g, h) \leq \tau$, where $ged(g, h)$ is the graph edit distance of graphs $g$ and $h$ defined in Definition 3.

Figure 2 shows a query graph $h$ and three data graphs $g_1, g_2$, and $g_3$. We can obtain that $ged(g_1, h) = 3$, $ged(g_2, h) = 4$, and $ged(g_3, h) = 3$. If the edit distance threshold $\tau = 3$, $g_1$ and $g_3$ are the required graphs.
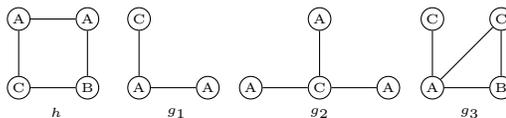


Figure 2: Query graph $h$ and data graphs $g_1$, $g_2$, and $g_3$.

The computation of graph edit distance is an NP-hard problem [22]. The state-of-the-art approaches like [16, 22,
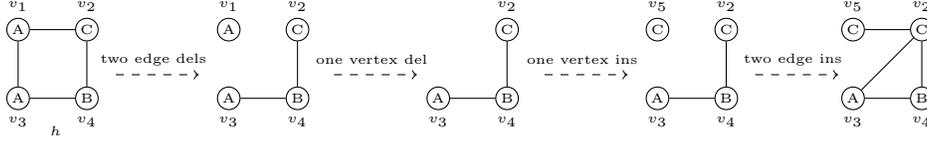
Figure 1: An edit path $P$ between graphs $g$ and $h$.

23, 24, 25] for graph similarity search use a filter-and-verify schema to speed up query process. In the filtering phase, it computes the candidate set $Cand = \{g : \xi(g, h) \leq \tau \text{ and } g \in G\}$, where $\xi(g, h)$ is the lower bound on $ged(g, h)$. In the verification phase, for each graph $g$ in $Cand$, it needs to compute $ged(g, h)$. Obviously, it is good for the size $|Cand|$ as small as possible.

In this paper, we propose two filters, i.e., degree-based $q$-gram counting filter and degree-sequence filter using the degree structures and label structures in a graph. Besides, we also use the following two simple but effective global filters, i.e., number count filter [22] and label count filter [24]. Number count filter is derived based upon the differences of the number of vertices and edges of comparing graphs and given by $dist_N(g, h) = ||V_g| - |V_h|| + ||E_g| - |E_h||$. Label count filter improves the number count filter by taking labels as well as structures into account and is given by $dist_L(g, h) = \max\{|V_g|, |V_h|\} - |\Sigma_{V_g} \cap \Sigma_{V_h}| + \max\{|E_g|, |E_h|\} - |\Sigma_{E_g} \cap \Sigma_{E_h}|$. By using all of them, we can obtain a candidate set as small as possible.

## 3. MULTIPLE FILTERS

### 3.1 Optimal Edit Path

Given two graphs $g$ and $h$, and an optimal edit path $P$ between them, we group the operations on $P$ into five sets of edit operations: vertex deletion group $P_{VD} = \{p_i : p_i = (u \rightarrow \epsilon) \in P\}$, vertex insertion group $P_{VI} = \{p_i : p_i = (\epsilon \rightarrow v) \in P\}$, vertex substitution group $P_{VS} = \{p_i : p_i = (u \rightarrow v) \in P\}$, edge deletion group $P_{ED} = \{p_i : p_i = (e(u,v) \rightarrow \epsilon) \text{ and } ((u \rightarrow \epsilon) \in P_{VD} \text{ or } (v \rightarrow \epsilon) \in P_{VD})\}$ consists of the edge deletions performed on the deleted vertices, and edge operation group $P_O$ consists of the edit operations performed on edges except for those in $P_{ED}$.

For an optimal edit path $P$ between $g$ and $h$ , the insertion/deletion/substitution edit operation on a vertex $v$ or an edge $e$ must happen only once, thus the edit operations in $P_{VD}$ are independent of each other. Therefore, we can obtain an edit path by arbitrarily arranging the edit operations in $P_{VD}$. In the rest of this paper, we use $P_{VD}$ to denote the edit operation set and the edit operation sequence interchangeably when there is no ambiguity. Similarly, $P_{VI}$, $P_{VS}$, $P_{ED}$ and $P_O$ could be also considered as the edit operation sets or paths. For an optimal edit path $P$, we can always obtain an optimal edit path $P' = P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS} \cdot P_O$ by arranging the edit operations in $P$. In the following section, we consider $P = P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS} \cdot P_O$ as the default optimal edit path for two given graphs.

LEMMA 1. *Given graphs $g$ and $h$, and an optimal edit path $P$ that transforms $h$ to $g$, then we have $|P_{VD}| = \max\{|V_h| - |V_g|, 0\}$ and $|P_{VI}| = \max\{|V_g| - |V_h|, 0\}$.*

PROOF. Let $P = P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS} \cdot P_O$ be an edit optimal path that transforms $h$ to $g$. Then we discuss the following three cases.

Case I. When $|V_h| = |V_g|$. To transform $h$ to $g$, the number of vertex deletions must be equal to that of vertex insertions, i.e., $|P_{VD}| = |P_{VI}|$. We prove $|P_{VD}| = |P_{VI}| = 0$ by contradiction. Assuming that $|P_{VD}| = |P_{VI}| = l \geq 1$, thus there must exist at least one vertex insertion and one deletion. Let $u$ be a deleted vertex and $v$ be a inserted vertex. We construct another edit path $P'$ by $P$ as follows. First, we substitute the label of $u$ with $\mu(v)$, and then perform these edit operations on $u$, which were performed on $v$ before in $P_O$. Finally, we maintain the rest edit operations in $P$. In other words, we replace $(u \rightarrow \epsilon)$ and $(\epsilon \rightarrow v)$ by $(u \rightarrow v)$. The length of $P'$ is $|P'| = |P| - 1 < |P|$, which contradicts the hypothesis that $P$ is an optimal edit path. Therefore, there exists no vertex deletions and insertions in $P$, i.e., $|P_{VD}| = |P_{VI}| = 0$.

Case II. When $|V_h| < |V_g|$. There exists at least $|V_g| - |V_h|$ vertex insertions in $P$. Let $h_1$ be the graph obtained by inserting $|V_g| - |V_h|$ vertices into $h$. According to the analysis in case I, no vertex deletions and insertions are needed in an optimal edit path that transforms $h_1$ to $g$. Thus, only $|V_g| - |V_h|$ vertex insertions are needed in $P$, i.e., $|P_{VI}| = |V_g| - |V_h|$ and $|P_{VD}| = 0$.

Case III. When $|V_h| > |V_g|$. The proof is similar to the proof of case II. We omit it here. □

### 3.2 Q-gram Counting Filters

DEFINITION 4 (DEGREE-BASED $q$-gram). *Let $D_v = (\mu(v), adj(v), d_v)$ be the degree structure of vertex $v$ in graph $g$, where $\mu(v)$ is the label of $v$, $adj(v)$ is the multiset of labels for edges adjacent to $v$ in $g$, and $d_v$ is the degree of $v$. The degree-based $q$-gram set of graph $g$ is defined as $D(g) = \{D_v : v \in V_g\}$.*

LEMMA 2. *Given two graphs $g$ and $h$, if $ged(g, h) \leq \tau$, then we have $|D(g) \cap D(h)| \geq 2\max\{|V_g|, |V_h|\} - |\Sigma_{V_g} \cap \Sigma_{V_h}| - 2\tau$.*

PROOF. First, we enumerate the effect of various edit operations on $D(g)$: (1) vertex insertion/deletion/substitution will affect one degree-based $q$-gram. (2) edge insertion/deletion/ substitution will affect two degree-based $q$-grams. Then, without loss of generality we assume that $|V_h| \leq |V_g|$ and prove Lemma 2 as follows.

Let $P = P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS} \cdot P_O$ be an optimal edit path that transforms $h$ to $g$, such that: $h \rightarrow h_1 \rightarrow g$, where $h_1$ is obtained by performing $P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS}$ on $h$, and $g$ is obtained by performing $P_O$ on $h_1$. By Lemma 1, we know that $|P_{VD}| = 0$ and $|P_{VI}| = |V_g| - |V_h|$. Since $P_{ED}$ consists of the edge deletions performed on the deleted vertices, we have $|P_{ED}| = 0$. To transform $h$ to $g$, $|P_{VS}|$ vertex substitutions are needed, thus $|P_{VS}| \geq |V_g| - (|\Sigma_{V_g} \cap \Sigma_{V_h}| + |P_{VI}|) =$

3

$|V_h| - |\Sigma_{V_g} \cap \Sigma_{V_h}|$. Since vertex insertion/substitution only affects one degree-based $q$-gram, we have $|D(g) \cap D(h)| \geq |D(g) \cap D(h_1)| - (|P_{VI}| + |P_{VS}|)$. Since $P_O$ only consists of the edit operations performed on edges and each of them affects two degree-based $q$-grams, we have $|D(g) \cap D(h_1)| \geq |V_g| - 2|P_O|$. Thus we have $|D(g) \cap D(h)| \geq |V_g| - 2|P_O| - (|P_{VI}| + |P_{VS}|) \geq 2|V_g| - |\Sigma_{V_g} \cap \Sigma_{V_h}| - 2\tau$. $\square$

DEFINITION 5 (LABEL-BASED $q$-gram). *The label-based $q$-gram set of graph $g$ is defined as $L(g) = \Sigma_{V_g} \cup \Sigma_{E_g}$, where $\Sigma_{V_g}$ and $\Sigma_{E_g}$ are the label multisets of $V_g$ and $E_g$, respectively.*

For the label-based $q$-gram, each edit operation affects one $q$-gram, thus we can obtain the label-based $q$-gram counting filter as follows. If $ged(g, h) \leq \tau$, then we have $|L(g) \cap L(h)| \geq \max\{|V_g|, |V_h|\} + \max\{|E_g|, |E_h|\} - \tau$. It is a rewritten form of the label count filter [24].

Figure 3 shows the degree-based $q$-gram and label-based $q$-gram sets of graphs shown in Figure 2. Note that the number on the left of each subgraph is the times of the $q$-gram occurring in the graph and we omit the degree value of each degree-based $q$-gram.
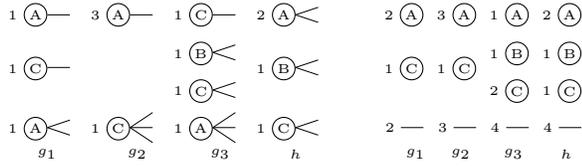


Figure 3: Degree-based $q$-gram (left) and label-based $q$-gram (right) sets.

We use an example to illustrate the degree-based $q$-gram and label-based $q$-gram counting filters. For the graphs $g_2$ and $h$ shown in Figure 2, if $\tau = 2$, by Lemma 2 we have $|D(g_2) \cap D(h)| = 0 < 2 \times \max\{4, 4\} - |\{A, A, B, C\} \cap \{A, A, A, C\}| - 2 \times 2 = 1$. Thus, $g_2$ will be filtered out. However, for the graph $g_1$, we have $|D(g_1) \cap D(h)| = 1 \geq 2 \times \max\{4, 3\} - |\{A, A, B, C\} \cap \{A, A, C\}| - 2 \times 2 = 1$ and hence $g_1$ will pass the filter. Similarly, only $g_1$ will be filtered out by the label-based $q$-gram counting filter Therefore, we can filter $g_1$ and $g_2$ out using the degree-based $q$-gram and label-based $q$-gram counting filters. However, for the graph $g_3$ shown in Figure 2, none of the above filters can filter it out. So, we propose another filter, called degree-sequence filter, which utilizes the degrees of vertices.

## 3.3 Degree-Sequence Filter

Let $\pi_g = [d_1, d_2, \ldots, d_{|V_g|}]$ be the degree vector of graph $g$, where $d_i$ is the degree of vertex $v_i$ in $g$. The degree sequence $\sigma_g$ of $g$ is a permutation of $d_1, d_2, \ldots, d_{|V_g|}$ satisfying $\sigma_g[i] \geq \sigma_g[j]$ for $i < j$. If $g$ is isomorphic to $h$, then we have $\sigma_g = \sigma_h$. Therefore, we can compute the lower bound on $ged(g, h)$ using $\sigma_g$ and $\sigma_h$.

DEFINITION 6 (DEGREE VECTOR DISTANCE). *Given two degree vectors $\pi_g$ and $\pi_h$ such that $|\pi_g| = |\pi_h|$. The distance between them is defined as $\Delta(\pi_g, \pi_h) = \lceil \sum_{\pi_h[i] \leq \pi_g[i]} (\pi_g[i] - \pi_h[i])/2 \rceil + \lceil \sum_{\pi_h[i] > \pi_g[i]} (\pi_h[i] - \pi_g[i])/2 \rceil$.*

LEMMA 3. *Let $x$ and $y$ be two degree vectors such that $|x| = |y| = n$ in non-increasing order. For any bijection function $f: \{1, \ldots, n\} \to \{1, \ldots, n\}$, we have $\Delta(x, z) \geq \Delta(x, y)$, where $z[i] = y[f(i)]$ for $1 \leq i \leq n$.*

PROOF. For degree vectors $x$ and $y$, let $s_1(x, y) = \sum_{x[i] \leq y[i]} (y[i] - x[i])$ and $s_2(x, y) = \sum_{x[i] > y[i]} (x[i] - y[i])$. We have $\Delta(x, y) = \lceil s_1(x, y)/2 \rceil + \lceil s_2(x, y)/2 \rceil$. We want to prove $s_1(x, z) \geq s_1(x, y)$ and $s_2(x, z) \geq s_2(x, y)$, where $z[i] = y[f(i)]$ for $1 \leq i \leq n$. We prove this claim for $x$ and $z$ by induction on the vector length $n$. And similar claim holds for $y$ and $z$.

For the base case $n = 1$, it is trivial that $s_1(x, z) \geq s_1(x, y)$ and $s_2(x, z) \geq s_2(x, y)$. For the inductive step, we assume that $s_1(x^k, z^k) \geq s_1(x^k, y^k)$ and $s_2(x^k, z^k) \geq s_2(x^k, y^k)$ for $n \leq k$ where $x^k = [x[1], \ldots, x[k]]$.

We then prove the claim holds for $n = k+1$. First, without loss of generality, we assume that $f(k + 1) = i$ $(i < k + 1)$, and $x[i] \geq y[i]$, thus we have $s_1(x^k, z^k) = s_1(x^{k-1}, z^{k-1})$ and $s_2(x^k, z^k) = s_2(x^{k-1}, z^{k-1}) + x[i] - y[i]$. Then we consider the following three cases.

Case I. When $x[i] \geq x[k + 1] \geq y[i] \geq y[k + 1]$.

$$s_1(x^{k+1}, z^{k+1}) = s_1(x^k, z^k)$$
$$\geq s_1(x^k, y^k) = s_1(x^{k+1}, y^{k+1})$$
$$s_2(x^{k+1}, z^{k+1}) = s_2(x^{k-1}, z^{k-1}) + x[i] - y[i] + x[k+1] - y[k+1]$$
$$= s_2(x^k, z^k) + x[k+1] - y[k+1]$$
$$\geq s_2(x^k, y^k) + x[k+1] - y[k+1]$$
$$= s_2(x^{k+1}, y^{k+1}).$$

Case II. When $x[i] \geq y[i] \geq x[k + 1] \geq y[k + 1]$.

$$s_1(x^{k+1}, z^{k+1}) = s_1(x^{k-1}, z^{k-1}) + y[i] - x[k+1]$$
$$= s_1(x^k, z^k) + y[i] - x[k+1]$$
$$\geq s_1(x^k, y^k) = s_1(x^{k+1}, y^{k+1}).$$
$$s_2(x^{k+1}, z^{k+1}) = s_2(x^{k-1}, z^{k-1}) + x[i] - y[k+1]$$
$$= s_2(x^k, z^k) - (x[i] - y[i]) + x[i] - y[k+1]$$
$$= s_2(x^k, z^k) + y[i] - y[k+1]$$
$$\geq s_2(x^k, y^k) + x[k+1] - y[k+1]$$
$$= s_2(x^{k+1}, y^{k+1}).$$

Case III. When $x[i] \geq y[i] \geq y[k + 1] \geq x[k + 1]$.

$$s_1(x^{k+1}, z^{k+1}) = s_1(x^{k-1}, z^{k-1}) + y[i] - x[k+1]$$
$$= s_1(x^k, z^k) + y[i] - x[k+1]$$
$$\geq s_1(x^k, y^k) + y[k+1] - x[k+1]$$
$$= s_1(x^{k+1}, y^{k+1}).$$
$$s_2(x^{k+1}, z^{k+1}) = s_2(x^{k-1}, z^{k-1}) + x[i] - y[k+1]$$
$$= s_2(x^k, z^k) - (x[i] - y[i]) + x[i] - y[k+1]$$
$$= s_2(x^k, z^k) + y[i] - y[k+1]$$
$$\geq s_2(x^k, y^k) = s_2(x^{k+1}, y^{k+1}).$$

Finally, we note that $s_1(x, z) \geq s_1(x, y)$ and $s_2(x, z) \geq s_2(x, y)$, and hence we have $\Delta(x, z) \geq \Delta(x, y)$. $\square$

LEMMA 4. *Given two graphs $g$ and $h$ with $|V_g| = |V_h|$, then we have $ged(g, h) \geq \Delta(\sigma_g, \sigma_h)$.*

PROOF. Let $f$ be the bijection from the vertices in $h$ to that in $g$ to ensure that the induced edit path is an optimal edit path. Assuming that $\sigma_h[i]$ and $\sigma_g[f(i)]$ be the respective degrees of a vertex $v$ in $h$ and the corresponding

vertex $u$ in $g$. If $\sigma_h[i] \leq \sigma_g[f(i)]$, we must insert at least $(\sigma_g[f(i)] - \sigma_h[i])$ edges on $v$; otherwise, we must delete at least $(\sigma_h[i] - \sigma_g[f(i)])$ edges. Since one edge insertion/deletion affects degrees of two vertices, we must insert at least $\lceil \sum_{\sigma_h[i] \leq \sigma_g[f(i)]} (\sigma_g[f(i)] - \sigma_h[i])/2 \rceil$ edges. Similarly, we also need to delete at least $\lceil \sum_{\sigma_h[i] > \sigma_g[f(i)]} (\sigma_h[i] - \sigma_g[f(i)])/2 \rceil$ edges. Thus, we have $ged(g, h) \geq \Delta(\sigma_h, \pi'_g)$, where $\pi'_g[i] = \sigma_g[f(i)]$ for $1 \leq i \leq |V_h|$. By Lemma 3, we have $ged(g, h) \geq \Delta(\sigma_h, \pi'_g) \geq \Delta(\sigma_g, \sigma_h)$. $\square$

LEMMA 5 (DEGREE-SEQUENCE FILTER). *Given two graphs $g$ and $h$, and an edit distance threshold $\tau$, if $ged(g, h) \leq \tau$, then we have $\tau \geq \max\{|V_g|, |V_h|\} - |\Sigma_{V_g} \cap \Sigma_{V_h}| + \lambda_e$, where*

$$\lambda_e = \begin{cases} \Delta(\sigma_g, \sigma_1) & \text{if } |V_h| \leq |V_g|; \\ \min_{h_1}\{|E_h| - \sum_j \sigma_{h_1}[j] + \Delta(\sigma_g, \sigma_{h_1})\} & \text{otherwise.} \end{cases}$$

*$\sigma_1 = [\sigma_h[1], \ldots, \sigma_h[|V_h|], 0, \ldots, 0_{|V_g| - |V_h|}]$ and $h_1$ is a subgraph of $h$ obtained by deleting $|V_h| - |V_g|$ vertices.*

PROOF. Let $P = P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS} \cdot P_O$ be an optimal edit path that converts $h$ to $g$, satisfying $h \to h_1 \to h_2 \to g$, where $h_1$ is obtained by performing $P_{ED} \cdot P_{VD}$ on $h$, $h_2$ is obtained by performing $P_{VI} \cdot P_{VS}$ on $h_1$ and $g$ is obtained by performing $P_O$ on $h_2$. Then we discuss the following two cases.

Case I. When $|V_h| \leq |V_g|$. We have $|P_{VD}| = 0$ and $|P_{VI}| = |V_g| - |V_h|$ and $|P_{ED}| = 0$ by Lemma 1. To transform $h_1$ to $h_2$, $|P_{VS}|$ vertex substitutions are needed in $P$, thus we have $|P_{VS}| \geq |V_g| - (|P_{VI}| + |\Sigma_{V_g} \cap \Sigma_{V_h}|) = |V_h| - |\Sigma_{V_g} \cap \Sigma_{V_h}|$. Since $h_2$ is obtained by performing $P_{ED} \cdot P_{VD} \cdot P_{VI} \cdot P_{VS}$ on $h$, we have $\sigma_{h_2} = [\sigma_h[1], \ldots, \sigma_h[|V_h|], 0, \ldots, 0_{|V_g| - |V_h|}]$. By Lemma 4, we have $|P_O| = ged(g, h_2) \geq \Delta(\sigma_g, \sigma_{h_2})$. Therefore $ged(g, h) = |P| = |P_{VI}| + |P_{VS}| + |P_O| \geq |V_g| - |\Sigma_{V_g} \cap \Sigma_{V_h}| + \Delta(\sigma_1, \sigma_g)$.

Case II. When $|V_h| > |V_g|$. We have $|P_{VI}| = 0$ and $|P_{VD}| = |V_h| - |V_g|$ by Lemma 1. To transform $h$ to $h_1$, the number of edge deletions in $P_{ED}$ is $|P_{ED}| = |E_h| - |E_{h_1}| = |E_h| - \sum_j \sigma_{h_1}[j]/2$. Since only $|P_{VS}|$ vertex substitutions are needed to transform $h_1$ to $h_2$, we have $|P_{VS}| \geq |V_h| - (|P_{VD}| + |\Sigma_{V_g} \cap \Sigma_{V_h}|) = |V_g| - |\Sigma_{V_g} \cap \Sigma_{V_h}|$ and $\sigma_{h_1} = \sigma_{h_2}$. By Lemma 4, we also have $|P_O| \geq \Delta(\sigma_g, \sigma_{h_2}) = \Delta(\sigma_g, \sigma_{h_1})$. Therefore $|P| = |P_{ED}| + |P_{VD}| + |P_{VS}| + |P_O| \geq \min_{h_1}\{|V_h| - |\Sigma_{V_g} \cap \Sigma_{V_h}| + |E_h| - \sum_j \sigma_{h_1}[j]/2 + \Delta(\sigma_g, \sigma_{h_1})\} = |V_h| - |\Sigma_{V_g} \cap \Sigma_{V_h}| + \min_{h_1}\{|E_h| - \sum_j \sigma_{h_1}[j]/2 + \Delta(\sigma_g, \sigma_{h_1})\}$. $\square$

We use an example to illustrate the degree-sequence filter. For the graphs $h$ and $g_3$ shown in Figure 2, we can compute $\sigma_h = [2, 2, 2, 2]$ and $\sigma_{g_3} = [3, 2, 2, 1]$. By Lemma 5, if $\tau = 2$, then we have $\max\{4, 4\} - |\{A, A, B, C\} \cap \{A, B, C, C\}| + \Delta(\sigma_h, \sigma_{g_3}) = 4 - 3 + \lceil (3-2)/2 \rceil + \lceil (2-1)/2 \rceil = 3 > 2$, then we can filter $g_3$ out.

## 4. REDUCED QUERY REGION

Given a database $G$, we consider each graph $g$ in $G$ as a point in the two-dimensional plane where the x-coordinate and y-coordinate denote the number of vertices and edges in $g$, respectively. Thus the graph database $G$ can be represented as a set of points $S = \{(|V_{g_j}|, |E_{g_j}|) : 1 \leq j \leq |G|\}$. These points form a rectangle area $A = [x_{min}, x_{max}] \times [y_{min}, y_{max}]$, where $x_{min} = \min_j\{|V_{g_j}|\}$, $x_{max} = \max_j\{|V_{g_j}|\}$, $y_{min} = \min_j\{|E_{g_j}|\}$ and $y_{max} = \max_j\{|E_{g_j}|\}$ for $1 \leq j \leq$

$|G|$. By partitioning $A$ into subregions, we can perform a query at a reduced query region.

Given an initial division point $(x_0, y_0)$ and a length $l$, we partition $A$ into disjoint subregions as follows. First, we construct the initial square subregion $A_{0,0}$ formed by the point set $\{(x, y) : |x - x_0| + |y - y_0| \leq l\}$. Then, we extend along the surrounding of $A_{0,0}$ to obtain subregions $A_{i,j}$ of the same size with $A_{0,0}$, where $i$ and $j$ denote the relative offsets with respect to $A_{0,0}$ in lines $y = x$ and $y = -x$, respectively. Finally, we repeat this process until all points in $A$ are exhausted. Then $A$ is partitioned into some disjoint subregions such that $A = \cup_{i,j} A_{i,j}$ and $A_{i,j} \cap A_{i',j'} = \varnothing$ for all $i \neq i'$ and $j \neq j'$. Note that $i$ and $j$ can be negative.

DEFINITION 7 (QUERY RECTANGLE AND REGION). *Given a query graph $h$ and an edit distance threshold $\tau$, query rectangle $A_h$ of $h$ is the rectangle formed by the point set of $\{(x, y) : |x - |V_h|| + |y - |E_h|| \leq \tau\}$. The query region $Q_h$ of $h$ is the union of all subregions intersecting with $A_h$, i.e., $Q_h = \cup_{i,j} A_{i,j}$ such that $A_{i,j} \cap A_h \neq \varnothing$.*

For graphs $g$ and $h$, if $ged(g, h) \leq \tau$, then we have $||V_g| - |V_h|| + ||E_g| - |E_h|| \leq \tau$. According to the definition of $A_h$, we know that $(|V_g|, |E_g|) \in A_h$. Since $Q_h = \cup_{i,j} A_{i,j}$ and $A_{i,j} \cap A_h \neq \varnothing$, we have $A_h \subseteq Q_h$. Therefore we have $(|V_g|, |E_g|) \in Q_h$ and hence can reduce the query region from $A$ to $Q_h$. In the example of Figure 4, we have $Q_h = \{A_{0,0}, A_{1,0}, A_{0,-1}, A_{1,-1}\}$ and then only need to perform the query at $Q_h$.
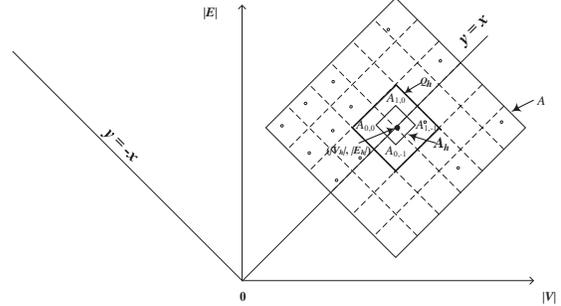


Figure 4: Illustration of $A_h$, $Q_h$ and $A$

For a two-dimensional point $(x, y)$, its coordinates in lines $y = x$ and $y = -x$ are $\frac{1}{\sqrt{2}}(x + y, y - x)$, thus its relative offsets with respect to $(x_0, y_0)$ are $d_x = \frac{1}{\sqrt{2}}((x + y) - (x_0 + y_0))$ and $d_y = \frac{1}{\sqrt{2}}((y - x) - (y_0 - x_0))$ in $y = x$ and $y = -x$, respectively. Since the side length of a subregion is $\frac{l}{\sqrt{2}}$, the respective relative offsets with respect to $A_{0,0}$ are $\lfloor \frac{d_x}{l/\sqrt{2}} \rfloor$ and $\lfloor \frac{d_y}{l/\sqrt{2}} \rfloor$ in $y = x$ and $y = -x$. Since the subregions in $Q_h$ are adjacent, we just need to find the boundaries of subregions intersecting with $A_h$ using the following formula.

$$Q_h = \cup_{i,j} A_{i,j} \text{ for all } i_1 \leq i \leq i_2 \text{ and } j_1 \leq j \leq j_2. \quad (1)$$

where $i_1 = \lfloor (|E_h| - \tau + |V_h| - (x_0 + y_0))/l \rfloor$ and $j_1 = \lfloor (|E_h| - \tau - |V_h| - (y_0 - x_0))/l \rfloor$ are the relative positions of the subregion in the lower left corner of $Q_h$ with respect to $A_{0,0}$ in $y = x$ and $y = -x$, respectively, $i_2 = \lfloor (|E_h| + \tau + |V_h| - (x_0 + y_0))/l \rfloor$ and $j_2 = \lfloor (|E_h| + \tau - |V_h| - (y_0 - x_0))/l \rfloor$ are the respective relative positions of the subregion in the top right corner of $Q_h$ with respect to $A_{0,0}$ in $y = x$ and $y = -x$.

# 5. SUCCINCT Q-GRAM TREE INDEX

Recall that we partitioned the region $A$ into some subregions and then obtained a reduced query region $Q_h$. In order to efficiently filter the graphs mapped into $Q_h$, we introduce a space-efficient index structure via succinct representation of the $q$-gram tree as follows.

## 5.1 Tree Structure

Let $\mathcal{U}_D$ and $\mathcal{U}_L$ be the sets of all distinct degree-based $q$-grams and label-based $q$-grams occurring in $G$, respectively, where $\mathcal{U}_D(i)$ and $\mathcal{U}_L(i)$ are the $i$th most frequently occurring degree-based $q$-gram and label-based $q$-gram in $G$, respectively. We use a four-tuple $LD = (F_D, F_L, n_v, n_e)$ to represent a graph $g$, where $n_v$ and $n_e$ are the number of vertices and edges in $g$, respectively, $F_D$ and $F_L$ are two arrays to store the degree-based $q$-gram and label-based $q$-gram sets $D(g)$ and $L(g)$, respectively, where $F_D[i]$ and $F_L[i]$ are the respective number of occurrences of the degree-based $q$-gram $\mathcal{U}_D(i)$ in $D(g)$ and the label-based $q$-gram $\mathcal{U}_L(i)$ in $L(g)$.

DEFINITION 8. *Given two four-tuples $LD$ and $LD'$, the union operator "$\sqcup$" of $LD$ and $LD'$ is defined as: $LD \sqcup LD' = (F_D \oplus F'_D, F_L \oplus F'_L, \min\{n_v, n'_v\}, \min\{n_e, n'_e\})$, where*

$$(F_D \oplus F'_D)[i] = \begin{cases} \max\{F_D[i], F'_D[i]\} & \text{if } i < \min\{|F_D|, |F'_D|\}; \\ F_D[i] & \text{if } |F'_D| \le i < |F_D|; \\ F'_D[i] & \text{if } |F_D| \le i < |F'_D|. \end{cases}$$

*and similar definition for $F_L \oplus F'_L$.*

Similarly, the union of multiple four-tuples can be defined recursively.

DEFINITION 9. *A $q$-gram tree is a balanced tree such that each leaf node stores the four-tuple $LD$ of the data graph $g$ and each internal node is the union of its child nodes.*

Figure 5 gives an example of a $q$-gram tree built on $g_1$, $g_2$, and $g_3$ shown in Figure 2.
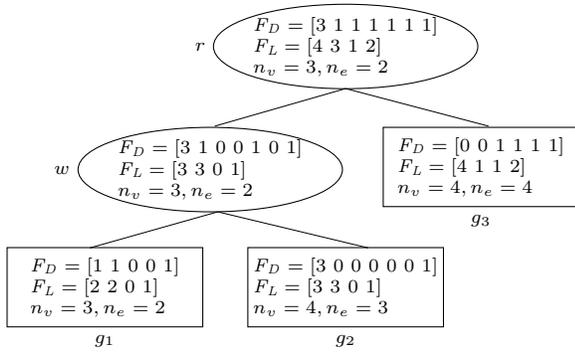


Figure 5: Example of a $q$-gram tree

## 5.2 Succinct Representation

The arrays $F_D$ and $F_L$ may contain lots of zeros, thus a succinct representation of them is a space-efficient way to store them. For a $q$-gram tree, we obtain its succinct representation by performing the following three steps. In the following sections, we refer $X$ to be $D$ or $L$.

(1) We use a bit vector $I_X$ and an array $V_X$ to represent $F_X$ as follows: if $F_X[j] = 0$ then we have $I_X[j] = 0$;

otherwise $I_X[j] = 1$. $V_X[j]$ represents the $j$th nonzero entry in $F_X$. For example, the array $F_D$ in the node $w$ shown in Figure 5 is $F_D = [3\ 1\ 0\ 0\ 1\ 0\ 1]$, then we use $(I_D, V_D) = ([1\ 1\ 0\ 0\ 1\ 0\ 1], [3\ 1\ 1\ 1])$ to represent $F_D$.

(2) We concatenate all bit vectors $I_X$ and arrays $V_X$ for all nodes from the root node to leaves in a depth-first traversal order to obtain a bit vector $B_X$ and an array $\Psi_X$, respectively. In addition, we also store the left and right boundaries $l_X$ and $r_X$ of $I_X$ for each node, respectively. For example, for the $q$-gram tree shown in Figure 5, we can obtain $B_D = [1\ 1\ 1\ 1\ 1\ 1\ 1\ 1\ 1\ 0\ 0\ 1\ 0\ 1\ 1\ 1\ 0\ 0\ 1\ 1\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 1\ 1\ 1]$ and $\Psi_D = [3\ 1\ 1\ 1\ 1\ 1\ 1\ 3\ 1\ 1\ 1\ 1\ 1\ 1\ 3\ 1\ 1\ 1\ 1\ 1]$.

(3) We divide $\Psi_X$ into fixed-length blocks of size $b$ and encode each block by choosing one from two different compression methods so that the encoded bit vector $S_X$ has the minimum space. One compression method uses the fixed-length encoding of $\lfloor \log b_{max} \rfloor + 1$ bits to encode each entry in a fixed-length encoding block, where $b_{max}$ is the maximum value in this block. The other method uses Elias $\gamma$ encoding to encode each entry in a $\gamma$-encoding block. Logarithms in this paper are in base 2 unless otherwise stated.

To support random access to $\Psi_X[j]$, we also need to store three auxiliary structures $SB_X$, $words_X$, and $flag_X$, where $SB_X$ stores the starting position of the encoding of each block in $S_X$; the bit vector $flag_X$ stores the encoding method used in each block such that $flag_X[k] = 1$ for the fixed-length encoding and $flag_X[k] = 0$ for the Elias $\gamma$ encoding for the $k$th block; $words_X$ stores the number of bits required for each entry in a fixed-length encoding block. We also build rank dictionaries over the bit vectors $B_X$ and $flag_X$ to obtain $rank_1(B_X, j)$ and $rank_1(flag_X, j)$ in constant time [7], where $rank_1(B_X, j)$ and $rank_1(flag_X, j)$ are the respective number of 1's up to $j$ in $B_X$ and $flag_X$.

Let $B_D$ and $B_L$ be the respective degree-based and label-based $q$-grams bit vectors, and $\Psi_D$ and $\Psi_L$ be the respective degree-based and label-based $q$-gram frequency arrays. We use four structures $S_D$, $SB_D$, $flag_D$, and $words_D$ to represent $\Psi_D$. Similarly, we use four structures $S_L$, $SB_L$, $flag_L$, and $words_L$ to represent $\Psi_L$. Figure 6 shows the succinct representation of the $q$-gram tree shown in Figure 5.
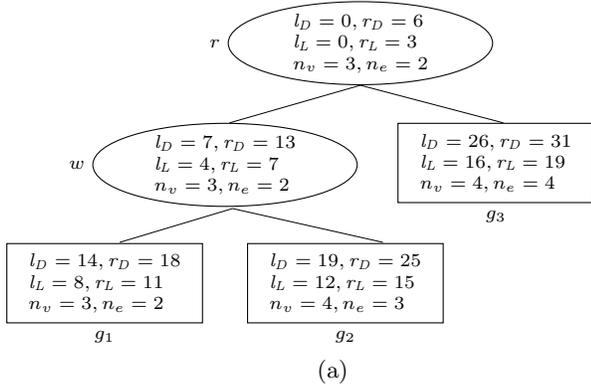
## 5.3 Access to $\Psi_X$

To access $\Psi_X[j]$, we first query $flag_X$ and $SB_X$ to determine the encoding method used and decoding position, respectively, and then decode $S_X$ from the decoding position. The last decoded value is $\Psi_X[j]$.

$$\Psi_X[j] = decompress(S_X, flag_X[\lfloor j/b \rfloor], SB_X[\lfloor j/b \rfloor], (j \bmod b) + 1) \quad (2)$$

where $b$ is the block size. The operation *decompress* performs a decoding on $S_X$. The encoding method and decoding position are determined by the second parameter $flag_X[\lfloor j/b \rfloor]$ and third parameter $SB_X[\lfloor j/b \rfloor]$ of *decompress*, respectively. $(j \bmod b) + 1$ is the number of times needed to be decoded.

For example, if we want to retrieve $\Psi_D[14]$ (suppose that the subscript starts from 0 and $b = 4$) shown in Figure 6, we find that $flag_D[\lfloor 14/b \rfloor] = flag_D[3] = 0$ and $SB_D[\lfloor 14/b \rfloor] = SB_D[3] = 16$. Thus starting from the 16th bit of $S_D$, we sequentially decode Elias $\gamma$ encoding three times and the last decoded value is $\Psi_D[14] = 3$.

| | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1111111 1100101 11001 1000001 001111 | | | | | | $B_D$ | | | | | | | | | | | | | | |

| $j$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\Psi_D$ | 3 | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 1 | 1 | 1 | 1 | 1 |
| $S_D$ | 011 1 1 1 | | | | | 1 1 1 011 | | | | | 1 1 1 1 | | | | 1 1 011 1 | | | | 1 1 1 1 | |
| $SB_D$ | 0 | | | | | 6 | | | | | 12 | | | | 16 | | | | 22 | |
| $flag_D$ | 0 | | | | | 0 | | | | | 1 | | | | 0 | | | | 1 | |
| $words_D$ | – | | | | | – | | | | | 1 | | | | – | | | | 1 | |

(b)

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1111 1101 1101 1101 1111 | | | | | | $B_L$ | | | | | | | | | | | |

| $j$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\Psi_L$ | 4 | 3 | 1 | 2 | 3 | 3 | 1 | 2 | 2 | 1 | 3 | 3 | 1 | 4 | 1 | 1 | 2 |
| $S_L$ | 100 011 001 010 | | | | 11 11 01 10 | | | | 10 01 11 11 | | | | 1 00100 1 1 | | | | 10 |
| $SB_L$ | 0 | | | | 12 | | | | 20 | | | | 28 | | | | 36 |
| $flag_L$ | 1 | | | | 1 | | | | 1 | | | | 0 | | | | 1 |
| $words_L$ | 3 | | | | 2 | | | | 2 | | | | – | | | | 2 |

(c)

Figure 6: Succinct representation of the *q-gram* tree

We use the following formula (3) to compute the original entry $F_X[i]$ in the node $w$.

$$F_X[i] = \begin{cases} 0 & \text{if } B_X[l_X + i] = 0; \\ \Psi_X[rank_1(B_X, l_X + i)] & \text{otherwise.} \end{cases} \quad (3)$$

where $0 \leq i \leq r_X - l_X$, and $l_X$ and $r_X$ are the left and right boundaries of $I_X$ for $w$, respectively. If $B_X[l_X + i] = 0$, then we have $F_X[i] = 0$ since $I_X[i] = B_X[l_X + i] = 0$; otherwise, we first compute the position of $F_X[i]$ in $\Psi_X$, i.e., $rank_1(B_X, l_X + i)$, and then use formula (2) to compute $\Psi_X[rank_1(B_X, l_X + i)]$, where $rank_1(B_X, l_X + i)$ is the number of 1's up to $l_X + i$ in $B_X$, which can be computed in constant time using a dictionary of $o(|B_X|)$ bits [7].

For example, to retrieve $F_D[0]$ in the node $g_2$ shown in Figure 5, we first obtain $B_D[l_D + 0] = B_D[19] = 1$, and then compute its position in $\Psi_D$ is $rank_1(B_D, 19) = 14$, thus we have $F_D[0] = \Psi_D[14] = 3$ by formula (2).

As discussed above, the core operation in a succinct *q-gram* tree is to calculate $\Psi_X[j]$ by formula (2), i.e., the *decompress* operation. In order to accelerate the *decompress* process, we use the look up table technique proposed in [5] to ensure that the *decompress* operation takes a constant time.

## 5.4 Space Analysis

In this section we analyze the space occupied by the succinct *q-gram* tree $T_{SQ}$ built on $G$. $T_{SQ}$ consists of three parts: the respective index structures for $\Psi_D$ and $\Psi_L$, and left and right boundaries, #vertices and #edges in each node of the tree. The former contains encoded sequence $S_X$ and corresponding auxiliary structures $B_X$, $SB_X$, $flag_X$ and $words_X$; The latter consists of $l_X$, $r_X$, $n_v$ and $n_e$ stored in each node of $T_{SQ}$, where $X$ denotes $D$ or $L$. An illustration of these structures is shown in Figure 6.

Let $v_m = \max_j\{|V_{g_j}|\}$, $e_m = \max_j\{|E_{g_j}|\}$ for $1 \leq j \leq |G|$, $n_D = |B_D|$, $n_L = |B_L|$, $b_D^m$ and $b_L^m$ be the maximum value in $\Psi_D$ and $\Psi_L$, respectively. For a degree-based *q-gram*, its maximum number of occurrence in a graph $g$ is $|V_g|$, thus we have $b_D^m \leq v_m$. Similarly, $b_L^m \leq \max\{v_m, e_m\}$. We first consider the space required by $l_X$, $r_X$, $n_v$ and $n_e$.

For any node of $T_{SQ}$, we use $\lfloor \log n_D \rfloor + 1$ bits to store $l_D$ and $r_D$, respectively, and $\lfloor \log n_L \rfloor + 1$ bits to store $l_L$ and $r_L$, respectively, since $l_D \leq r_D \leq n_D$ and $l_L \leq r_L \leq n_L$. We also use respective $\lfloor \log v_m \rfloor + 1$ and $\lfloor \log e_m \rfloor + 1$ bits to store $n_v$ and $n_e$, since $n_v \leq v_m$ and $n_e \leq e_m$. For an average fan-out of $d$ for each node in $T_{SQ}$ with $|G|$ leaf nodes, the total number of nodes in $T_{SQ}$ is bounded by $\sum_{h=0}^{\log_d |G|} \frac{|G|}{d^h} \leq \frac{d|G|}{d-1}$. Thus, we can use $\lfloor \log \frac{d|G|}{d-1} \rfloor + 1$ bits to store each child pointer of a node in $T_{SQ}$. Thus, the total number of bits required by $l_D$, $r_D$, $l_L$, $r_L$, $n_v$, $n_e$ and pointers for all nodes in $T_{SQ}$ is bounded by

$$\frac{d|G|}{d-1}(2(\lfloor \log n_D \rfloor + 1) + 2(\lfloor \log n_L \rfloor + 1) + \lfloor \log v_m \rfloor + 1 +$$

$$\lfloor \log e_m \rfloor + 1 + \lfloor \log \frac{d|G|}{d-1} \rfloor + 1)$$

$$\leq \frac{d|G|}{d-1}(2\log(n_D n_L) + \log(v_m e_m) + \log \frac{d|G|}{d-1} + 7).$$

We then consider the space required by $S_X$, $B_X$, $SB_X$, $flag_X$ and $words_X$.

First we analyze the space needed by the encoded sequence $S_X$. Let $N^g$ and $N^f$ be the respective collection of blocks with $\gamma$ encoding and fixed-length encoding, and $|\gamma(b_i)|$ and $|f(b_i)|$ be the respective number of bits needed to encode the $i$th block $b_i$ using $\gamma$ encoding and fixed-length encoding. By our hybrid encoding scheme, the number of bits required by $S_X$ is bounded by

$$\sum_{i=1}^{|\Psi_X|/b} \min\{|\gamma(b_i)|, |f(b_i)|\}$$

$$= \sum_{i \in N^g} |\gamma(b_i)| + \sum_{i \in N^f} |f(b_i)| \leq \sum_{i \in N^g} |f(b_i)| + \sum_{i \in N^f} |f(b_i)|$$

$$\leq \sum_{i \in N^g \cup N^f} b(\lfloor \log b_X^m \rfloor + 1) \leq \frac{|\Psi_X|}{b} b(\lfloor \log b_X^m \rfloor + 1)$$

$$\leq |\Psi_X| \log b_X^m + |\Psi_X|.$$

where the first inequality is due to the fact that $|\gamma(b_i)| \leq |f(b_i)|$ when $i \in N^g$. The number of bits required to encode block $b_i$ of $\Psi_X$ using fixed-length encoding is bounded by $b(\lfloor \log b_X^m \rfloor + 1)$. The third inequality is due to the fact that $|N^g| + |N^f| = |\Psi_X|/b$, where $b$ is the block size.

Second, we analyze the space required by auxiliary structures $B_X$, $SB_X$, $flag_X$ and $words_X$.

For bit vector $B_X$, the total number of bits required to store it and its rank dictionary is $|B_X| + o(|B_X|)$ bits, where $o(|B_X|)$ is the space in bits required by the rank dictionary built on $B_X$ [7].

For $SB_X$, the space needed is $\frac{|\Psi_X|}{b}(\log(|\Psi_X|\log b_X^m + |\Psi_X|) + 1)$ in bits in the worst case since each entry needs $\lfloor\log(|\Psi_X|\log b_X^m + |\Psi_X|)\rfloor + 1$ bits and there are $|\Psi_X|/b$ blocks.

For $flag_X$, it is trivial that the total number of bits required is $|\Psi_X|/b + o(|\Psi_X|/b)$ bits, since each block takes one bit and there are total $|\Psi_X|/b$ blocks. The rank dictionary built on $flag_X$ needs $o(|\Psi_X|/b)$ bits.

For $words_X$, the space used is bounded by $\frac{|\Psi_X|}{b}(\lfloor\log b_X^m\rfloor + 1)$, since each entry requires $\lfloor\log b_X^m\rfloor + 1$ bits to store and there are $|\Psi_X|/b$ entries in the worst case.

Putting all space needed for auxiliary structures $B_X$, $SB_X$, $flag_X$ and $words_X$ together, we then obtain

$$|B_X| + o(|B_X|) + \frac{|\Psi_X|}{b}\log(|\Psi_X|\log b_X^m + |\Psi_X|) +$$
$$\frac{|\Psi_X|}{b}\log b_X^m + 3\frac{|\Psi_X|}{b} + o(\frac{|\Psi_X|}{b})$$
$$= |B_X| + o(|B_X|) + o(|\Psi_X|), \text{ for } b = \log^2|\Psi_X|.$$

By adding $|\Psi_X|\log b_X^m + |\Psi_X|$ bits required by $S_X$ to the space required by auxiliary structures, we obtain that the space is $|B_X| + o(|B_X|) + |\Psi_X|\log b_X^m + |\Psi_X| + o(|\Psi_X|)$ bits.

By summing up all space for $T_{SQ}$ and replacing $X$ with $D$ or $L$, we obtain that the succinct $q$-$gram$ tree $T_{SQ}$ takes $\frac{d|G|}{d-1}(2\log(n_D n_L) + \log(v_m e_m) + \log\frac{d|G|}{d-1} + 7) + n_D + o(n_D) + n_L + o(n_L) + |\Psi_D|(\log v_m + 1) + |\Psi_L|(\log\max\{v_m, e_m\} + 1) + o(|\Psi_D|) + o(|\Psi_L|)$ bits of space.

## 6. QUERY PROCESSING

Our query process consists of two phrases. We first compute the reduced query region $Q_h$ by formula (1), and then perform the query on the succinct $q$-$gram$ trees built on the graphs mapped into $Q_h$.

### 6.1 Query on Succinct q-gram Tree

We introduce the query method on the succinct $q$-$gram$ tree $T$ in this section.

LEMMA 6. *Let $C_D$ and $C_L$ be the respective number of common degree-based and label-based $q$-grams between any internal node $w$ of $T$ and the query graph $h$, if $C_D < \max\{n_v, |V_h|\} - 2\tau$ or $C_L < \max\{n_v, |V_h|\} + \max\{n_e, |E_h|\} - \tau$, then we can safely prune all child nodes of $w$, where $n_v$ and $n_e$ are the number of vertices and edges in $w$, respectively.*

PROOF. Let $x.LD = (x.F_D, x.F_L, x.n_v, x.n_e)$ denote the four-tuple of $x$, where $x$ is a node in $T$. For any internal node $w$ and a query graph $h$, we have $C_D = \sum_i\min\{w.F_D[i], h.F_D[i]\}$. According to Definition 8, for a child node $w_j$ of $w$ we have $w_j.F_D[i] \leq w.F_D[i]$. Therefore, for a descendent leaf node (i.e., graph) $g$ of $w$, we have $|D(g) \cap D(h)| = \sum_i\min\{g.F_D[i], h.F_D[i]\} \leq \cdots \leq \sum_i\min\{w_j.F_D[i], h.F_D[i]\} \leq \sum_i\min\{w.F_D[i], h.F_D[i]\} = C_D$. Similarly, we also have $w.n_v \leq w_j.n_v \leq \cdots \leq g.n_v = |V_g|$ and $w.n_e \leq w_j.n_e \leq \cdots \leq g.n_e = |E_g|$. If $C_D < \max\{w.n_v, |V_h|\} - 2\tau$, then we have $|D(g) \cap D(h)| \leq C_D < \max\{w.n_v, |V_h|\} - 2\tau \leq \max\{|V_g|, |V_h|\} - 2\tau \leq 2\max\{|V_g|, |V_h|\} - |\Sigma_{V_g} \cap \Sigma_{V_h}| - 2\tau$ and can safely prune graph $g$ by Lemma 3. Similarly, if $C_L < \max\{w.n_v, |V_h|\} +$

$\max\{w.n_e, |E_h|\} - \tau$, then we have $|L(g) \cap L(h)| < \max\{|V_g|, |V_h|\} + \max\{|E_g|, |E_h|\} - \tau$ and then can safely prune $g$ by the label-based $q$-$gram$ counting filter. So, we can safely prune all child nodes of $w$. □

Algorithm 1 gives the query algorithm on $T$, where $r$ is the root node, $(F_D, F_L, n_v, n_e)$ is the four-tuple of $w$, $l_D$ and $r_D$ are the left and right boundaries of $I_D$ for $w$ in $T$, respectively.

---

**Algorithm 1:** $searchQTree(T, h, \tau)$

---

**Input**: $T, h, \tau$
**Output**: $Cand$

1   $Cand \leftarrow \varnothing$
2   Compute the four-tuple $LD' = (F_D', F_L', n_v', n_e')$ and degree sequence $\sigma_h$ of $h$
3   $searchTree(r, LD', \sigma_h, h, \tau)$
4   **return** $Cand$
5   **Procedure** $searchTree(w, LD', \sigma_h, h, \tau)$
6     $C_L \leftarrow \sum_i\min\{F_L[i], F_L'[i]\}$
7     **if** $C_L \geq \max\{n_v, |V_h|\} + \max\{n_e, |E_h|\} - \tau$ **then**
8       $C_D \leftarrow \sum_i\min\{F_D[i], F_D'[i]\}$
9       **if** $C_D \geq \max\{n_v, |V_h|\} - 2\tau$ **then**
10         **if** $w$ *is an internal node* **then**
11           **for** *each child $w_i$ of $w$* **do**
12             $searchTree(w_i, LD', \sigma_h, h, \tau)$
13         **if** $C_D \geq 2\max\{|V_h|, |V_w|\} - |\Sigma_{V_w} \cap \Sigma_{V_h}| - 2\tau$ **then**
14           **for** $i \leftarrow l_D$ **to** $r_D$ **do**
15             $F_w[i - l_D] \leftarrow F_D[i - l_D]$
16           Obtain the degree sequence $\sigma_w$ of $w$
17           Compute the lower bound $\xi$
18           **if** $\xi \leq \tau$ **then**
19             $Cand \leftarrow Cand \cup \{w\}$

---

In Algorithm 1, we first compute the four-tuple $LD'$ and degree sequence $\sigma_h$ of $h$ in line 2, respectively, and then perform the search processing $searchTree$ starting from a node $w$ initialized to $r$, the root node of $T$ as follows. First, we determine whether a node $w$ needs to be pruned based upon Lemma 6 in lines 6–12. In lines 6 and 8, we compute the number of common label-based $q$-grams $C_L$ and degree-based $q$-grams $C_D$ between $w$ and $h$, respectively. Note that, each entry $F_D[i]$ and $F_L[i]$ are compressed in $\Psi_D$ and $\Psi_L$, respectively, thus we need to use formula (3) to compute them. If $C_L < \max\{n_v, |V_h|\} + \max\{n_e, |E_h|\} - \tau$ or $C_D < \max\{n_v, |V_h|\} - 2\tau$, we prune $w$; otherwise each subtree of $w$ will be accessed in lines 11–12. Then, we determine whether a node $w$ needs to be pruned based upon Lemma 2 in line 13, i.e., the degree-based $q$-gram counting filter. If $C_D < 2\max\{n_v, |V_h|\} - |\Sigma_{V_w} \cap \Sigma_{V_h}| - 2\tau$, we prune $w$, where $|\Sigma_{V_w} \cap \Sigma_{V_h}|$ is the number of common vertex labels between $w$ and $h$ obtained while computing $C_L$. Finally, we first obtain the degree sequence $\sigma_w$ and then determine whether $w$ needs to be pruned based upon Lemma 5, i.e., the degree-sequence filter. In lines 14–15, we first compute the array $F_w$ storing the degree-based $q$-gram set of $w$ and then obtain $\sigma_w$ using $F_w$ and $T_D$ in line 16, where $T_D$ is a table storing the mapping between a degree-based $q$-gram and its identifier. If $\xi > \tau$, then

we prune $w$; otherwise, it passes all filters to become a candidate.

## 6.2 Query Algorithm

Algorithm 2 gives the whole query algorithm, where $(x_0, y_0)$ is the initial division point, $l$ is the subregion length and $T_{i,j}$ is the succinct $q$-$gram$ tree built on the graphs mapped into the subregion $A_{i,j}$.

---

**Algorithm 2:** $search(h, \tau, x_0, y_0, l)$

---
**Input**: $h, \tau, x_0, y_0, l$
**Output**: $Cand$
1   $Cand \leftarrow \varnothing$
2   $Q_h \leftarrow \cup_{i,j} A_{i,j}$ for all $i_1 \leq i \leq i_2$ and $j_1 \leq j \leq j_2$
3   **foreach** $A_{i,j} \subseteq Q_h$ **do**
4      $C_{i,j} \leftarrow searchQTree(T_{i,j}, h, \tau)$
5      $Cand \leftarrow Cand \cup C_{i,j}$
6   **return** $Cand$

---

In Algorithm 2, we first compute the query region $Q_h$ by formula (1) in line 2, where $i_1 = \lfloor (|E_h| - \tau + |V_h| - (x_0 + y_0))/l \rfloor$, $i_2 = \lfloor (|E_h| + \tau + |V_h| - (x_0 + y_0))/l \rfloor$, $j_1 = \lfloor (|E_h| - \tau - |V_h| - (y_0 - x_0))/l \rfloor$ and $j_2 = \lfloor (|E_h| + \tau - |V_h| - (y_0 - x_0))/l \rfloor$. Then we only need to perform the query on the $q$-$gram$ trees $T_{i,j}$ built on these subregions $A_{i,j}$ satisfying $A_{i,j} \subseteq Q_h$ in lines 3–5. For each candidate graph $g$ in $Cand$, we can use the methods in [10, 14, 24] to compute the edit distance between $g$ and $h$ to seek for the required graphs.

## 7. EXPERIMENTAL RESULTS

In this section, we evaluate the performance of our proposed method and compare it with C-Star [22], GSimJoin [24] and Mixed [26] on the real and synthetic datasets.

## 7.1 Datasets and Settings

We choose several real and synthetic datasets to test the performance of the above approaches in our experiment, described as follows:

(1) AIDS[1]. It is a DTP AIDS antivirus screen compound dataset from the Development and Therapeutics Program in NCI/NIH to discover compounds capable of inhibiting the HIV virus. It contains 42687 chemical compounds. We generate the labeled graphs from these chemical compounds and omit Hydrogen atoms as did in [18].

(2) PubChem[2]. It is a NIH funded project to record experimental data of chemical interactions with biological systems. It contains more than 50 million chemical compounds until today. We randomly select 25 million chemical compounds to make up the large dataset PubChem-25M used in this experiment.

(3) Synthetic. The synthetic datasets are generated by the synthetic graph data generator GraphGen[3]. The synthetic generator can create a labeled and undirected graph dataset. It allows us to specify various parameters, including the dataset size, the average graph density $\rho = \frac{2|E|}{|V|(|V|-1)}$, the number of edges in a graph, and the number of distinct vertex and edge labels in the dataset, respectively. In order

---

[1] http://dtp.nci.nih.gov/docs/aids/aidsdata.html
[2] http://pubchem.ncbi.nlm.nih.gov/
[3] http://www.cse.ust.hk/graphgen/

Table 1: Statistics of the three data sets

| DataSet | $|G|$ | $|V|$ | $|E|$ | $|\Sigma_V|$ | $|\Sigma_E|$ |
|---|---|---|---|---|---|
| AIDS | 42687 | 25.6 | 27.5 | 62 | 3 |
| S100K.E30.D50.L5 | 100,000 | 11.02 | 30 | 5 | 2 |
| PubChem-25M | 25,000,000 | 23.4 | 25.2 | 101 | 3 |

to evaluate the performance of the above approaches on the density graphs, we generate the dataset S100K.E30.D50.L5, which means that this dataset contains 100000 graphs; the average density of each graph is 50%; the number of edges in each graph is 30; and the number of distinct vertex and edge labels are 5 and 2, respectively.

For each dataset, we randomly select 50 graphs from it as its query graphs. Table 1 summarizes some general characteristics of the three datasets described above.

We have conducted all experiments on a HP Z800 PC with a 2.67 GHz CPU and 24GB memory, running Ubuntu 12.04 operating system. We implemented our algorithm in C++, with $-O3$ to compile and run. For GSimJoin, we set $p = 4$ for the sparse graphs in datasets AIDS and PubChem-25M, and $p = 3$ for the density graphs in dataset S100K.E30.D50.L5, which are the recommended values [24]. In the following sections, we refer MSQ-Index to our index structure and set the subregion length $l = 4$ and block size $b = 16$, respectively.

## 7.2 Index Construction and Space Usage

In this section, we introduce extensive experiments to evaluate index construction performance of C-Star, GSimJoin, Mixed and MSQ-Index.

### 7.2.1 Evaluating Our Index

In order to evaluate the effectiveness of our hybrid encoding, we compare it with fixed-length encoding, Elias $\delta$ encoding, Golomb encoding, and Elias $\gamma$ encoding.

For each dataset described in Table 1, we show the number of bits on the average required by each entry in $\Psi_D$ and $\Psi_L$ in Table 2 when applying fixed-length encoding ($f$), Elias $\delta$ encoding ($\delta$), Golomb encoding ($g$), Elias $\gamma$ encoding ($\gamma$) and hybrid encoding ($h$) to $\Psi_D$ and $\Psi_L$, where S100K and Pub-25M stand for S100K.E30.D50.L5 and PubChem-25M, respectively. Each entry in $\Psi_D$ and $\Psi_L$ uses about between 4 and 6 bits on the tested data, which is much smaller than that used to represent an entry in the previous state-of-the-art indexing methods compared in this paper.

Table 2: Average number of bits

| $Datasets$ | $\Psi_D$ | | | | | $\Psi_L$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $f$ | $g$ | $\delta$ | $\gamma$ | $h$ | $f$ | $g$ | $\delta$ | $\gamma$ | $h$ |
| AIDS | 4.51 | 4.25 | 3.68 | 3.57 | **3.51** | 5.97 | 6.23 | 6.17 | 6.13 | **5.93** |
| S100K | 3.33 | 4.04 | 3.31 | 3.18 | **3.08** | 4.01 | 4.75 | 4.5 | 4.27 | **3.88** |
| Pub-25M | 4.57 | 4.19 | 3.62 | 3.61 | **3.36** | 5.73 | 6.07 | 5.85 | 5.79 | **5.31** |

Among all the encoding methods shown in Table 2, hybrid encoding gives the minimum space. Compared with fixed-length encoding, the average number of bits required for hybrid encoding decreases by about 10%.

In Table 3, we report the storage space of the $q$-$gram$ tree $T_Q$ and its succinct representation $T_{SQ}$ built on the above three datasets. For $T_Q$, we decompose its storage space into three parts $S_a$, $S_b$ and $S_c$, where $S_a$ is the storage

Table 3: Storage space (MByte)

| Datasets | $T_Q$ | | | $T_{SQ}$ | | |
|---|---|---|---|---|---|---|
| | $S_a$ | $S_b$ | $S_c$ | $S_a'$ | $S_b'$ | $S_c'$ |
| AIDS | 0.29 | 6.09 | 2.11 | 0.51 | 0.39 | 0.27 |
| S100K | 0.53 | 5.45 | 7.78 | 1.13 | 0.29 | 0.44 |
| Pub-25M | 343.58 | 3669.91 | 2014.12 | 598.1 | 303.83 | 237.21 |

space of $n_v, n_e$ and pointers of all nodes, and $S_b$ and $S_c$ are the storage space of $F_D$ and $F_L$ of all nodes, respectively. Correspondingly, $S_a'$ is the storage space of $n_v$, $n_e$, $l_D$, $r_D$, $l_L$, $r_L$ and pointers of all nodes, shown in Figure 6(a). $S_b'$ is total storage space of $B_D$, $S_D$, $SB_D$, $words_D$ and $flag_D$, shown in Figure 6(b), and $S_c'$ is total storage space of $B_L$, $S_L$, $SB_L$, $words_L$ and $flag_L$, shown in Figure 6(c).

From Table 3, we know that $S_b$ and $S_c$ take up most amount of storage space of $T_Q$, thus a succinct representation of $F_D$ and $F_L$ of all nodes is an efficient way to reduce the storage space of $T_Q$. Compared with $S_b$ and $S_c$, both $S_b'$ and $S_c'$ can be reduced by more than 90%. This is because that (1) only nonzero entries are needed to encode in the succinct representation; (2) our hybrid encoding will greatly reduce the number of bits required for each nonzero entry. Compared with the storage space of $T_Q$ (the sum of $S_a$, $S_b$ and $S_c$), the storage space of $T_{SQ}$ (the sum of $S_a'$, $S_b'$ and $S_c'$) can be reduced by more than 80%. Thus, the succinct representation of $q$-$gram$ tree can greatly reduce the storage space.

### 7.2.2 Comparing with Existing Indexes

We vary the size of datasets to evaluate the index storage space and construction time, and show the results in Figure 7. Regarding the index size, Mixed consumes the most amount of space in AIDS and PubChem-25M since it has to store all branch and disjoint structures. However, GSimJoin does not perform well in S100K.E30.D50.L5 since the number of paths increases exponentially in the dense graphs. MSQ-Index performs the best and its index size is only 5% of that of Mixed and 15% of that of C-Star. This is because that (1) the total number of degree structures and label structures is less than the number of tree structures and paths; (2) the entries in the succinct $q$-$gram$ tree are compressed for efficient storage. For the large dataset PubChem-25M, Mixed, GSimJoin and C-Star cannot properly run for the memory error when the dataset size is more than 15M, while the index size of MSQ-Index is about 1.2GB, which achieves an excellent performance.

By Figure 7, we know that C-Star has the shortest index construction time. This is because that it only needs to enumerate all star structures in each data graph without any complex index. Although MSQ-Index is stored in a succinct form, its construction time is shorter than GSimJoin and Mixed. For the large dataset PubChem-25M, it can be built done in 1 hour.

## 7.3 Filter Performance

In this section, we evaluate the query performance of all tested methods on the datasets AIDS and S100K.E30.D50.L5 and PubChem-25M on two metrics: # of candidates passed the filtering and overall processing time. For overall processing, we further divide it into two parts: indexing processing time and candidate verification time.

We fix the datasets and vary the edit distance threshold $\tau$

from 1 to 5 to evaluate the filter capability and response time. Figure 8 shows the the average candidate size and total response time (i.e., the filtering time plus the verification time) of different methods for the fifty query graphs. Note that, we combine the heuristic estimate function $h(x)$ in [24] into the software provided by Riesen et al. [14] to compute the exact graph edit distance in the verifcation phase for C-Star, Mixed and MSQ-Index, except for GSimJoin has implemented it in their executable binary file.

Regarding the candidate size, we can know that our method has the smallest candidate size in most case. GSimJoin and C-Star do not perform well because that both tree structures and paths have much more overlapping. In S100K.E30.D50.L5, Mixed performs the best when $\tau \leq 3$ and our method has a close candidate size with it. For the large dataset PubChem-25M, only our method can properly run since only it can be built done in our environment.

For the response time of C-Star (denoted by "C"), GSimJoin (denoted by "G"), Mixed (denoted by "M") and MSQ-Index (denoted by "S"), we know that C-Star consumes the longest filtering time because it needs to construct a minimum weighted bipartite graph between each data graph and the query graph. Even though GSimJoin shows a better filtering time in AIDS, it produces a large candidate set than MSQ-Index, making the total response time large than MSQ-Index. Compared with Mixed, MSQ-Index can achieve 1.6x speedup on AIDS and 3.8x speedup on S100K.E30.D50.L5 on the average. In addition, although MSQ-Index are compressed for efficient storage, it can provide good filtering efficiency especially when $\tau$ is small, such as the total filtering time of MSQ-Index is less than 5s in PubChem-25M when $\tau = 1$.

## 7.4 Scalability

In this section, we evaluate the scalability performance of C-Star, GSimJoin, Mixed and MSQ-Index on the real and synthetic datasets.

### 7.4.1 Varying $|V_h|$

We vary the query graph size from 10 to 60, and fix the size of PubChem-25M be 5M and $\tau = 3$, respectively, to evaluate the effect of the query graph size on the query performance. Figure 9 shows the distribution of graphs in $G$, where x-axis is the graph size, i.e., the number of vertices and y-axis is the number of graphs of the same size in the dataset. Figure 10 shows the average candidate size and total filtering time for the fifty query graphs, respectively.
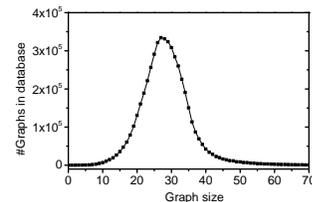


Figure 9: Graph distribution.

By Figure 9, we know that the distribution of data graphs in $G$ is close to a normal distribution and the number of graphs whose size near 30 is relatively large. Thus, the average candidate size of all tested (excepting GSimJoin for
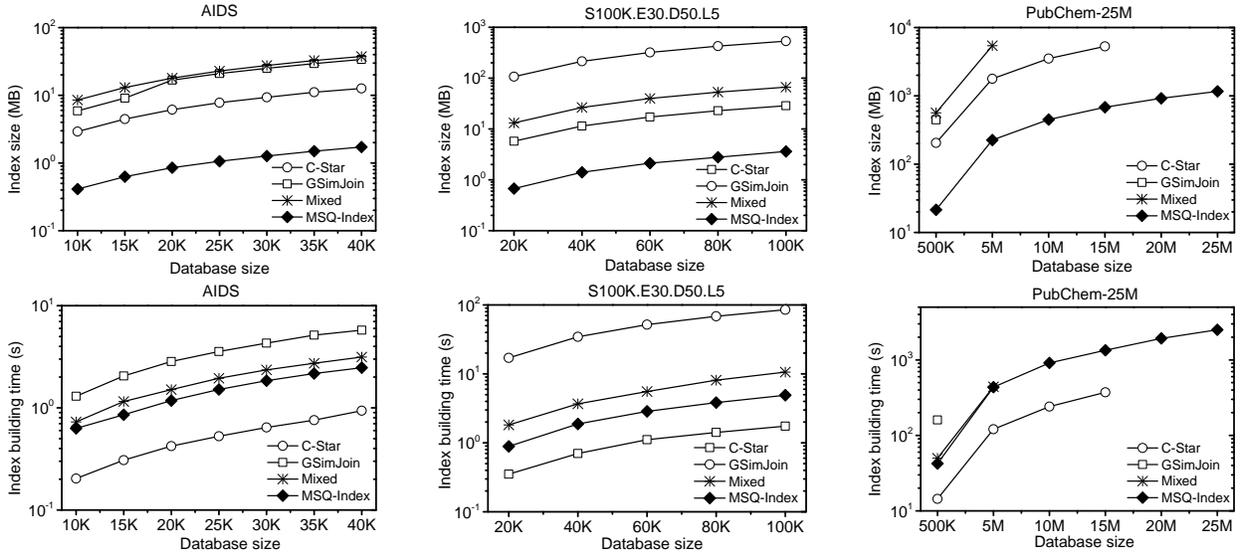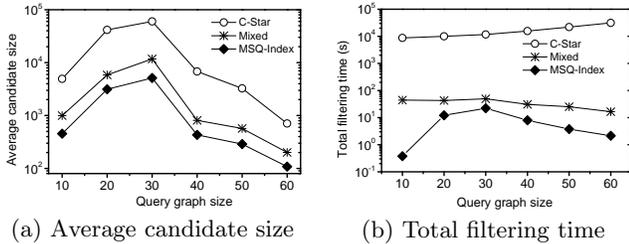
Figure 7: Index size and building time on the three datasets.

the memory error) methods first increase and then decrease, and achieves the maximum when the query graph size is 30.

By Figure 10(b), we know that MSQ-Index has the shortest filtering time. Compared with Mixed, MSQ-Index can achieve 8–40x speedup when the query graph size is less than 20 or more than 50. The reason is that the number of data graphs whose size near 20 or 50 in the dataset is relatively small by Figure 9, resulting in the query region $Q_h$ containing few data graphs.



(a) Average candidate size  (b) Total filtering time

Figure 10: Scalability *vs.* $|V_h|$.

### 7.4.2 Varying $|G|$

We fix $\tau = 5$ and vary the size of PubChem-25M from 500K (kilo) to 25M (million) to evaluate the effect of the dataset size. Figure 11 shows the average candidate size and the total response time for the fifty graphs. Among all tested methods, MSQ-Index has the smallest candidate size and the shortest response time. When the dataset size is 10M, GSimJoin and Mixed cannot properly run for the memory error, and the verification time of C-Star is longer than 48 hours, making all of them not be suitable for such large dataset. Only MSQ-Index can easily scale to cope with it.

### 7.4.3 Varying $|\Sigma_V|$

We fix $\tau = 5$, the dataset size be 100K, the average graph density $\rho = 50\%$, the number of edges in each data graph be 30, respectively, and then produce a group of synthetic



(a) Average candidate size  (b) Total response time

Figure 11: Scalability *vs.* $|G|$.

datasets to evaluate the effect of the number of labels. Figure 12 shows the average candidate size of all tested methods. By Figure 12, we know that the average candidate size decreases as the number of vertex labels increases. This is because that more information can be used to filter.
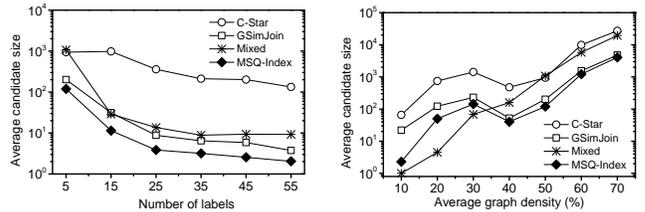


Figure 12: Scalability *vs.* $|\Sigma_V|$  Figure 13: Scalability *vs.* $\rho$.

### 7.4.4 Varying $\rho$

We fix $\tau = 5$, the dataset size be 100K, the number of edges and vertex labels in each data graph be 30 and 5, respectively, and then produce a group of synthetic datasets to evaluate the effect of the average density. Figure 13 shows the average candidate size of all tested methods. It shows that the candidate size of all tested methods increases as $\rho$ increases when $\rho \geq 40\%$. This is because that all tested methods only using the local structures have a weak filter ability for the the density graphs.
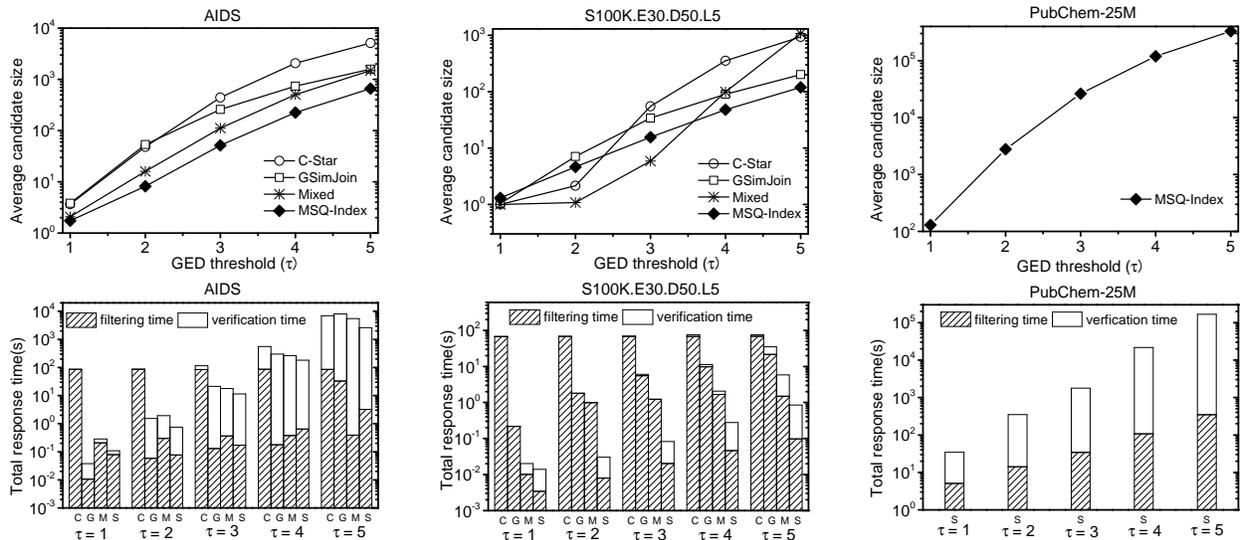
Figure 8: Average candidate size and total response time on the three datasets.

## 8. RELATED WORKS

Recently, graph similarity search has received considerable attention. $\kappa$-AT [16] and GSimJoin [24] are two major *q-gram* counting filters. In $\kappa$-AT, a *q-gram* is defined as a tree consisting of a vertex $v$ and the paths whose length no longer than $\kappa$ starting from $v$. However, GSimJoin considered the simple path whose length is $p$ as a *q-gram*. The principle of the *q-gram* counting filter is stated as follows: if $ged(g, h) \leq \tau$, graphs $g$ and $h$ must share at least $\max\{|Q(g)| - \gamma_g \cdot \tau, |Q(h)| - \gamma_h \cdot \tau\}$ common *q-grams*, where $Q(g)$ and $Q(h)$ denote the multisets of *q-grams* in $g$ and $h$, respectively, $\gamma_g$ and $\gamma_h$ are the maximum number of *q-grams* that can be affected by an edit operation, respectively. C-Star [22] and Mixed [25, 26] are two mapping distance-based filters. The lower bounds are $L_S(g, h) = \frac{s_m(g, h)}{max\{4, max\{d_g, d_h\}+1\}}$ and $L_B(g, h) = \frac{b_m(g, h)}{2}$, respectively, where $s_m(g, h)$ and $b_m(g, h)$ are the mapping distances derived based on the minimum weighted bipartite graphs between the star and branch structures of $g$ and $h$, respectively, and $d_g$ and $d_h$ are the respective maximum degrees in $g$ and $h$. SEGOS [17] introduced a two-level index structure to speed up the filtering process, which has the same filter ability with C-star. Pars [23] divided each data graph $g$ into $\tau + 1$ non-overlapping substructures, and pruned the graph $g$ if there exists no substructure that is subgraph isomorphic to $h$. The above methods show different performance on different databases and we can hardly prove the merits of them theoretically [4].

In the verification phase, $A^*$ algorithm [12] is widely used to compute the exact graph edit distance. Zhao et.al [24] and Gouda et.al [9] designed different heuristic estimate functions to improve $A^*$. Note that, we only focus on the filtering phase in this paper.

When the database contains millions of graphs, many existing approaches cannot properly run. gWT [21] utilized the Weisfeiler-Lehman (WL) kernel [11] function to compute the similarity between two graphs and constructed a wavelet tree [13] to speed up the query processing. Chen et.al [20] built the index structure on a hadoop [8] cluster to search on the large database. Unlike previous methods, we propose the first succinct index structure for this problem for efficient storage. Our index can scale to cope with the large dataset of millions of graphs.

## 9. CONCLUSIONS AND FUTURE WORK

We present an space-efficient index structure for the graph similarity search problem, whose encoded sequence $S_X$ requires $|\Psi_X| \log b_X^m + |\Psi_X|$ bits, where $X$ denotes $D$ or $L$. Our index structure incorporates succinct data structures and hybrid encoding to significantly reduce the index space usage while at the same time keeping fast query performance. Each entry in $\Psi_X$ requires about between 4 and 6 bits on our data, which is much smaller than that used to represent an entry in the compared indexing methods in this paper. However, there is still room for improvement on this space bound of $S_X$. The design of a representation of the *q-gram* tree that achieves the entropy-compressed space bound while still preserving query efficiency is left as a future work.

## 10. ACKNOWLEDGMENTS

## 11. ADDITIONAL AUTHORS

## 12. REFERENCES

[1] D.Justice and A.Hero. A binary linear programming formulation of the graph edit distance. *IEEE Trans. Pattern Anal Mach Intell.*, 28(8):1200–1214, 2006.

[2] M. L. Fernández and G. Valiente. A graph distance metric combining maximum common subgraph and minimum common supergraph. *Pattern Recognit Lett.*, 22(6):753–758, 2001.

[3] H. Fröhlich, J. K. Wegner, and F. Sieker. Optimal assignment kernels for attributed molecular graphs. In *ICML*, pages 225–232, 2005.

[4] K. Gouda and M. Arafa. An improved global lower bound for graph edit similarity search. *Pattern Recognit Lett.*, 58:8–14, 2015.

[5] H.Huo, L.Chen, J.S.Vitter, and Y.Nekrich. A practical implementation of compressed suffix arrays with applications to self-indexing. In *DCC*, pages 292–301, 2014.

[6] H.Shang, K.Zhu, X.Lin, Y.Zhang, and R.Ichise. Similarity search on supergraph containment. In *ICDE*, pages 903–914, 2010.

[7] G. Jacobson. *Succinct data structures*. Carnegie Mellon University, 1989.

[8] J.Dean and S.Ghemawat. Mapreduce: simplified data processing on large clusters. *Commun. ACM.*, 51(1):107–113, 2008.

[9] K.Gouda and M.Hassaan. CSI-GED: An efficient approach for graph edit similarity computation. In *ICDE*, pages 256–275, 2016.

[10] K.Riesen, S.Fankhauser, and H.Bunke. Speeding up graph edit distance computation with a bipartite heuristic. In *MLG*, pages 21–24, 2007.

[11] N.Shervashidze and K.M.Borgwardt. Fast subtree kernels on graphs. In *NIPS*, pages 1660–1668, 2009.

[12] P.E.Hart, N.J.Nilsson, and B.Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE Trans.SSC.*, 4(2):100–107, 1968.

[13] R.Grossi, A.Gupta, and J. S.Vitter. High-order entropy-compressed text indexes. In *SODA*, pages 841–850, 2003.

[14] K. Riesen, S. Emmenegger, and H. Bunke. A novel software toolkit for graph edit distance computation. In *GbRPR*, pages 142–151, 2013.

[15] A. Robles-Kelly and R. H. Edwin. Graph edit distance from spectral seriation. *IEEE Trans. Pattern Anal Mach Intell.*, 27(3):365–378, 2005.

[16] G. Wang, B. Wang, X. Yang, and G. Yu. Efficiently indexing large sparse graphs for similarity search. *IEEE Trans. Knowl Data Eng.*, 24(3):440–451, 2012.

[17] X. Wang, X. Ding, A. K. H. Tung, S. Ying, and H. Jin. An efficient graph indexing method. In *ICDE*, pages 210–221, 2012.

[18] X. Wang, A. Smalter, J. Huan, and H. Gerald. G-hash: towards fast kernel-based similarity search in large graph databases. In *EDBT*, pages 472–480, 2009.

[19] X. Yan, P. S. Yu, and J. Han. Graph indexing: a frequent structure-based approach. In *SIGMOD*, pages 335–346, 2004.

[20] Y.Chen, X.Zhao, G.Bin, C.Xiao, and C.H.Cui. Practising scalable graph similarity joins in mapreduce. In *BigData*, pages 112–119, 2014.

[21] Y.Tabei and K.Tsuda. Kernel-based similarity search in massive graph databases with wavelet trees. In *SDM*, pages 154–163, 2011.

[22] Z. Zeng, A. K. H. Tung, J. Wang, J. Feng, and L. Zhou. Comparing stars: On approximating graph edit distance. *PVLDB*, 2(1):25–36, 2009.

[23] X. Zhao, C. Xiao, X. Lin, Q. Liu, and W. Zhang. A partition-based approach to structure similarity search. *PVLDB*, 7(3):169–180, 2013.

[24] X. Zhao, C. Xiao, X. Lin, and W. Wang. Efficient graph similarity joins with edit distance constraints. In *ICDE*, pages 834–845, 2012.

[25] W. Zheng, L. Zou, X. Lian, D. Wang, and D. Zhao. Graph similarity search with edit distance constraint in large graph databases. In *CIKM*, pages 1595–1600, 2013.

[26] W. Zheng, L. Zou, X. Lian, D. Wang, and D. Zhao. Efficient graph similarity search over large graph databases. *IEEE Trans. Knowl Data Eng.*, 27(4):964–978, 2015.