

# Effective Community Search for Large Attributed Graphs

Yixiang Fang, Reynold Cheng, Siquang Luo, Jiafeng Hu  
Department of Computer Science, The University of Hong Kong, Pokfulam Road, Hong Kong  
{yxfang, ckcheng, sqluo, jhu}@cs.hku.hk

## ABSTRACT

Given a graph  $G$  and a vertex  $q \in G$ , the *community search* query returns a subgraph of  $G$  that contains vertices related to  $q$ . Communities, which are prevalent in *attributed graphs* such as social networks and knowledge bases, can be used in emerging applications such as product advertisement and setting up of social events. In this paper, we investigate the *attributed community query* (or AC-Q), which returns an *attributed community* (AC) for an *attributed graph*. The AC is a subgraph of  $G$ , which satisfies both *structure cohesiveness* (i.e., its vertices are tightly connected) and *keyword cohesiveness* (i.e., its vertices share common keywords). The AC enables a better understanding of how and why a community is formed (e.g., members of an AC have a common interest in music, because they all have the same keyword “music”). An AC can be “personalized”; for example, an ACQ user may specify that an AC returned should be related to some specific keywords like “research” and “sports”.

To enable efficient AC search, we develop the CL-tree index structure and three algorithms based on it. We evaluate our solutions on four large graphs, namely Flickr, DBLP, Tencent, and DBpedia. Our results show that ACs are more effective and efficient than existing community retrieval approaches. Moreover, an AC contains more precise and personalized information than that of existing community search and detection methods.

## 1. INTRODUCTION

Due to the recent developments of gigantic social networks (e.g., Flickr, Facebook, and Twitter), the topic of *attributed graphs* has attracted attention from industry and research communities [29, 3, 6, 14, 16, 33, 17, 10]. An attributed graph is essentially a graph associated with text strings or keywords. Figure 1 illustrates an attributed graph, where each vertex represents a social network user, and its keywords describe the interest of that user.

In this paper, we investigate the *attributed community query* (or ACQ). Given an attributed graph  $G$  and a vertex  $q \in G$ , the ACQ returns one or more subgraphs of  $G$  known as *attributed communities* (or ACs). An AC is a kind of *community*, which consists of vertices that are closely related [26, 5, 4, 15, 22, 11]. Particularly,

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org).

*Proceedings of the VLDB Endowment*, Vol. 9, No. 12  
Copyright 2016 VLDB Endowment 2150-8097/16/08.

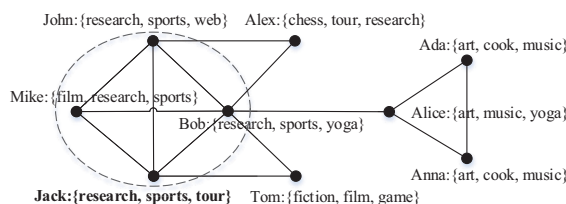


Figure 1: Attributed graph and AC (circled).

Table 1: Classification of works in community retrieval.

Graph Type	Community detection (CD)	Community search (CS)
Non-attributed	[22, 11]	[26, 5, 4, 15, 18]
Attributed	[34, 21, 20, 29, 23]	<b>ACQ (This paper)</b>

an AC satisfies *structure cohesiveness* (i.e., its vertices are closely linked to each other) and *keyword cohesiveness* (i.e., its vertices have keywords in common). Figure 1 illustrates an AC (circled), which is a connected subgraph with vertex degree 3; its vertices {Jack, Bob, John, Mike} have two keywords (i.e., “research” and “sports”) in common.

**Prior works.** The problems related to retrieving communities from a graph can generally be classified into *community detection* (CD) and *community search* (CS). In general, CD algorithms aim to retrieve all communities for a graph [22, 11, 34, 21, 20, 29, 23]. These solutions are not “query-based”, i.e., they are not customized for a query request (e.g., a user-specified query vertex). Moreover, they can take a long time to find all the communities for a large graph, and so they are not suitable for quick or *online* retrieval of communities. To solve these problems, CS solutions have been recently developed [26, 5, 4, 15]. These approaches are query-based, and are able to derive communities in an “online” manner. However, existing CS algorithms assume *non-attributed* graphs, and only use the graph structure information to find communities. The ACQ is a class of CS problem for attributed graphs. As we will show, the use of keyword information can significantly improve the effectiveness of the communities retrieved. Table 1 summarizes some representative existing works in this area.

**Features of ACs.** We now present more details about ACs.

• **Ease of interpretation.** As demonstrated in Figure 1, an AC contains tightly-connected vertices with similar contexts or backgrounds. Thus, an ACQ user can focus on the common keywords or features of these vertices (e.g., the vertices of the AC in this example contain “research” and “sports”, reflecting that all members of this AC like research and sports). We call the set of common

keywords among AC vertices the *AC-label*. In our experiments, the AC-labels facilitate understanding of the vertices that form the AC.

The design of ACs allows it to be used in setting up of social events. For example, if a Twitter member has many keywords about traveling (e.g., he posted a lot of photos about his trips, with keywords), issuing an ACQ with this member as the query vertex may return other members interested in traveling, because their vertices also have keywords related to traveling. A group tour can then be recommended to these members.

• **Personalization.** The user of an ACQ can control the semantics of the AC, by specifying a set of  $S$  of keywords. Intuitively,  $S$  decides the meaning of the AC based on the user’s need. If we let  $q=Jack$  and  $S=\{\text{“research”}\}$ , the AC is formed by  $\{Jack, Bob, John, Mike, Alex\}$ , who are all interested in research. Let us consider another example in the DBLP bibliographical network, where each vertex’s attribute is represented by the top-20 frequent keywords in their publications. Let  $q=Jim\ Gray$ . If  $S$  is the set of keywords  $\{\text{transaction, data, management, system, research}\}$ , we obtain the AC in Figure 2(a), which contains six prominent database researchers closely related to Jim. On the other hand, when  $S$  is  $\{\text{sloan, digital, sky, survey, SDSS}\}$ , the ACQ yields another AC in Figure 2(b), which indicates the seven scientists involved in the SDSS project<sup>1</sup>. Thus, with the use of different keyword sets  $S$ , different “personalized” communities can be obtained.

Existing CS algorithms, which do not handle attributed graphs, may not produce the two ACs above. For example, the CS algorithm in [26] returns the community with *all* the 14 vertices shown in Figures 2(a) and (b). The main reasons are: (1) these vertices are heavily linked with Jim; and (2) the keywords are not considered. In contrast, the use of set  $S$  in the ACQ places these vertices into two communities, containing vertices that are cohesive in terms of *structure* and *keyword*. This allows a user to focus on the important vertices that are related to  $S$ . For example, using the AC of Figure 2(a), a database conference organizer can invite speakers who have a close relationship with Jim.

The personalization feature is also useful in marketing. Suppose that Mary, a yoga lover, is a customer of a gym. An ACQ can be issued on a social network, with Mary as the query vertex and  $S=\{\text{“yoga”}\}$ . Since members of the AC contain the keyword “yoga”, they can be the gym’s advertising targets. On the other hand, current CS algorithms may return a community that contains one or more vertices without the keyword “yoga”. It is not clear whether the corresponding user of this vertex is interested in yoga.

• **Online evaluation.** Similar to other CS solutions, we have developed efficient ACQ algorithms for large graphs, allowing ACs to be generated quickly upon a query request. On the contrary, existing CD algorithms [34, 23, 21, 20] that generate all communities for a graph are often considered to be offline solutions, since they are often costly and time-consuming, especially on very large graphs.

**Technical challenges and our contributions.** We face two important questions: (1) What should be a sound definition of an AC? (2) How to evaluate ACQ efficiently? For the first question, we define an AC based on the *minimum degree*, which is one of the most common structure cohesiveness metrics [22, 11, 26, 5]. This measure requires that every vertex in the community has a degree of  $k$  or more. We formulate the keyword cohesiveness as maximizing the number of shared keywords in keyword set  $S$ . The shared keywords naturally reveal the common features among vertices (e.g., common interest of social network users). We can also use these shared keywords to explain how a community is formed.

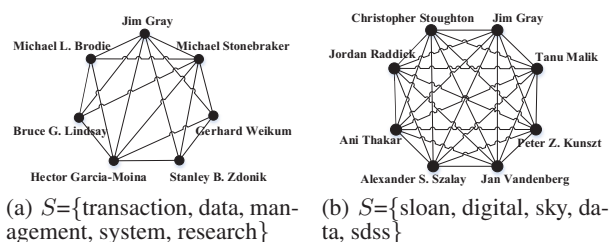


Figure 2: Two ACs of Jim Gray.

The second question is not easy to answer, because the attributed graph  $G$  to be explored can be very large, and the (structure and keyword) cohesiveness criteria can be complex to handle. A simple way is first to consider all the possible keyword combinations, and then return the subgraphs, which satisfy the minimum degree constraint and have the most shared keywords. This solution, which requires the enumeration of all the subsets of  $q$ ’s keyword set, has a complexity exponential to the size  $l$  of  $q$ ’s keyword set. In our experiments, for some queries,  $l$  can be up to 30, resulting in the consideration of  $2^{30} = 1,073,741,824$  subsets of  $q$ . The algorithm is impractical, especially when  $q$ ’s keyword set is large.

We observe the *anti-monotonicity* property, which states that given a set  $S$  of keywords, if it appears in every vertex of an AC, then for every subset  $S'$  of  $S$ , there exists an AC in which every vertex contains  $S'$ . We use this intuition to propose better algorithms. We further develop the *CL-tree*, an index that organizes the vertex keyword data in a hierarchical structure. The CL-tree has a space and construction time complexity linear to the size of  $G$ . We have developed three different ACQ algorithms based on the CL-tree, and they are able to achieve a superior performance.

We have performed extensive experiments on four large real graph datasets (namely Flickr, DBLP, Tencent, and DBpedia). We found that a large number of common keywords appear across vertices in our graph datasets. In DBLP, for instance, an AC with one common keyword contains over 5,000 vertices on average; an AC with two common keywords contains over 700 vertices. Hence, using shared keywords among vertices as keyword cohesiveness makes sense. We have also studied how to quantify the quality of a community, based on occurrence frequencies of keywords and similarity between the keyword sets of two vertices. We conducted a detailed case study on DBLP. These results confirm the superiority of the AC over the communities returned by existing community detection and community search algorithms, in terms of community quality. The performance of our best algorithm is 2 to 3 order-of-magnitude better than solutions that do not use the CL-tree. Another advantage of our approaches is that they organize and search vertex keywords for ACs effectively, achieving a higher efficiency than existing community search solutions (that do not use vertex keywords in the community search process).

**Organization.** We review the related work in Section 2, and define the ACQ problem formally in Section 3. Section 4 presents the basic solutions, and Section 5 discusses the CL-tree index. We present the query algorithms in Section 6. Our experimental results are reported in Section 7. We conclude in Section 8.

## 2. RELATED WORK

**Community detection (CD).** A large class of studies aim to discover or *detect* all the communities from an entire graph. Table 1 summarises these works. Earlier solutions, such as [22, 11], employ link-based analysis to obtain these communities. However,

<sup>1</sup>URL of the SDSS project: <http://www.sdss.org>.

they do not consider the textual information associated with graphs. Recent works focus on attributed graphs, and use clustering techniques to identify communities. For instance, Zhou et al. [34] considered both links and keywords of vertices to compute the vertices' pairwise similarities, and then clustered the graph. Ruan et al. [23] proposed a method called CODICIL. This solution augments the original graphs by creating new edges based on content similarity, and then uses an effective graph sampling to boost the efficiency of clustering. We will compare ACQ with this method experimentally.

Another common approach is based on topic models. In [21, 20], the Link-PLSA-LDA and Topic-Link LDA models jointly model vertices' content and links based on the LDA model. In [29], the attributed graph is clustered based on probabilistic inference. In [24], the topics, interaction types and the social connections are considered for discovering communities. CESNA [31] detects overlapping communities by assuming communities "generate" both the link and content. A discriminative approach [32] has also been considered for community detection. As discussed before, CD algorithms are generally slow, as they often consider the pairwise distance/similarity among vertices. Also, it is not clear how they can be adapted to perform online ACQ. In this paper, we propose online algorithms for finding communities on attributed graphs.

**Community search (CS).** Another class of solutions aims to obtain communities in an "online" manner, based on a query request. For example, given a vertex  $q$ , several existing works [26, 5, 18, 4, 15] have developed fast algorithms to obtain a community for  $q$ . To measure the structure cohesiveness of a community, the *minimum degree* is often used [26, 5, 18]. Sozio et al. [26] proposed the first algorithm Global to find the  $k$ -core containing  $q$ . Cui et al. [5] proposed Local, which uses local expansion techniques to enhance the performance of Global. We will compare these two solutions in our experiments. Other definitions, including  $k$ -clique [4] and  $k$ -truss [15], have also been considered for searching communities. A recent work [18] finds communities with high influence. These works assume non-attributed graphs, and overlook the rich information of vertices that come with attributed graphs. As we will see, performing CS on attributed graphs is better than on non-attributed graphs.

**Graph keyword search.** Given an attributed graph  $G$  and a set  $Q$  of keywords, graph keyword search solutions output a tree structure, whose nodes are vertices of  $G$ , and the union of these vertices' keyword sets is a superset of  $Q$  [3, 6, 16]. Recent work studies the use of a subgraph of  $G$  as the query output [17]. These works are substantially different from the ACQ problem. First, they do not specify query vertices as required by the ACQ problem. Second, the tree or subgraph produced do not guarantee structure cohesiveness. Third, keyword cohesiveness is not ensured; there is no mechanism that enforces query keywords to be shared among the keyword sets of all query output's vertices. Thus, graph keyword search solutions are not designed to find ACs.

### 3. THE ACQ PROBLEM

We now discuss the attributed graph model, the  $k$ -core, and the AC. In the CS and CD literature, most existing works assume that the underlying graph is undirected [26, 18, 29, 23]. We also suppose that an attributed graph  $G(V, E)$  is undirected, with vertex set  $V$  and edge set  $E$ . Each vertex  $v \in V$  is associated with a set of keywords,  $W(v)$ . Let  $n$  and  $m$  be the corresponding sizes of  $V$  and  $E$ . The degree of a vertex  $v$  of  $G$  is denoted by  $deg_G(v)$ . Table 2 lists the symbols used in the paper.

A community is often a subgraph of  $G$  that satisfies *structure cohesiveness* (i.e., the vertices contained in the community are linked

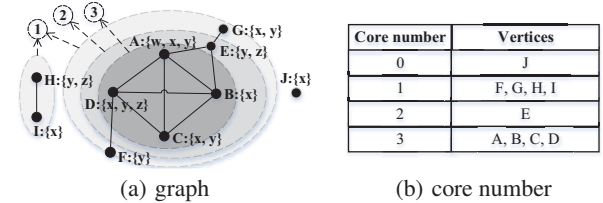
**Table 2: Symbols and meanings.**

Symbol	Meaning
$G(V, E)$	An attributed graph with vertex set $V$ and edge set $E$
$W(v)$	The keyword set of vertex $v$
$deg_G(v)$	The degree of vertex $v$ in $G$
$G[S']$	The largest connected subgraph of $G$ s.t. $q \in G[S']$ , and $\forall v \in G[S'], S' \subseteq W(v)$
$G_k[S']$	The largest connected subgraph of $G$ s.t. $q \in G_k[S']$ , and $\forall v \in G_k[S'], deg_{G_k[S']}(v) \geq k$ and $S' \subseteq W(v)$

to each other in some way). A common notion of structure cohesiveness is that the *minimum degree* of all the vertices that appear in the community has to be  $k$  or more [26, 25, 2, 7, 5, 18]. This is used in the  $k$ -core and the AC. Let us discuss the  $k$ -core first.

**DEFINITION 1 ( $k$ -CORE [25, 2]).** Given an integer  $k$  ( $k \geq 0$ ), the  $k$ -core of  $G$ , denoted by  $H_k$ , is the largest subgraph of  $G$ , such that  $\forall v \in H_k, deg_{H_k}(v) \geq k$ .

We say that  $H_k$  has an order of  $k$ . Notice that  $H_k$  may not be a connected graph [2], and its connected components, denoted by  $k$ -cores, are usually the "communities" returned by  $k$ -core search algorithms.



**Figure 3: Illustrating the  $k$ -core and the AC.**

**EXAMPLE 1.** In Figure 3(a),  $\{A, B, C, D\}$  is both a 3-core and a 3-core. The 1-core has vertices  $\{A, B, C, D, E, F, G, H, I\}$ , and is composed of two 1-core components:  $\{A, B, C, D, E, F, G\}$  and  $\{H, I\}$ . The number  $k$  in each circle represents the  $k$ -core contained in that ellipse.

Observe that  $k$ -cores are "nested" [2]: given two positive integers  $i$  and  $j$ , if  $i < j$ , then  $H_j \subseteq H_i$ . In Figure 3(a),  $H_3$  is contained in  $H_2$ , which is nested within  $H_1$ .

**DEFINITION 2 (CORE NUMBER).** Given a vertex  $v \in V$ , its core number, denoted by  $core_G[v]$ , is the highest order of a  $k$ -core that contains  $v$ .

A list of core numbers and their respective vertices for Example 1 are shown in Figure 3(b). In [2], an  $O(m)$  algorithm was proposed to compute the core number of every vertex.

We now formally define the ACQ problem as follows.

**PROBLEM 1 (ACQ).** Given a graph  $G(V, E)$ , a positive integer  $k$ , a vertex  $q \in V$  and a set of keywords  $S \subseteq W(q)$ , return a set  $\mathcal{G}$  of graphs, such that  $\forall G_q \in \mathcal{G}$ , the following properties hold:

- **Connectivity.**  $G_q \subseteq G$  is connected and contains  $q$ ;
- **Structure cohesiveness.**  $\forall v \in G_q, deg_{G_q}(v) \geq k$ ;
- **Keyword cohesiveness.** The size of  $L(G_q, S)$  is maximal, where  $L(G_q, S) = \cap_{v \in G_q} (W(v) \cap S)$  is the set of keywords shared in  $S$  by all vertices of  $G_q$ .

We call  $G_q$  the *attributed community* (or AC) of  $q$ , and  $L(G_q, S)$  the *AC-label* of  $G_q$ . In Problem 1, the first two properties are also specified by the *k-core* of a given vertex  $q$  [26]. The *keyword cohesiveness* (Property 3), which is unique to Problem 1, enables the retrieval of communities whose vertices have common keywords in  $S$ . We use  $S$  to impose semantics on the AC produced by Problem 1. By default,  $S = W(q)$ , which means that the AC generated should have keywords common to those associated with  $q$ . If  $S \subset W(q)$ , it means that the ACQ user is interested in forming communities that are related to some (but not all) of the keywords of  $q$ . A user interface could be developed to display  $W(q)$  to the user, allowing her to include the desired keywords into  $S$ . For example, in Figure 3(a), if  $q=A$ ,  $k=2$  and  $S=\{w, x, y\}$ , the output of Problem 1 is  $\{A, C, D\}$ , with AC-label  $\{x, y\}$ , meaning that these vertices share the keywords  $x$  and  $y$ .

We require  $L(G_q, S)$  to be maximal in Property 3, because we wish the AC(s) returned only contain(s) the most related vertices, in terms of the number of common keywords. Let us use Figure 3(a) to explain why this is important. Using the same query ( $q=A, k=2, S=\{w, x, y\}$ ), without the “maximal” requirement, we can obtain communities such as  $\{A, B, E\}$  (which do not share any keywords),  $\{A, B, D\}$ , or  $\{A, B, C\}$  (which share 1 keyword). Note that there does not exist an AC with AC-label being exactly  $\{w, x, y\}$ . Our experiments (Section 7) show that imposing the “maximal” constraint yields the best result. Thus, we adopt Property 3 in Problem 1. If there is no AC whose vertices share one or more keywords (i.e.,  $|L(G_q, S)|=0$ ), we return the subgraph of  $G$  that satisfies Properties 1 and 2 only.<sup>2</sup>

There are other candidates for structure cohesiveness (e.g.,  $k$ -truss,  $k$ -clique) and *keyword cohesiveness* (e.g., Jaccard similarity and string edit distance). An AC can also be defined in different ways. For example, an ACQ user may specify that an AC returned must have vertices that contain a specific set of keywords. An interesting direction is to extend ACQ to support for these criteria, and study their effectiveness.

## 4. BASIC SOLUTIONS

For ease of presentation, we say that  $v$  contains a set  $S'$  of keywords, if  $S' \subseteq W(v)$ . We use  $G[S']$  to denote the largest connected subgraph of  $G$ , where each vertex contains  $S'$  and  $q \in G[S']$ . We use  $G_k[S']$  to denote the largest connected subgraph of  $G[S']$ , in which every vertex has degree being at least  $k$  in  $G_k[S']$ . We call  $S'$  a *qualified keyword set* for the query vertex  $q$  on the graph  $G$ , if  $G_k[S']$  exists.

Given a query vertex  $q$ , a straightforward method to answer ACQ performs three steps. First, all non-empty subsets of  $S$ ,  $S_1, S_2, \dots, S_{2^l-1}$  ( $l=|S|$ ), are enumerated. Then, for each subset  $S_i$  ( $1 \leq i \leq 2^l-1$ ), we verify the existence of  $G_k[S_i]$  and compute it when it exists (We postpone to discuss the details). Finally, we output the subgraphs having the most shared keywords among all  $G_k[S_i]$ .

One major drawback of the straightforward method is that we need to compute  $2^l - 1$  subgraphs (i.e.,  $G_k[S_i]$ ). For large values of  $l$ , the computation overhead renders the method impractical, and we do not further consider this method in the paper. To alleviate this issue, we propose the following two-step framework.

### 4.1 Two-Step Framework

The two-step framework is mainly based on the following *anti-monotonicity* property.

<sup>2</sup>In practice, the query user can be alerted by the system when there is no sharing among the vertices.

**LEMMA 1 (ANTI-MONOTONICITY).** *Given a graph  $G$ , a vertex  $q \in G$  and a set  $S$  of keywords, if there exists a subgraph  $G_k[S]$ , then there exists a subgraph  $G_k[S']$  for any subset  $S' \subseteq S$ .*

All the proofs of lemmas studied in this paper can be found in the full version [30]. The anti-monotonicity property allows us to stop examining all the super sets of  $S'$  ( $S' \subseteq S$ ), once we have verified that  $G_k[S']$  does not exist. The basic solution begins with examining the set,  $\Psi_1$ , of size-1 candidate keyword sets, i.e., each candidate contains a single keyword of  $S$ . It then repeatedly executes the following two key steps, to retrieve the size-2 (size-3, ...) qualified keyword subsets until no qualified keyword sets are found.

• **Verification.** For each candidate  $S'$  in  $\Psi_c$  (initially  $c=1$ ), mark  $S'$  as a qualified set if  $G_k[S']$  exists.

• **Candidate generation.** For any two current size- $c$  qualified keyword sets which only differ in one keyword, union them as a new expanded candidate with size- $(c+1)$ , and put it into set  $\Psi_{c+1}$ , if all its subsets are qualified, by Lemma 1.

Among the above steps, the key issue is how to compute  $G_k[S']$ . Since  $G_k[S']$  should satisfy the *structure cohesiveness* (i.e., minimum degree at least  $k$ ) and *keyword cohesiveness* (i.e., every vertex contains keyword set  $S'$ ). Intuitively, we have two approaches to compute  $G_k[S']$ : either searching the subgraph satisfying degree constraint first, followed by further refining with keyword constraints (called `basic-g`); or vice versa (called `basic-w`). These two algorithms form our baseline solutions. Their pseudocodes are presented in the appendix of the full version [30].

## 5. CL-TREE INDEX

The major limitation of `basic-g` and `basic-w` is that they need to find the *k-cores* and do keyword filtering repeatedly. This makes the community search very inefficient. To achieve higher query efficiency, we propose a novel index, called **CL-tree** (Core Label tree), which organizes both the *k-cores* and keywords into a tree structure. Based on the index, the efficiency of answering ACQ and its variants can be improved significantly. We first introduce the index in Section 5.1, and then propose two index construction methods in Section 5.2.

### 5.1 Index Overview

The CL-tree index is built based on the key observation that cores are nested. Specifically, a  $(k+1)$ -core must be contained in a  $k$ -core. The rationale behind is, a subgraph has a minimum degree at least  $k+1$  implies that it has a minimum degree at least  $k$ . Thus, all *k-cores* can be organized into a tree structure<sup>3</sup>. We illustrate this in Example 2.

**EXAMPLE 2.** *Consider the graph in Figure 3(a). All the *k-cores* can be organized into a tree as shown in Figure 4(a). The height of the tree is 4. For each tree node, we attach the core number and vertex set of its corresponding *k-core*.*

From the tree structure in Figure 4(a), we conclude that, if a  $(k+1)$ -core (denoted as  $\mathcal{C}_{k+1}$ ) is contained in a  $k$ -core (denoted as  $\mathcal{C}_k$ ), then there is a tree node corresponding to  $\mathcal{C}_{k+1}$  and its parent node corresponds to  $\mathcal{C}_k$ . Besides, the height of the tree is at most  $k_{max} + 1$ , where  $k_{max}$  is the maximum core number.

The tree structure in Figure 4(a) can be stored compactly, as shown in Figure 4(b). The key observation is that, for any internal node  $p$  in the tree, the vertex sets of its child nodes are the subsets of  $p$ 's vertex set, because of the inclusion relationship. To save space cost, we can remove the redundant vertices that are shared by

<sup>3</sup>We use “node” to mean “CL-tree node” in this paper.

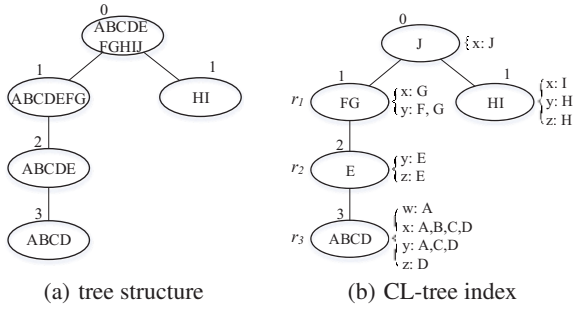


Figure 4: An example CL-tree index.

$p$ 's child nodes from  $p$ 's vertex set. After such removal, we obtain a compressed tree, where each graph vertex appears only once. This structure constitutes the CL-tree index, the nodes of which are further augmented by inverted lists (Figure 4(b)). For each keyword  $e$  that appears in a CL-tree node, a list of IDs of vertices whose keyword sets contain  $e$  is stored. For example, in node  $r_3$ , the inverted list of keyword  $y$  contains  $\{A, C, D\}$ . As discussed later, given a keyword set  $T$ , these inverted lists allow efficient retrieval of vertices whose keyword sets contain  $T$ . To summarize, each CL-tree node contains four elements:

- *coreNum*: the core number of the  $k$ - $\widehat{core}$ ;
- *vertexSet*: a set of graph vertices;
- *invertedList*: a list of  $\langle key, value \rangle$  pairs, where the *key* is a keyword contained by vertices in *vertexSet* and the *value* is the list of vertices in *vertexSet* containing *key*;
- *childList*: a list of child nodes.

Figure 4(b) depicts the CL-tree index for the example graph in Figure 3(a), the elements of each tree node are labeled explicitly. Using the CL-tree, the following two key operations used by our query algorithms (Section 6), can be performed efficiently.

- **Core-locating.** Given a vertex  $q$  and a core number  $c$ , find the  $k$ - $\widehat{core}$  with core number  $c$  containing  $q$ , by traversing the CL-tree.
- **Keyword-checking.** Given a  $k$ - $\widehat{core}$ , find vertices which contain a given keyword set, by intersecting the inverted lists of keywords contained in the keyword set.

**Remarks.** The CL-tree can also support  $k$ - $\widehat{core}$  queries on general graphs without keywords. For example, it can be applied to finding  $k$ - $\widehat{core}$  in previous community search methods [26].

**Space cost.** Since each graph vertex appears only once and each keyword only needs constant space cost, the space cost of keeping such an index is  $O(\widehat{l} \cdot n)$ , where  $\widehat{l}$  denotes the average size of  $W(v)$  over  $V$ . Thus, the space cost is linear to the size of  $G$ .

## 5.2 Index Construction

To build the CL-tree index, we propose two methods, *basic* and *advanced*, as presented in Section 5.2.1 and 5.2.2.

### 5.2.1 The Basic Method

As  $k$ - $\widehat{cores}$  of a graph are nested naturally, it is straightforward to build the CL-tree recursively in a top-down manner. Specifically, we first generate the root node for 0-core, which is exactly the entire graph. Then, for each  $k$ - $\widehat{core}$  of 1-core, we generate a child node for the root node. After that, we only remain vertices with core numbers being 0 in the root node. Then for each child node, we can generate its child nodes in the similar way. This procedure is executed recursively until all the nodes are well built.

Algorithm 1 illustrates the pseudocodes. We first do  $k$ -core decomposition using the linear algorithm [2], and obtain an array

---

### Algorithm 1 Index construction: basic

---

```

1: function BUILDINDEX( $G(V, E)$ )
2:    $core_G[] \leftarrow k$ -core decomposition on  $G$ ;
3:    $k \leftarrow 0, root \leftarrow (k, V)$ ;
4:   BUILDNODE( $root, 0$ );
5:   build an inverted list for each tree node;
6:   return  $root$ ;
7: function BUILDNODE( $root, k$ )
8:    $k \leftarrow k + 1$ ;
9:   if  $k \leq k_{max}$  then
10:    obtain  $U_k$  from  $root$ ;
11:    compute the connected components for the induced
graph on  $U_k$ ;
12:    for each connected component  $C_i$  do
13:      build a tree node  $p_i \leftarrow (k, C_i.vertexSet)$ ;
14:      add  $p_i$  into  $root.childList$ ;
15:      remove  $C_i$ 's vertex set from  $root.vertexSet$ ;
16:    BUILDNODE( $p_i, k$ );

```

---

$core_G[]$  (line 2), where  $core_G[i]$  denotes the core number of vertex  $i$  in  $G$ . We denote the maximal core number by  $k_{max}$ . Then, we initialize the root node by the core number  $k=0$  and  $V$  (line 3). Next, we call the function BUILDNODE to build its child nodes (line 4). Finally, we build an inverted list for each tree node and obtain a well built CL-tree (lines 5-6).

In BUILDNODE, we first update  $k$  and obtain the vertex set  $U_k$  from  $root.vertexSet$ , which is a set of vertices with core numbers being at least  $k$ . Then we find all the connected components from the subgraph induced by  $U_k$  (lines 8-11). Since each connected component  $C_i$  corresponds to a  $k$ - $\widehat{core}$ , we build a tree node  $p_i$  with core number  $k$  and the vertex set of  $C_i$ , and then link it as a child of  $root$  (lines 12-14). We also update  $root$ 's vertex set by removing vertices (line 15), which are shared by  $C_i$ . Finally, we call the BUILDNODE function to build  $p_i$ 's child nodes recursively until all the tree nodes are created (line 16).

**Complexity analysis.** The  $k$ -core decomposition can be done in  $O(m)$ . The inverted lists of each node can be built in  $O(\widehat{l} \cdot n)$ . In function BUILDNODE, we need to compute the connected components with a given vertex set, which costs  $O(m)$  in the worst case. Since the recursive depth is  $k_{max}$ , the total time cost is  $O(m \cdot k_{max} + \widehat{l} \cdot n)$ . Similarly, the space complexity is  $O(m + \widehat{l} \cdot n)$ .

### 5.2.2 The Advanced Method

While the *basic* method is easy to implement, it meets efficiency issues when both the given graph size and its  $k_{max}$  value are large. For instance, when given a clique graph with  $n$  vertices (*i.e.*, edges exist between every pair of nodes), the value of  $k_{max}$  is  $n-1$ . Therefore, the time complexity of the *basic* method could be  $O((m + \widehat{l}) \cdot n)$ , which may lead to low efficiency for large-scale graphs. To enable more efficient index construction, we propose the *advanced* method, whose time and space complexities are almost linear with the size of the input graph.

The *advanced* method builds the CL-tree level by level in a bottom-up manner. Specifically, the tree nodes corresponding to larger core numbers are created prior to those with smaller core numbers. For ease of presentation, we divide the discussion into two main steps: creating tree nodes and creating tree edges.

**1. Creating tree nodes.** We observe that, if we acquire the vertices with core numbers at least  $c$  and denote the induced subgraph on the vertices as  $T_c$ , then the connected components of  $T_c$  have one-to-one correspondence to the  $c$ - $\widehat{cores}$ . A simple algorithm would be, searching connected components for  $T_c$  ( $0 \leq c \leq k_{max}$ ) independently, followed by creating one node for each dis-

tinct component. This algorithm apparently costs  $O(k_{max} \cdot m)$  time, as computing connected components takes linear time.

However, we can do better if we can incrementally update the connected components in a level by level manner (*i.e.*, maintain the connected components of  $T_{c+1}$  from those of  $T_c$ ). We note that, such a node creation process is feasible by exploiting the classical *union-find forest* [1]. Generally speaking, the union-find forest enables efficient maintenance of connected components of a graph when edges are incrementally added. Using union-find forest to maintain connected components follows a process of edge examination. Initially, each vertex is regarded as a connected component. Then, edges are examined one by one. During the examine process, two components are merged together when encounters an edge connecting them. To achieve an efficient merge of components, the vertices in the component form a tree. The tree root acts as the representative vertex of the component. As such, merging two components is essentially linking two root vertices together. To guarantee the CL-tree nodes are formed in a bottom-up manner, we assign an examine priority to each edge. The priority is defined by the larger value of the two core numbers corresponding to the two end vertices of an edge. The edges associated to vertices with larger core numbers are examined first.

**2. Creating tree edges.** Tree edges are also inserted during the graph edge examination process. In particular, when we examine a vertex  $v$  with a set,  $B$ , of its neighbors, whose core numbers are larger than  $core_G[v]$ , we require that the tree node containing  $v$  should link to the tree node containing the vertex, whose core number is the smallest among all the vertices in  $B$ . Nevertheless, the classical union-find forest is not able to maintain such information. To address this issue, we thus propose an auxiliary data structure, called **Anchored Union-Find** (details of AUF are in [30]), based on the classical union-find forest. We first define *anchor vertex*.

**DEFINITION 3 (ANCHOR VERTEX).** *Given a connected subgraph  $G' \subseteq G$ , the anchor vertex is the vertex with core number being  $\min\{core_G[v] | v \in G'\}$ .*

The AUF is an extension of union-find forest, in which each tree has an anchor vertex, and it is attached to the root node. In CL-tree, for any node  $p$  with corresponding  $k$ -core  $C_k$ , its child nodes correspond to the  $k$ -cores, which are contained by  $C_k$  and have core numbers being the most close to the core number of node  $p$ . This implies that, when building the CL-tree in a bottom-up manner, we can maintain the anchor vertices for the  $k$ -cores dynamically, and they can be used to link nodes with their child nodes. In addition, we maintain a vertex-node map, where the key is a vertex and the value is the tree node contains this vertex, for locating tree nodes. The pseudocodes and analysis are reported in the full version [30].

**Complexity analysis.** With our proposed AUF, we can reduce the complexity of CL-tree construction to  $O(m \cdot \alpha(n))$ , where  $\alpha(n)$ , the inverse Ackermann function, is less than 5 for all remotely practical values of  $n$  [1].

**EXAMPLE 3.** *Figure 5 depicts an example graph with 14 vertices  $A, \dots, N$ .  $V_i$  denotes the set of vertices whose core numbers are  $i$ . When  $k=3$ , we first generate two leaf nodes  $p_1$  and  $p_2$ , then update the AUF, where roots' anchor vertices are in the round brackets. When  $k=2$ , we first generate node  $p_3$ , then link it to  $p_1$ , and then update the AUF forest. When  $k=1$ , we first generate nodes  $p_4$  and  $p_5$ . Specifically, to find the child nodes of  $p_4$ , we first find its neighbor  $A$ , then find  $A$ 's parent  $B$  using current AUF forest. Since the anchor vertex of  $B$  is  $E$  and  $E$  points to  $p_3$  in the inverted array, we add  $p_3$  into  $p_4$ 's child List. When  $k=0$ , we generate  $p_6$  and finish the index construction.*

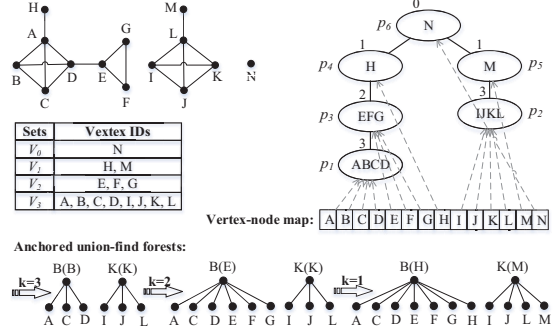


Figure 5: An index built by advanced method.

**Index maintenance.** We now briefly discuss an *incremental* version of the CL-tree construction algorithm, which can handle the changes of keywords and graph edges, without rebuilding the CL-tree from scratch. To insert (or remove) a keyword of a vertex, we just need to update the *invertedList* of the CL-tree node containing the vertex. To insert (or remove) a graph edge, we can borrow the results from [19], which discusses how to maintain a  $k$ -core. We plan to investigate this issue more extensively in the future.

## 6. QUERY ALGORITHMS

In this section, we present three query algorithms based on the CL-tree index. Based on how we verify the candidate keyword sets, we divide our algorithms into incremental algorithms (from examining smaller candidate sets to larger ones) and decremental algorithm (from examining larger candidate sets to smaller ones). We propose two incremental algorithms called **Inc-S (Incremental Space efficient)** and **Inc-T (Incremental Time efficient)**, to trade off between the memory consumption and the computational overhead. The decremental algorithm is called **Dec (Decremental)**. Our interesting finding is that, while **Dec** seems not intuitive, it ranks as the most efficient one. **Inc-S** and **Inc-T** are presented in Section 6.1. **Dec** is introduced in Section 6.2.

### 6.1 The Incremental Algorithms

While the high-level idea of incremental algorithms resembles the basic solutions (see Section 4), **Inc-S** and **Inc-T** advance them with the exploitation of the CL-tree. Specifically, they can always verify the existence of  $G_k[S']$  within a subgraph of  $G$  instead of the entire graph  $G$ . More interestingly, the subgraph for such verifications shrinks when the candidate set  $S'$  expands. Therefore, a large sum of redundant computation is cut off during the verification. We present **Inc-S** and **Inc-T** in Sections 6.1.1 and 6.1.2.

#### 6.1.1 Inc-S Algorithm

We first introduce a new concept, called **subgraph core number**, which is geared to the main idea of **Inc-S**.

**DEFINITION 4 (SUBGRAPH CORE NUMBER).** *The core number of a subgraph  $G'$  of  $G$ ,  $core_G[G']$ , is defined as  $\min\{core_G[v] | v \in G'\}$ .*

**Inc-S** follows the two-step framework (*verification* and *candidate generation*) introduced in Section 4. With the CL-tree, we improve the verification step as follows.

- **Core-based verification.** For each newly generated size- $(c+1)$  candidate keyword set  $S'$  expanded from size- $c$  sets  $S_1$  and  $S_2$ , mark  $S'$  as a qualified set if  $G_k[S']$  exists in a subgraph of core number  $\max\{core_G[G_k[S_1]], core_G[G_k[S_2]]\}$ .

The core-based verification guarantees that, with the expansion of the candidate keyword sets, the verification becomes faster as it only needs to examine the existence of  $G_k[S']$  in a smaller  $k$ - $\widehat{\text{core}}$  (Recall that cores with large core numbers are nested in the cores with small core numbers). The correctness of such shrunk verification range is guaranteed by the following lemma.

LEMMA 2. *Given two subgraphs  $G_k[S_1]$  and  $G_k[S_2]$  of a graph  $G$ , for a new keyword set  $S'$  generated from  $S_1$  and  $S_2$  (i.e.,  $S' = S_1 \cup S_2$ ), if  $G_k[S']$  exists, then it must appear in a  $k$ - $\widehat{\text{core}}$  with core number at least*

$$\max\{\text{core}_G[G_k[S_1]], \text{core}_G[G_k[S_2]]\}. \quad (1)$$

The verification process can be further accelerated by checking the numbers of vertices and edges, as indicated by Lemma 3.

LEMMA 3. *Given a connected graph  $G(V, E)$  with  $n=|V|$  and  $m=|E|$ , if  $m - n < \frac{k^2-k}{2} - 1$ , there is no  $k$ - $\widehat{\text{core}}$  in  $G$ .*

This lemma implies that, for a connected subgraph  $G'$ , whose edge and vertex numbers are  $m$  and  $n$ , if  $m - n < \frac{k^2-k}{2} - 1$ , then we cannot find  $G_k[S']$  from  $G'$ .

We present `INC-S` in Algorithm 2. The input is a CL-tree rooted at *root*, a query vertex  $q$ , a positive integer  $k$  and a keyword set  $S$ . We apply `core-locating` on the CL-tree to locate the internal nodes whose corresponding  $k$ - $\widehat{\text{cores}}$  contain  $q$  (line 2). Note that their core numbers are in the range of  $[k, \text{core}_G[q]]$ , as required by the structure cohesiveness. Then, we set  $l=0$ , indicating the sizes of current keyword sets, and initialize a set  $\Psi$  of  $\langle S', c \rangle$  pairs, where  $S'$  is a set containing a keyword from  $S$  and  $c$  is the initial core number  $k$  (line 3). Note that we skip those keywords, which are in  $S$ , but not in  $W(q)$ . In the while loop (lines 4-18), for each  $\langle S', c \rangle$  pair, we first perform `keyword-checking` to find  $G[S']$  using the keyword inverted lists of the subtree rooted at node  $r_c$ . If we cannot ensure that  $G[S']$  does not contain a  $k$ - $\widehat{\text{core}}$  by Lemma 3, we then find  $G_k[S']$  from  $G[S']$  (lines 8-9). If  $G_k[S']$  exists, we put  $S'$  with its core number into the set  $\Phi_l$  (lines 10-11). Next, if  $\Phi_l$  is nonempty, we generate new candidates by calling `GENECAND`( $\Phi_l$ ), which is detailed in the full version [30]. For each candidate  $S'$  in  $\Psi$ , we compute the core number using Lemma 2 and update it as a pair in  $\Psi$  (lines 12-17); otherwise, we stop (line 18). Finally, we output the communities of the latest verified keyword sets (line 19).

---

#### Algorithm 2 Query algorithm: `INC-S`

---

```

1: function QUERY( $G, \text{root}, q, k, S$ )
2:   find subtree root nodes  $r_k, r_{k+1}, \dots, r_{\text{core}_G[q]}$ ;
3:   initialize  $l=0, \Psi$  using  $S$ ;
4:   while true do
5:      $l \leftarrow l + 1; \Phi_l \leftarrow \emptyset$ ;
6:     for each  $\langle S', c \rangle \in \Psi$  do
7:       find  $G[S']$  under the root  $r_c$ ;
8:       if  $G[S']$  is not pruned by Lemma 3 then
9:         find  $G_k[S']$  from  $G[S']$ ;
10:      if  $G_k[S']$  exists then
11:         $\Phi_l.\text{add}(\langle S', \text{core}_G[G_k[S']] \rangle)$ ;
12:      if  $\Phi_l \neq \emptyset$  then
13:         $\Psi \leftarrow \text{GENECAND}(\Phi_l)$ ;
14:        for each  $S'$  in  $\Psi$  do
15:          if  $S'$  is generated from  $S_1$  and  $S_2$  then
16:             $c \leftarrow \max\{\text{core}_G[G_k[S_1]], \text{core}_G[G_k[S_2]]\}$ ;
17:             $\Psi.\text{update}(S', \langle S', c \rangle)$ ;
18:          else break;
19:   output the communities of keyword sets in  $\Phi_{l-1}$ ;

```

---

EXAMPLE 4. *Consider the graph in Figure 3(a) and its index in Figure 4(b). Let  $q=A, k=1$  and  $S=\{w, x, y\}$ . By Algorithm 2, we first find 3 root nodes  $r_1, r_2$  and  $r_3$ . In the first while loop, we find 2 qualified keyword sets  $\{x\}$  and  $\{y\}$  with core numbers being 3 and 1. By Lemma 2, we only need to verify the new candidate keyword set  $\{x, y\}$  under node  $r_3$ .*

#### 6.1.2 `INC-T` Algorithm

We begin with a lemma which inspires the design of `INC-T`.

LEMMA 4. *Given two keyword sets  $S_1$  and  $S_2$ , if  $G_k[S_1]$  and  $G_k[S_2]$  exist, we have*

$$G_k[S_1 \cup S_2] \subseteq G_k[S_1] \cap G_k[S_2]. \quad (2)$$

This lemma implies, if  $S'$  is generated from  $S_1$  and  $S_2$ , we can find  $G_k[S']$  from  $G_k[S_1] \cap G_k[S_2]$  directly. Since every vertex in  $G_k[S_1] \cap G_k[S_2]$  contains both  $S_1$  and  $S_2$ , we do not need to consider the keyword constraint again when finding  $G_k[S']$ .

Based on Lemma 4, we introduce a new algorithm `INC-T`. Different from `INC-S`, `INC-T` maintains  $G_k[S']$  rather than  $\text{core}_G[G_k[S']]$  for each qualified keyword set  $S'$ . As we will demonstrate later, `INC-T` is more effective for shrinking the subgraphs containing the ACs, and thus more efficient. As a trade-off for better efficiency, `INC-T` consumes more memory as it needs to store a list of subgraph  $G_k[S']$  in memory.

Algorithm 3 presents `INC-T`. We first apply `core-locating` to find the  $k$ - $\widehat{\text{core}}$  containing  $q$  from the CL-tree (line 2). Then, we set  $l=0$ , indicating the sizes of current keyword sets, and initialize a set  $\Psi$  of  $\langle S', \widehat{G} \rangle$  pairs, where  $S'$  is a set containing a keyword from  $S$  and  $\widehat{G}$  is the  $k$ - $\widehat{\text{core}}$ . The while loop (lines 4-18) is similar with that of `INC-S`. The main differences are that: (1) for each qualified keyword set  $S'$ , `INC-T` keeps  $G_k[S']$  in memory (line 11); and (2) for each candidate keyword set  $S'$  generated from  $S_1$  and  $S_2$ , `INC-T` finds  $G_k[S']$  from  $G_k[S_1] \cap G_k[S_2]$  directly without further keyword verification (lines 6-9, 16).

---

#### Algorithm 3 Query algorithm: `INC-T`

---

```

1: function QUERY( $G, \text{root}, q, k, S$ )
2:   find the  $k$ - $\widehat{\text{core}}$ , which contains  $q$ ;
3:   initialize  $l=0, \Psi$  using  $S$ ;
4:   while true do
5:      $l \leftarrow l + 1; \Phi_l \leftarrow \emptyset$ ;
6:     for each  $\langle S', \widehat{G} \rangle \in \Psi$  do
7:       find  $G[S']$  from  $\widehat{G}$ ;
8:       if  $G[S']$  is not pruned by Lemma 3 then
9:         find  $G_k[S']$  from  $G[S']$ ;
10:      if  $G_k[S']$  exists then
11:         $\Phi_l.\text{add}(\langle S', G_k[S'] \rangle)$ ;
12:      if  $\Phi_l \neq \emptyset$  then
13:         $\Psi \leftarrow \text{GENECAND}(\Phi_l)$ ;
14:        for each  $S' \in \Psi$  do
15:          if  $S'$  is generated from  $S_1$  and  $S_2$  then
16:             $\widehat{G} \leftarrow G_k[S_1] \cap G_k[S_2]$ ;
17:             $\Psi_l.\text{update}(S', \langle S', \widehat{G} \rangle)$ ;
18:          else break;
19:   output the communities of keyword sets in  $\Phi_{l-1}$ ;

```

---

EXAMPLE 5. *Continue the graph and query ( $q=A, k=1, S=\{w, x, y\}$ ) in Example 4. By `INC-T`, we first find  $G_1[\{x\}]$  and  $G_1[\{y\}]$ , whose vertex sets are  $\{A, B, C, D\}$  and  $\{A, C, D, E, F, G\}$ . Then to find  $G_1[\{x, y\}]$ , we only need to search it from the subgraph, induced by the vertex set  $\{A, C, D\}$ .*

## 6.2 The Decremental Algorithm

The decremental algorithm, denoted by `Dec`, differs from previous query algorithms not only on the generation of candidate keyword sets, but also on the verification of candidate keyword sets.

**1. Generation of candidate keyword sets.** `Dec` exploits the key observation that, if  $S'$  ( $S' \subseteq S$ ) is a qualified keyword set, then there are at least  $k$  of  $q$ 's neighbors containing set  $S'$ . This is because every vertex in  $G_k[S']$  must have degree at least  $k$ . This observation implies, we can generate all the candidate keyword sets directly by using the query vertex  $q$  and  $q$ 's neighbors, without touching other vertices.

Specifically, we consider  $q$  and  $q$ 's neighbor vertices. For each vertex  $v$ , we only select the keywords, which are contained by  $S$  and at least  $k$  of its neighbors. Then we use these selected keywords to form an itemset, in which each item is a keyword. After this step, we obtain a list of itemsets. Then we apply the well-studied frequent pattern mining algorithms (e.g., Apriori [12] and FP-Growth [13]) to find the frequent keyword combinations, each of which is a candidate keyword set. Since our goal is to generate keyword combinations shared by at least  $k$  neighbors, we set the minimum support as  $k$ . In this paper, we use the well-known FP-Growth algorithm [13].

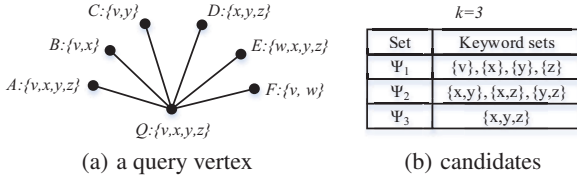


Figure 6: An example of candidate generation.

**EXAMPLE 6.** Consider a query vertex  $Q$  ( $k=3, S=\{v, x, y, z\}$ ) with 6 neighbors in Figure 6(a), where the selected keywords of each vertex are listed in the curly braces. By FP-Growth, 8 candidate keyword sets are generated, as shown in Figure 6(b).  $\Psi_i$  denotes the set of keyword sets with sizes being  $i$ .

**2. Verification of candidate keyword sets.** As candidates can be obtained using  $S$  and  $q$ 's neighbors directly, we can verify them either incrementally as that in `Inc-S`, or in a decremental manner (larger candidate keyword sets first and smaller candidate keyword sets later). In this paper, we choose the latter manner. The rationale behind is that, for any two keyword sets  $S_1 \subseteq S_2$ , the number of vertices containing  $S_2$  is usually smaller than that of  $S_1$ , which implies  $S_2$  can be verified more efficiently than  $S_1$ .

Based on the above discussions, we design `Dec` as shown in Algorithm 4. We first generate candidate keyword sets using  $S$  and  $q$ 's neighbors by FP-Growth algorithm (line 2). Then, we apply `core-locating` to find the root (with core number  $k$ ) of the subtree from CL-tree, whose corresponding  $k$ -core contains  $q$  (line 3). Next, we traverse the subtree rooted at  $r_k$  and find vertices which share keywords with  $q$  (line 4).  $R_i$  denote the sets of vertices sharing  $i$  keywords with  $q$ . Then, we initialize  $l$  as  $h$  (line 5), as we verify keyword sets with the largest size  $h$  first. We maintain a set  $\hat{R}$  dynamically, which contains vertices sharing at least  $l$  keywords with  $q$  (line 6). In the while loop, we examine candidate keyword sets in the decremental manner. For each candidate  $S' \in \Psi_l$ , we first try to find  $G[S']$ , then find  $G_k[S']$ , and put  $G_k[S']$  into  $Q$  if it exists (lines 8-11). Finally, if we have found at least one qualified community, we stop at the end of this loop and output  $Q$ ; otherwise, we examine smaller candidate keyword sets in next loop.

Algorithm 4 Query algorithm: `Dec`

```

1: function QUERY( $G, root, q, k, S$ )
2:   generate  $\Psi_1, \Psi_2, \dots, \Psi_h$  using  $S$  and  $q$ 's neighbors;
3:   find the subtree root node  $r_k$ ;
4:   create  $R_1, R_2, \dots, R_h$  by using subtree rooted at  $r_k$ ;
5:    $l \leftarrow h; Q \leftarrow \emptyset$ ;
6:    $\hat{R} \leftarrow R_h \cup \dots \cup R_{h'}$ ;
7:   while  $l \geq 1$  do
8:     for each  $S' \in \Psi_l$  do
9:       find  $G[S']$  from the subgraph induced on  $\hat{R}$ ;
10:      find  $G_k[S']$  from  $G[S']$ ;
11:      if  $G_k[S']$  exists then  $Q.add(G_k[S'])$ ;
12:      if  $Q=\emptyset$  then
13:         $l \leftarrow l - 1$ ;
14:         $\hat{R} \leftarrow \hat{R} \cup R_l$ ;
15:      else break;
16:   output communities in  $Q$ ;

```

## 7. EXPERIMENTS

We now present the experimental results. Section 7.1 discusses the setup. We discuss the results in Sections 7.2 and 7.3.

### 7.1 Setup

We consider four real datasets. For *Flickr*<sup>4</sup> [27], a vertex represents a user, and an edge denotes a “follow” relationship between two users. For each vertex, we use the 30 most frequent tags of its associated photos as its keywords. For *DBLP*<sup>5</sup>, a vertex denotes an author, and an edge is a co-authorship relationship between two authors. For each author, we use the 20 most frequent keywords from the titles of her publications as her keywords. In the *Tencent* graph provided by the KDD contest 2012<sup>6</sup>, a vertex is a person, an organization, or a microblog group. Each edge denotes the friendship between two users. The keyword set of each vertex is extracted from a user’s profile. For the *DBpedia*<sup>7</sup>, each vertex is an entity, and each edge is the relationship between two entities. The keywords of each entity are extracted by the Stanford Analyzer and Lemmatizer. Table 3 shows the number of vertices and edges, the  $k_{max}$  value, a vertex’s average degree  $\hat{d}$ , and its keyword set size  $\hat{l}$ .

Table 3: Datasets used in our experiments.

Dataset	Vertices	Edges	$k_{max}$	$\hat{d}$	$\hat{l}$
Flickr	581,099	9,944,548	152	17.11	9.90
DBLP	977,288	3,432,273	118	7.02	11.8
Tencent	2,320,895	50,133,369	405	43.2	6.96
DBpedia	8,099,955	71,527,515	95	17.66	15.03

To evaluate ACQs, we set the default value of  $k$  to 6. The input keyword set  $S$  is set to the whole set of keywords contained by the query vertex. For each dataset, we randomly select 300 query vertices with core numbers of 6 or more, which ensures that there is a  $k$ -core containing each query vertex. Each data point is the average result for these 300 queries. We implement all the algorithms in Java, and run experiments on a machine having a quad-core Intel i7-3770 processor, and 32GB of memory, with Ubuntu installed.

### 7.2 Results on Effectiveness

<sup>4</sup><https://www.flickr.com/>

<sup>5</sup><http://dblp.uni-trier.de/xml/>

<sup>6</sup><http://www.kddcup2012.org/c/kddcup2012-track1>

<sup>7</sup><http://dbpedia.org/datasets>



We now study the effectiveness of ACQ, and compare it with existing CD and CS methods. We then discuss a case study.

### 7.2.1 ACQ Effectiveness

We first define two measures, namely CMF and CPJ, for evaluating the keyword cohesiveness of the communities. Let  $C(q)=\{C_1, C_2, \dots, C_{\mathcal{L}}\}$  be the set of  $\mathcal{L}$  communities returned by an algorithm for a query vertex  $q \in V$  (Note that  $S=W(q)$ ).

- **Community member frequency (CMF)**: this is inspired by the classical document frequency measure. Consider a keyword  $x$  of  $q$ 's keyword set  $W(q)$ . If  $x$  appears in most of the vertices (or members) of a community  $C_i$ , then we regard  $C_i$  to be highly cohesive. The CMF uses the occurrence frequencies of  $q$ 's keywords in  $C_i$  to determine the degree of cohesiveness. Let  $f_{i,h}$  be the number of vertices of  $C_i$  whose keyword sets contain the  $h$ -th keyword of  $W(q)$ . Then,  $\frac{f_{i,h}}{|C_i|}$  is the relative occurrence frequency of this keyword in  $C_i$ . The CMF is the average of this value over all keywords in  $W(q)$ , and all communities in  $C(q)$ :

$$CMF(C(q)) = \frac{1}{\mathcal{L} \cdot |W(q)|} \sum_{i=1}^{\mathcal{L}} \sum_{h=1}^{|W(q)|} \frac{f_{i,h}}{|C_i|} \quad (3)$$

Notice that  $CMF(C(q))$  ranges from 0 to 1. The higher its value, the more cohesive is a community.

- **Community pair-wise Jaccard (CPJ)**: this is based on the similarity between the keyword sets of any pair of vertices of community  $C_i$ . We adopt the Jaccard similarity, which is commonly used in the IR literature. Let  $C_{i,j}$  be the  $j$ -th vertex of  $C_i$ . The CPJ is then the average similarity over all pairs of vertices of  $C_i$ , and all communities of  $C(q)$ :

$$CPJ(C(q)) = \frac{1}{\mathcal{L}} \sum_{i=1}^{\mathcal{L}} \left[ \frac{1}{|C_i|^2} \sum_{j=1}^{|C_i|} \sum_{k=1}^{|C_i|} \frac{|W(C_{i,j}) \cap W(C_{i,k})|}{|W(C_{i,j}) \cup W(C_{i,k})|} \right] \quad (4)$$

The  $CPJ(C(q))$  value has a range of 0 and 1. A higher value of  $CPJ(C(q))$  implies better cohesiveness.

**1. Effect of common keywords.** We examine the impact of the AC-label length (i.e., the number of keywords shared by all the vertices of the AC) on keyword cohesiveness. We collect ACs containing one to five keywords, and then group the ACs according to their AC-label lengths. The average CMF and CPJ value of each group is shown in Figure 7. For all the datasets, when the AC-label lengths increase, both CMJ and CPJ value rises. This justifies the use of the maximal AC-label length as the criterion of returning an AC in our ACQ Problem.

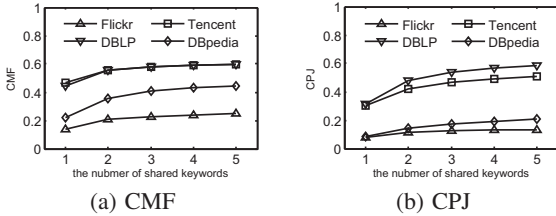


Figure 7: AC-label length.

**2. Comparison with existing CD methods** As mentioned ahead, the existing CD methods for attributed graph can be adapted for community search. We choose to adapt CODICIL [23] for comparison. The main reasons are: (1) it has been tested on the ever reported largest attributed graph (vertex number:3.6M); (2) it

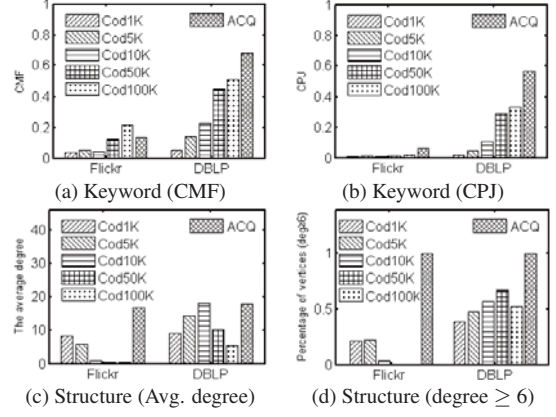


Figure 8: Comparing with community detection method.

identifies communities of comparable or superior quality than those of many existing methods like [21, 32]; and (3) it runs faster than many existing methods. Since CODICIL needs users to specify the number of clusters expected, we set the numbers as 1K, 5K, 10K, 50K and 100K. The corresponding adapted algorithms are named as Cod1K, ..., Cod100K respectively. Other parameter settings are the same as those in [23]. We first run these algorithms offline to obtain all the communities. Given a query vertex  $q$ , they return communities containing  $q$  as the results.

We consider both keyword and structure for evaluating community quality. (1) **Keyword**: Figures 8(a) and (b) show that ACQ (implemented by Dec) always performs the best, in terms of CMF and CPJ. (2) **Structure**: As CODICIL has no guarantee on vertices' minimum degrees, it is unfair to compare them using this metric. We intuitively compare their structure cohesiveness by reporting the average degree of the vertices in the communities and the percentage of vertices having degrees of 6 or more. When the number of clusters in CODICIL is too large or too small, the structure cohesiveness becomes weak, as shown in Figures 8(c) and (d). ACQ always performs better than CODICIL, even when its number of cluster is well set (e.g., Cod10K and Cod50K on DBLP dataset).

**3. Comparison with existing CS methods.** The existing methods mainly focus on non-attributed graphs. We implement two state-of-the-art methods: Global [26] and Local [5]. Both of them use the metric minimum degree, we thus focus on the keyword cohesiveness. Figure 9 shows the CMF and CPJ values for the four datasets. We can see that the keyword cohesiveness of ACQ is superior to both Global and Local, because ACQ considers vertex keywords, while Global and Local do not.

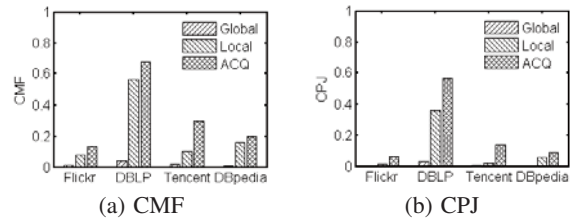


Figure 9: Comparing with community search methods.

### 7.2.2 A Case Study

We next perform a case study on the DBLP dataset, in which we consider two renowned researchers in database and data mining:

Jim Gray and Jiawei Han. We use  $k = 4$  here. We use Cod50K to represent CODICIL for further analysis. We mainly consider the input query keyword set  $S$ , keywords and sizes of communities.

**1. Effect of  $S$ .** Figure 10 shows two ACs of Jiawei (AC-labels are shown in the captions), where the query keyword set  $S$  are set as  $\{\text{analysis, mine, data, information, network}\}$  and  $\{\text{mine, data, pattern, database}\}$  respectively. These two groups of Jiawei’s collaborators are involved in graph analysis (Figure 10(a)) and pattern mining (Figure 10(b)). Although these researchers all have close co-author relationship with Jiawei, the use of the input keyword set  $S$  enables the identification of communities with different research themes. For Jim, we can obtain similar results as discussed in Section 1 (Figure 2). While for CODICIL, it is not clear how to consider the keyword set  $S$ , and we thus do not show the results.

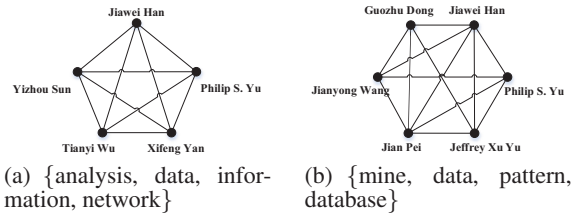


Figure 10: Jiawei Han’s ACs.

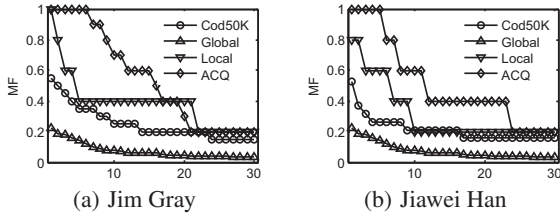


Figure 11: Frequency distribution of keywords.

Table 4: # distinct keywords of communities.

Researcher	Cod50K	Global	Local	ACQ
Jim Gray	134	139,881	60	44
Jiawei Han	140	139,881	58	54

**2. Keyword analysis.** We analyze the frequency distribution of keywords in their communities. Specifically, given a keyword  $w_h$ , we define the member frequency (MF) of  $w_h$  as:  $MF(w_h, C(q)) = \frac{1}{L} \sum_{i=1}^L \frac{f_{i,h}}{|C_i|}$ . The MF measures the occurrence of a keyword in  $C(q)$ . For each  $C_q$  generated by an algorithm, we select 30 keywords with the highest MF values. We report the MF of each keyword in descending order of their MF values in Figure 11. We see that ACQ has the highest MF values for the top 20 keywords. Thus, the keywords associated with the communities generated by ACQ tend to repeat among the community members.

The number of distinct keywords of ACQ communities is also the fewest, as shown in Table 4. For example, the  $k$ -core returned by Global has over 139K distinct keywords, about 2,300 times more than that returned by ACQ (less than 60 keywords). While the semantics of the  $k$ -core can be difficult to understand, the small number of distinct keywords of AC makes it easier to understand why the community is so formed. We further report the keywords

with the 6 highest MF values in Jiawei’s communities in Table 5. We can see that, the top-6 keywords of ACQ are highly related to the input query keyword set, while keywords of Global and Local tend to be less related to the query keyword set, and thus they cannot be used to characterize the communities specifically related to Jiawei. The overall results show that, ACQ performs better than other methods.

Table 5: Top-6 keywords (Jiawei Han).

Algo.	Keywords
Cod50K	information, mine, data, cube, text, network
Global	use, system, model, network, analysis, data
Local	scalable, topical, text, phrase, corpus, mine
ACQ	mine, analysis, data, information, network, heterog

**3. Effect of  $k$  on community size.** We vary the value of  $k$  and report the average size of communities in Figure 12. We can see that the communities returned by Global are extremely large (more than  $10^5$ ), which can make them difficult for a query user to analyze. The community size of Local increases sharply when  $k=8$ . In this situation, Local returns the same community as Global. The size of an AC is relatively insensitive to the change of  $k$ , as AC contains around a hundred vertices for a wide range of values of  $k$ .

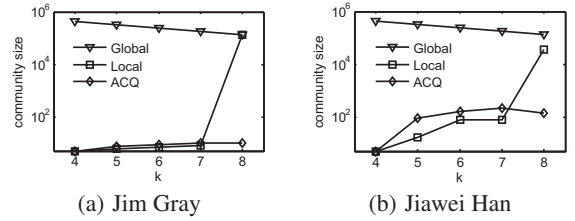


Figure 12: Community size.

### 7.3 Results on Efficiency

For each dataset, we randomly select 20%, 40%, 60% and 80% of its vertices, and obtain four subgraphs induced by these vertex sets. For each vertex, we randomly select 20%, 40%, 60% and 80% of its keywords, and obtain four keyword sets.

**1. Index construction.** Figures 13(a)-13(d) compare the efficiency of Basic and Advanced. We study their main parts, which build the tree without considering keywords. We denote them by Basic- and Advanced-. Notice that Advanced performs consistently faster, and scales better, than Basic. When the subgraph size increases, the performance gap between Advanced and Basic is enlarged. Similar results can be observed between Advanced- and Basic-. In addition, we also run the CD method CODICIL, which takes 32 mins, 2 mins, 1 day, and 3+ days (we stop it after running 3 days) to cluster the vertices of Flickr, DBLP, Tencent and DBpedia offline respectively.

**2. Efficiency of CS methods.** Figures 14(a)-14(d) compares our best algorithm Dec with existing CS methods. We see that Local performs faster than Global for most cases. Also, Dec, which uses the CL-tree index, is the fastest.

**3. Effect of  $k$ .** Figures 14(e)-14(h) compare the query efficiency under different  $k$ . A lower  $k$  renders a larger subgraph, so as the time costs, for all the algorithms. Note that basic-g performs faster than basic-w, but are slower than index-based algorithms. Inc-T performs better than Inc-S, and Dec performs the best. The performance gaps decrease as  $k$  increases.

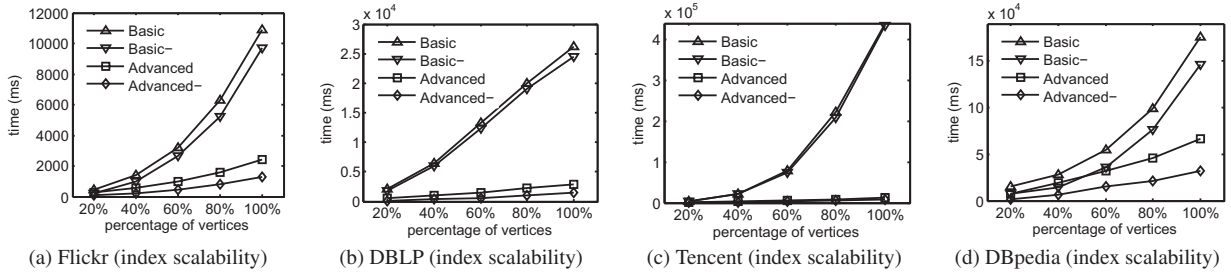


Figure 13: Efficiency results of index construction.

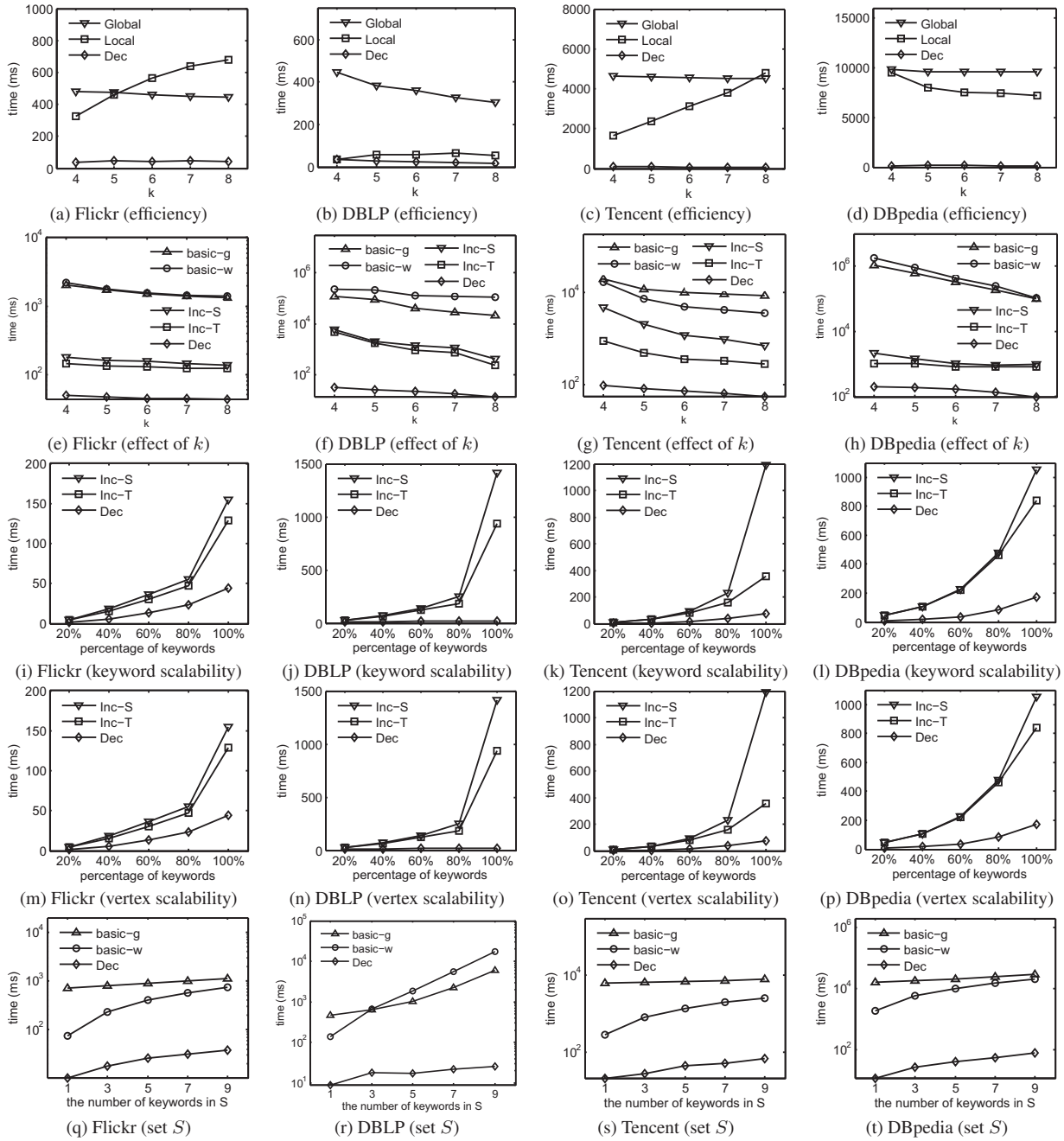


Figure 14: Efficiency results of community search.

**4. ACQ scalability w.r.t. keyword.** Figures 14(i)-14(l) examine scalability over the fraction of keywords for each vertex. All the vertices are considered. The running times of the algorithms increase as more keywords are involved. `Dec` performs the best.

**5. ACQ scalability w.r.t. vertex.** Figures 14(m)-14(p) report the scalability over different fraction of vertices. All the keywords of each vertex are considered. Again, `Dec` scales the best.

**6. Effect of size of  $S$ .** For each query vertex, we randomly select 1, 3, 5, 7 and 9 keywords to form the query keyword set  $S$ . As `Dec` performs better than `Inc-S` and `Inc-T`, we mainly compare `Dec` with the baseline solutions. Figures 14(q)-14(t) show that the cost of all algorithms increase with the  $|S|$ . Also, `Dec` is 1 to 3 order-of-magnitude faster than `basic-g` and `basic-w`.

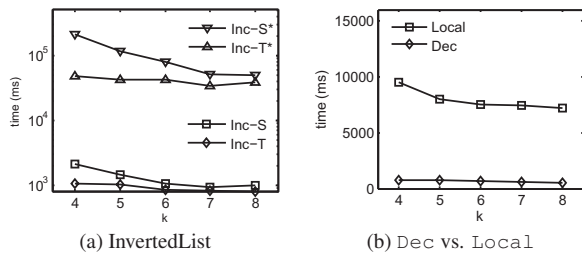


Figure 15: More experimental results on DBpedia.

Next, we present additional results about DBpedia, the largest dataset used in our experiments. Results for other datasets are similar, and are reported in our report [30] due to space constraints.

**7. Effect of invertedList.** We implement `Inc-S*` and `Inc-T*`, which are respective variants of `Inc-S` and `Inc-T`, without the invertedList on each CL-tree node. Figure 15(a) shows that `Inc-S` (`Inc-T`) is 2 orders of magnitude faster than `Inc-S*` (`Inc-T*`). The keyword-checking operation, which uses the invertedList, is frequently performed. The use of invertedList thus improves the performance of our algorithms significantly.

**8. Non-attributed graphs.** We also test `Dec` and `Local` on non-attributed graphs, by ignoring the keywords of the graph dataset. Figure 15(b) shows that `Dec` is always faster than `Local`. In `Dec`, the cores are organized into the CL-tree, and since its height is limited, the core-locating operation is efficient.

## 8. CONCLUSIONS

An AC is a community that exhibits structure and keyword cohesiveness. To facilitate ACQ evaluation, we develop the CL-tree index and its query algorithms. Our experimental results show that ACs are easier to interpret than those of existing community detection/search methods, and they can be “personalized”. Our solutions are also faster than existing community search algorithms.

We will study the use of other measures of structure cohesiveness (e.g.,  $k$ -truss,  $k$ -clique) and keyword cohesiveness (e.g., Jaccard similarity and string edit distance) in the ACQ definition. We will also investigate how the directions of edges will affect the formation of an AC. We will examine how graph pattern matching techniques [28, 8, 9] can be extended to find ACs. An interesting research direction is to study how to automatically generate a meaningful graph pattern that reflects a real community, and how to use these patterns to find ACs.

## Acknowledgments

Reynold Cheng, Yixiang Fang, Siqiang Luo, and Jiafeng Hu were supported by the Research Grants Council of Hong Kong (RGC Project HKU 17205115) and HKU (Project 104004129).

## 9. REFERENCES

- [1] [https://en.wikipedia.org/wiki/Disjoint-set\\_data\\_structure](https://en.wikipedia.org/wiki/Disjoint-set_data_structure).
- [2] V. Batagelj and M. Zaversnik. An  $o(m)$  algorithm for cores decomposition of networks. *arXiv*, 2003.
- [3] G. Bhalotia et al. Keyword searching and browsing in databases using banks. In *ICDE*, 2002.
- [4] W. Cui, Y. Xiao, H. Wang, Y. Lu, and W. Wang. Online search of overlapping communities. In *SIGMOD*, 2013.
- [5] W. Cui, Y. Xiao, H. Wang, and W. Wang. Local search of communities in large graphs. In *SIGMOD*, 2014.
- [6] B. Ding, J. X. Yu, S. Wang, L. Qin, X. Zhang, and X. Lin. Finding top- $k$  min-cost connected trees in databases. In *ICDE*, 2007.
- [7] S. N. Dorogovtsev et al.  $K$ -core organization of complex networks. *Physical review letters*, 2006.
- [8] W. Fan, J. Li, S. Ma, N. Tang, Y. Wu, and Y. Wu. Graph pattern matching: from intractable to polynomial time. *PVLDB*, 2010.
- [9] W. Fan, X. Wang, Y. Wu, and J. Xu. Association rules with graph patterns. *PVLDB*, 8(12):1502–1513, 2015.
- [10] Y. Fang, H. Zhang, Y. Ye, and X. Li. Detecting hot topics from twitter: A multiview approach. *Journal of Information Science*, 40(5):578–593, 2014.
- [11] S. Fortunato. Community detection in graphs. *Physics Reports*, 486(3):75–174, 2010.
- [12] J. Han, M. Kamber, and J. Pei. *Data mining: concepts and techniques*. Elsevier, 2011.
- [13] J. Han, J. Pei, and Y. Yin. Mining frequent patterns without candidate generation. In *SIGMOD*, 2000.
- [14] H. He, H. Wang, J. Yang, and P. S. Yu. Blinks: ranked keyword searches on graphs. In *SIGMOD*, 2007.
- [15] X. Huang, H. Cheng, L. Qin, W. Tian, and J. X. Yu. Querying  $k$ -truss community in large and dynamic graphs. In *SIGMOD*, 2014.
- [16] V. Kacholia et al. Bidirectional expansion for keyword search on graph databases. In *VLDB*, 2005.
- [17] M. Kargar and A. An. Keyword search in graphs: Finding  $r$ -cliques. *PVLDB*, 4(10):681–692, 2011.
- [18] R.-H. Li, L. Qin, J. X. Yu, and R. Mao. Influential community search in large networks. In *PVLDB*, 2015.
- [19] R.-H. Li, J. X. Yu, and R. Mao. Efficient core maintenance in large dynamic graphs. *TKDE*, 2014.
- [20] Y. Liu, A. Niculescu-Mizil, and W. Gryc. Topic-link lda: joint models of topic and author community. In *ICML*, 2009.
- [21] R. M. Nallapati, A. Ahmed, E. P. Xing, and W. W. Cohen. Joint latent topic models for text and citations. In *KDD*, 2008.
- [22] M. Newman et al. Finding and evaluating community structure in networks. *Physical review E*, 2004.
- [23] Y. Ruan, D. Fuhry, and S. Parthasarathy. Efficient community detection in large networks using content and links. In *WWW*, 2013.
- [24] M. Sachan et al. Using content and interactions for discovering communities in social networks. In *WWW*, 2012.
- [25] S. B. Seidman. Network structure and minimum degree. *Social networks*, 5(3):269–287, 1983.
- [26] M. Sozio and A. Gionis. The community-search problem and how to plan a successful cocktail party. In *KDD*, 2010.
- [27] B. Thomee et al. The new data and new challenges in multimedia research. *arXiv:1503.01817*, 2015.
- [28] H. Tong, C. Faloutsos, B. Gallagher, and T. Eliassi-Rad. Fast best-effort pattern matching in large attributed graphs. In *KDD*, 2007.
- [29] Z. Xu, Y. Ke, Y. Wang, H. Cheng, and J. Cheng. A model-based approach to attributed graph clustering. In *SIGMOD*, 2012.
- [30] Y. Fang et al. Effective community search for large attributed graphs (technical report). <http://www.cs.hku.hk/research/techreps/document/TR-2016-01.pdf>.
- [31] J. Yang, J. McAuley, and J. Leskovec. Community detection in networks with node attributes. In *ICDM*, 2013.
- [32] T. Yang, R. Jin, Y. Chi, and S. Zhu. Combining link and content for community detection: a discriminative approach. In *KDD*, 2009.
- [33] J. X. Yu, L. Qin, and L. Chang. Keyword search in databases. *Synthesis Lectures on Data Management*, 2009.
- [34] Y. Zhou, H. Cheng, and J. X. Yu. Graph clustering based on structural/attribute similarities. *VLDB*, 2009.