

SQAC: Scalable Querying of Attribute-Constrained (α, β) -Cores over Large Bipartite Graphs

Xin Deng
Hunan University
Changsha, China
xindeng@hnu.edu.cn

Peng Peng
Hunan University
Changsha, China
hnu16pp@hnu.edu.cn

Baoqing Sun
Hunan University
Changsha, China
xueqing@hnu.edu.cn

Shuo Dai
Hunan University
Changsha, China
ds0230@hnu.edu.cn

Zheng Qin
Hunan University
Changsha, China
zqin@hnu.edu.cn

Lijun Chang
The University of Sydney
Sydney, Australia
Lijun.Chang@sydney.edu.au

Xuemin Lin
Shanghai Jiaotong University
Shanghai, China
xuemin.lin@sjtu.edu.cn

ABSTRACT

Many important real-world networks can be effectively modeled as bipartite graphs, where vertex attributes convey critical semantic information essential for graph analysis. As a fundamental subgraph structure in bipartite graphs, (α, β) -cores have attracted extensive attention and been widely used. However, existing studies have largely ignored the attribute property in bipartite graphs, which limits attribute-aware applications in the real world. In this paper, we formulate a new problem of querying attribute-constrained (α, β) -cores over bipartite graphs, and study index-based methods. We first propose a vertex-based core index (VC-Index), which maintains the minimal attribute-constrained set pairs for each vertex and each possible (α, β) pair. To improve the query efficiency, we then design an attribute-based core index (AC-Index) that computes an index for each possible attribute set. However, these two indexes suffer from either unsatisfactory query performance or excessive space overhead. Therefore, we further construct the minimized AC-Index, named MAC-Index, which significantly reduces the index size by leveraging the containment relationship of the attribute set for the (α, β) -core, while guaranteeing efficient query processing. Extensive experiments over real-world datasets confirm the efficiency and effectiveness of our index-based algorithms.

PVLDB Reference Format:

Xin Deng, Peng Peng, Baoqing Sun, Shuo Dai, Zheng Qin, Lijun Chang, and Xuemin Lin. SQAC: Scalable Querying of Attribute-Constrained (α, β) -Cores over Large Bipartite Graphs. PVLDB, 14(1): XXX-XXX, 2026. doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/XueQingSBQ/SQAC>.

1 INTRODUCTION

Bipartite graphs have been widely used in the real world, such as purchase networks over custom-product relation [5, 38, 40] and

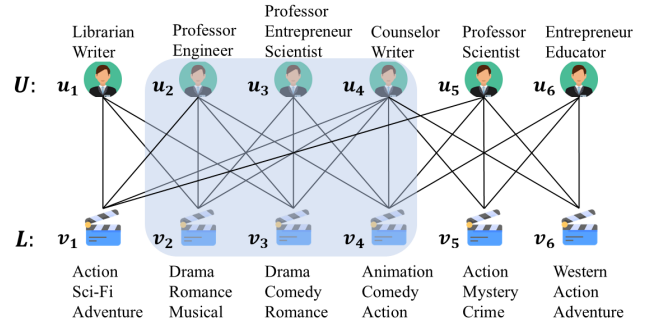


Figure 1: A user-movie network

citation networks over author-paper relation [8, 18]. A bipartite graph is represented as $G = (U, L, E)$, with U and L being two disjoint vertex sets that represent different types of entities. An edge $e \in E$ connects a vertex in U and another in L . Analysis of bipartite graphs is of great significance, and computing (α, β) -core is one of the fundamental problems over bipartite graphs [6, 9, 23, 29, 30, 43, 48]. The (α, β) -core is defined as a maximal subgraph of the given bipartite graph G whose vertices in the upper layer (i.e., U) have a degree of at least α and vertices in the lower layer (i.e., L) have a degree of at least β . Efficient querying of (α, β) -cores has wide-ranging applications, such as group recommendation [24, 34], community detection [43, 46], and fraudster detection [5, 8].

However, existing studies over (α, β) -core mainly emphasize the structure information of bipartite graphs, such as (α, β) -core decomposition [21, 26, 29, 30] and (α, β) -core maintenance [6, 31], while disregarding the attribute properties of vertices. Many real-world bipartite graphs are vertex-attributed, in which attributes can reflect important semantic information. For example, in E-commerce websites such as Amazon and Alibaba, customers and products form a bipartite graph, where an edge indicates the purchase relationship between a customer and a product. The attribute of a customer may indicate age or occupation, and the attribute of a product may represent the category or popularity.

In this paper, we study a new problem of Scalable Querying of Attribute-constrained (α, β) -Cores (SQAC, for short) over bipartite graphs. Given two integers α, β and two query attribute sets Q_U, Q_L , SQAC aims to identify an attribute-aware (α, β) -core H from an input graph G in which each vertex $u \in U(H)$ (resp. $v \in L(H)$) has at least one attribute from Q_U (resp. Q_L).

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097. doi:XX.XX/XXX.XX

Application. Efficiently analyzing attribute-aware (α, β) -core is of great importance. We present two real applications as follows, and the corresponding case studies are presented in Section 6.5.

① *Personalized group recommendation.* Consider a user-movie network G of IMDB (<https://www.imdb.com>) in Figure 1, where an edge indicates a user's preference for a movie. A vertex in U represents a user, and the associated attribute indicates its occupation. A vertex in V represents a movie, and the associated attributes indicate its genres. Assume a new user u_0 joins IMDB with the objective of connecting with faculty members (e.g., "professor" or "counselor") who share similar tastes in movies (e.g., "Drama", "Comedy" or "Romance"), our SQAC can return a tailored cohesive group H (marked in blue) for u_0 , which is composed of the users $U(H) = \{u_2, u_3, u_4\}$ and the movies $V(H) = \{v_2, v_3, v_4\}$. From Figure 1, we can see that, although both H and G are $(3, 3)$ -core, H is evidently more appropriate for the new user u_0 , since both members and movies in H are highly compatible with u_0 's preferences.

② *Accelerate the querying of attributed community.* Another significant application of SQAC is to serve as a building block for complicated dense subgraph models. For example, Xu et al. [46] formulate an attributed (α, β) -community model, which aims to identify a connected (α, β) -core H containing the query vertex, and the vertices of H in the same layer share the most attributes under given query attribute sets. Apparently, the attributed (α, β) -community is a subset of our SQAC under the same query conditions, as our SQAC problem relaxes the strict attribute constraints on vertices. There is a state-of-the-art (SOTA) method, called Inc [46], proposed for the attributed (α, β) -community. An essential step of its query processing is to compute the unit (α, β) -cores by traversing the whole graph. On one hand, we can obtain the unit (α, β) -cores by scanning the resulting subgraph of our SQAC, rather than the whole input graph, which significantly narrows the search scope. On the other hand, if the query result of SQAC is empty, clearly we do not need to execute the subsequent steps since the community cannot exist either, avoiding unnecessary time costs. Therefore, our SQAC can function as an optimized procedure for the attributed community search, which significantly speeds up query efficiency.

Challenges. Efficiently computing the SQAC is challenging. Firstly, although we can derive a naive online approach of SQAC based on the existing online approach [21] for the computation of (α, β) -core, such an online approach suffers from extremely poor query efficiency due to the time-consuming peeling process. Its time complexity is presented in Table 1. Secondly, when we resort to index-based approaches, a basic idea is to construct the existing Bicore-Index [29] for each possible attribute set pair, where the Bicore-Index has been proposed for the (α, β) -core query over non-attributed bipartite graphs. Although this index exhibits optimal query performance, it requires $O(m \cdot 2^{|\mathcal{A}_U|+|\mathcal{A}_L|})$ space for storage and $O(\delta \cdot m \cdot 2^{|\mathcal{A}_U|+|\mathcal{A}_L|})$ time for construction, where m represents the number of edges and \mathcal{A}_U (resp. \mathcal{A}_L) denotes the set of attribute types in U (resp. L); this exponential complexity is not acceptable in real applications. Thirdly, the (α, β) -core has two correlated parameters that need to be considered simultaneously. As a result, analysis of the relationship between (α, β) -cores is more difficult than k -cores [7], which requires a two-dimensional method rather than one-dimensional monotonicity. Furthermore, when all possible

Table 1: Comparison of different solutions. n' (resp. m') denotes the number of vertices (resp. edges) in the intermediate subgraph G' , \mathcal{A}_U (resp. \mathcal{A}_L) denotes the set of attribute types in $U(G)$ (resp. $L(G)$), $|R|$ denotes the size of query result, μ denotes the average number of MASs for each vertex and each possible (α, β) pair in VC-Index, γ denotes the average number of vertices for each vertex block in AC-Index, D_{max} denotes the maximum vertex degree in G , \bar{d} denotes the average vertex degree in G , and ρ denotes the average occurrence frequency of all vertices in G within MAC-Index.

Index	Query time	Index space	Construction time
Online	$O(m)$		
I_{ABi}	$O(R)$	$O(m \cdot 2^{ \mathcal{A}_U + \mathcal{A}_L })$	$O(\delta \cdot m \cdot 2^{ \mathcal{A}_U + \mathcal{A}_L })$
I_{VC}	$O(n \cdot \mu \cdot D_{max} + m')$	$O(n \cdot D_{max}^2 \cdot \mu)$	$O(n \cdot D_{max}^2 \cdot \mu)$
I_{AC}	$O(n' + m')$	$O(m \cdot (2^{ \mathcal{A}_U } + 2^{ \mathcal{A}_L }))$	$O(n \cdot D_{max}^2 \cdot \bar{d} \cdot \log \gamma)$
I_{MAC}	$O(n' + m')$	$O(\rho \cdot n)$	$O(n \cdot D_{max}^2 \cdot \bar{d} \cdot \log n)$

pairs of attribute sets are considered, the containment relationship among (α, β) -cores becomes intensely complicated.

Our Approach. In response to the challenges outlined above, we first propose a Vertex-based Core index (VC-Index, for short) to answer the SQAC. We construct a VC-Index by maintaining some pivotal attribute sets for each vertex and each possible integer pair (α, β) . However, this method is inefficient in terms of query performance, since it requires checking the index entries of each vertex of G . To improve the query performance, we propose an Attribute-based Core index (AC-Index, for short), which computes the Bicore-Index for each possible attribute set, optimizing the query performance to be related to each query rather than the input graph size. While the AC-Index offers outstanding query efficiency, it incurs exponential space complexity. Furthermore, we design a superior-optimized index called Minimized attribute-based Core index (MAC-Index, for short), which compresses redundant information in the AC-Index leveraging the containment relationship of the attribute set for (α, β) -cores. The MAC-Index significantly optimizes the space cost from $O(m \cdot (2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|}))$ of the AC-Index to $O(\rho \cdot n)$, while ensuring the query efficiency comparable to the AC-Index. Here, ρ denotes the average occurrence frequency of all vertices in G (ρ is a small integer and $\rho \ll n$ in practice).

Contributions. Our contributions are summarized as follows:

- **The SQAC problem.** To the best of our knowledge, we are the first work to propose and solve the attribute-constrained (α, β) -core queries over large attributed bipartite graphs.
- **Two basic index structures.** We propose two basic indexes, called VC-Index and AC-Index, which can correctly answer the query of our proposed SQAC and provide better query efficiency than online query processing.
- **A superior index structure.** We design a superior index structure, called MAC-Index, which can accurately describe the containment property of attribute set for the (α, β) -core. MAC-Index not only significantly reduces the index size from $O(m \cdot (2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|}))$ to $O(\rho \cdot n)$ (ρ is a small integer in practice), but also guarantees efficient query performance.
- **Efficient algorithm for index construction and maintenance.** We also propose an efficient algorithm for constructing MAC-Index, and effective techniques to enable the maintenance of MAC-Index.

- **Extensive experiments.** Comprehensive performance studies on real datasets demonstrate the efficiency and effectiveness of our approaches proposed in this paper.

2 BACKGROUND

In this section, we first give the formal definition of concepts and problem statement (Section 2.1). We then present an online query approach and a naive index for the SQAC problem (Sections 2.2-2.3).

2.1 Preliminaries

DEFINITION 2.1 (ATTRIBUTED BIPARTITE GRAPH). An attributed bipartite graph is defined as $G = (V = (U \cup L), E, \mathcal{A}_U, \mathcal{A}_L)$, where U, L are two disjoint vertex sets and $E \subseteq U \times L$ represents the edge set. Each vertex $u \in V$ is associated with a set of attributes denoted by $A(u)$. We use $a_i(u)$ to denote the i -th attribute of u . \mathcal{A}_U (resp. \mathcal{A}_L) denotes the set of attribute types in U (resp. L), i.e., $\mathcal{A}_U = \bigcup_{u \in U} A(u)$. Each edge $e \in E$ between two vertices $u \in U$ and $v \in L$ is denoted by (u, v) . We may also use $V(G), U(G), L(G), E(G), \mathcal{A}_U(G), \mathcal{A}_L(G)$ to denote $V, U, L, E, \mathcal{A}_U, \mathcal{A}_L$ of G , respectively.

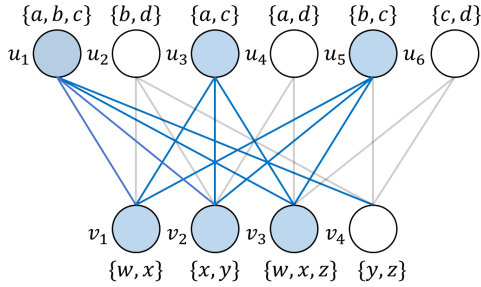


Figure 2: An attributed bipartite graph G

Given an attributed bipartite graph $G = (V = (U \cup L), E, \mathcal{A}_U, \mathcal{A}_L)$ in Figure 2, we define the attribute-induced subgraph of G over an attribute set pair (A_U, A_L) as $G[A_U, A_L]$, where each vertex $u \in U(G[A_U, A_L])$ (resp. $v \in L(G[A_U, A_L])$) contains at least one attribute $a_i(u) \in A_U$ (resp. A_L), i.e., $U(G[A_U, A_L]) = \{u \in U(G) \mid A(u) \cap A_U \neq \emptyset\}$. Particularly, we refer to $G[A_U, \mathcal{A}_L]$ with $A_U \in \mathcal{A}_U$ as a subgraph of G induced by the attribute set A_U . We have a symmetrical statement for $G[\mathcal{A}_U, A_L]$.

Example 2.1. Figure 2 presents an attributed bipartite graph G . Given an attribute set pair $(A_U = \{a, c\}, A_L = \{w, x\})$, the attribute-induced subgraph $G[\{a, c\}, \{w, x\}] = \{U = \{u_1, u_3, u_4, u_5, u_6\}, L = \{v_1, v_2, v_3\}\}$, as shown in Figure 3.

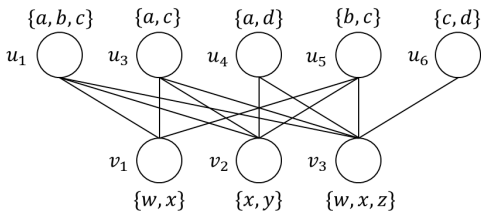


Figure 3: The attribute-induced subgraph $G[\{a, c\}, \{w, x\}]$

We use n and m to denote the number of vertices and edges in G , respectively. The degree of a vertex u in G is denoted as

$\deg(u, G)$. In addition, we use $d_{\max}^U(G)$ (resp. $d_{\max}^L(G)$) to denote the maximum degree among all vertices in $U(G)$ (resp. $L(G)$). We also use $D_{\max}(G)$ to denote the maximum degree among all vertices in G , i.e., $D_{\max}(G) = \max\{d_{\max}^U(G), d_{\max}^L(G)\}$.

DEFINITION 2.2 ((α, β)-CORE [21]). Given a bipartite graph $G = (U, L, E)$ and an integer pair (α, β) , the (α, β) -core of G , denoted by $C_{\alpha, \beta}(G)$, is the maximal subgraph such that each vertex in $U(C_{\alpha, \beta})$ has a degree of at least α and each vertex in $L(C_{\alpha, \beta})$ has a degree of at least β within $C_{\alpha, \beta}(G)$. We may also use $C_{\alpha, \beta}$ for $C_{\alpha, \beta}(G)$.

There are some typical properties for (α, β) -core [29, 31]: (1) (α, β) -cores have partial nested relationship, i.e., for two (α, β) -core $C_{\alpha, \beta}$ and (α', β') -core $C_{\alpha', \beta'}$, the $C_{\alpha, \beta}$ is **contained** in the $C_{\alpha', \beta'}$ if $(\alpha' < \alpha, \beta' \leq \beta)$ or $(\alpha' \leq \alpha, \beta' < \beta)$. (2) An (α, β) -core may be disconnected, but its connected components are usually the "communities" that have been extensively studied.

DEFINITION 2.3 (BI-CORE NUMBER). Given a bipartite graph G and a vertex u , an integer pair (α, β) is a bi-core number of u if the (α, β) -core $C_{\alpha, \beta}$ contains u and there does not exist any other (α', β') -core $C_{\alpha', \beta'}$ containing u such that $C_{\alpha', \beta'}$ is contained in $C_{\alpha, \beta}$.

Note that, for a vertex u , there may exist multiple bi-core numbers in the graph G . All bi-core numbers can be computed by a peeling process that iteratively removes the vertices with the smallest degree [21] (a.k.a. (α, β) -core decomposition).

DEFINITION 2.4 (THE SQAC PROBLEM). Given an attributed bipartite graph $G = (V = (U \cup L), E, \mathcal{A}_U, \mathcal{A}_L)$, two integers α, β and two query attribute sets Q_U, Q_L , SQAC aims to identify an attribute-aware (α, β) -core H from the graph G in which each vertex $u \in U(H)$ (resp. $v \in L(H)$) has at least one attribute from Q_U (resp. Q_L).

Example 2.2. Consider the attributed bipartite graph G in Figure 2. Given two integers $\alpha = 3, \beta = 3$ and two query attribute sets $Q_U = \{a, c\}, Q_L = \{w, x\}$, the result of SQAC is $\{U = \{u_1, u_3, u_5\}, L = \{v_1, v_2, v_3\}\}$, marked in blue in the graph.

2.2 Online Query Approach

Existing work [21] has proposed an efficient online algorithm for the query of (α, β) -core over bipartite graphs, which can be easily adapted to the SQAC problem. The pseudocode is presented in Algorithm 1. Specifically, given two integers α, β and two query attribute sets Q_U, Q_L , we first construct the attribute-induced subgraph $G[Q_U, Q_L]$, and then iteratively remove the vertices that violate degree constraints from $G[Q_U, Q_L]$ (i.e., peeling process) in Lines 3-12. The search process terminates when the remaining vertices form an (α, β) -core. The time complexity is bounded by $O(m)$. Although this approach can correctly answer the SQAC, it is difficult to apply in large-scale datasets due to time-consuming query processing, prompting us to explore efficient indexing approaches.

2.3 Naive Index

In this subsection, we design a naive index for the proposed SQAC, named Attributed Bicore-Index (ABi-Index, for short), denoted by \mathcal{I}_{ABi} , which is adapted from the existing work Bicore-Index [29]. Before introducing ABi-Index, let's discuss the Bicore-Index for the query of non-attributed (α, β) -core. For convenience, we use $\beta_{\max, \alpha}(u, G)$ (resp. $\alpha_{\max, \beta}(u, G)$) denote the maximum value of β

Algorithm 1: Online Query

Input: An attributed bipartite graph G , two integers α, β and two query attribute sets Q_U, Q_L

Output: The result of the SQAC

```

1  $G' \leftarrow$  Construct the attribute-induced subgraph  $G[Q_U, Q_L]$ ;
2 while  $V(G') \neq \emptyset$  do
3   foreach  $u \in U(G')$  do
4     if  $\deg(u, G') < \alpha$  then
5       Remove  $u$  and its incident edges from  $G'$ ;
6     foreach  $v \in N(u, G')$  do
7        $\deg(v, G') \leftarrow \deg(v, G') - 1$ ;
8   foreach  $v \in L(G')$  do
9     if  $\deg(v, G') < \beta$  then
10      Remove  $v$  and its incident edges from  $G'$ ;
11     foreach  $u \in N(v, G')$  do
12        $\deg(u, G') \leftarrow \deg(u, G') - 1$ ;
13 if all vertices in  $G'$  satisfy the degree constraints then Break;
14 return  $G'$ ;

```

for the vertex u and a specific α such that u is contained in the (α, β) -core of the graph G .

2.3.1 Bicore-Index. Given a bipartite graph G , the Bicore-Index of G , denoted by \mathcal{I}_{Bi} , is a three-level tree structure with two parts for vertices in U and L , respectively, denoted by \mathcal{I}_{Bi}^U and \mathcal{I}_{Bi}^L . We present only \mathcal{I}_{Bi}^U as follows, since \mathcal{I}_{Bi}^L is symmetrical to \mathcal{I}_{Bi}^U .

- The first level is an array containing d_{max}^U pointers to the arrays at the second level, where d_{max}^U denotes the maximum degree among all vertices in U .
- The second level is a sub-table containing d_{max}^U arrays. The length of the α -th array is equal to the maximum β value (i.e., if the (α, β) -core exists).
- The third level consists of vertex blocks. Each vertex block $\mathcal{I}_{Bi}^U(\alpha, \beta)$ is associated with an integer pair (α, β) and contains the vertices $u \in U$ with $\beta_{max, \alpha}(u) = \beta$.

The size of \mathcal{I}_{Bi} is $O(m)$ and the construction time is $O(\delta \cdot m)$ [29]. For query processing, it retrieves the (α, β) -core by collecting the vertices in any vertex block $\mathcal{I}_{Bi}^U(\alpha, \beta')$ (resp. $\mathcal{I}_{Bi}^L(\alpha', \beta)$) with $\beta' \geq \beta$ (resp. $\alpha' \geq \alpha$) based on the partial nest relationship of (α, β) -core. The query time is $O(|R|)$, where $|R|$ denotes the size of query result.

Example 2.3. Figure 4 presents the \mathcal{I}_{Bi} for the subgraph in Figure 3, which also illustrates the procedure of computing (3, 3)-core. Processing steps are shown in bold arrows and the visited elements are marked in blue.

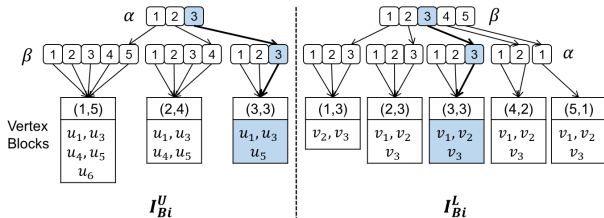


Figure 4: The Bicore-Index \mathcal{I}_{Bi} for $G[\{a, c\}, \{w, x\}]$ in Figure 3

2.3.2 ABi-Index. Given that the SOTA solution Bicore-Index is proposed for the query of non-attributed (α, β) -core, a naive index

for our SQAC is to compute the Bicore-Index for each possible attribute-induced subgraph. We refer to this index as the Attributed Bicore-Index (ABi-Index, for short), denoted by \mathcal{I}_{ABi} . Although ABi-Index can achieve optimal query time cost $O(|R|)$, it takes $O(m \cdot 2^{|\mathcal{A}_U|+|\mathcal{A}_L|})$ space for storage and $O(\delta \cdot m \cdot 2^{|\mathcal{A}_U|+|\mathcal{A}_L|})$ time for construction, which suffers from extremely poor scalability due to its exponential complexity.

3 VERTEX-BASED CORE INDEX

In this section, we first analyze the containment property of attribute set for the (α, β) -cores and introduce the concept of minimal attribute-constrained set (Section 3.1), based on which we then propose a vertex-based core index (Section 3.2). Finally, we present the query process in Section 3.3.

3.1 Minimal Attribute-constrained Set

Before introducing the concept of minimal attribute-constrained set, let's discuss an important lemma as follows:

Lemma 3.1. *Given an attributed bipartite graph $G = (V = (U \cup L), E, \mathcal{A}_U, \mathcal{A}_L)$ and two integers α, β , for two subgraphs $G[A_U, A_L]$ and $G[A'_U, A'_L]$, if $A_U \subseteq A'_U$ and $A_L \subseteq A'_L$, then the (α, β) -core $C_{\alpha, \beta}(G[A_U, A_L])$ is contained in the (α, β) -core $C_{\alpha, \beta}(G[A'_U, A'_L])$.*

PROOF. Since A_U and A_L are the subsets of A'_U and A'_L , respectively, it follows that $G[A_U, A_L]$ is a subgraph of $G[A'_U, A'_L]$. Hence, the lemma holds. \square

Lemma 3.1 illustrates the containment relationship of the attribute set for the (α, β) -core. However, considering two types of attribute sets simultaneously is inherently more complicated than considering just one. Consequently, we derive a special case (i.e., Lemma 3.2) of Lemma 3.1 as follows.

Lemma 3.2. *Given an attributed bipartite graph $G = (V = (U \cup L), E, \mathcal{A}_U, \mathcal{A}_L)$ and two integers α, β , for two subgraphs $G[A_U, \mathcal{A}_L]$ and $G[A'_U, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, A_L]$ and $G[\mathcal{A}_U, A'_L]$), if $A_U, A'_U \subseteq \mathcal{A}_U$ and $A_U \subset A'_U$, then the (α, β) -core $C_{\alpha, \beta}(G[A_U, \mathcal{A}_L])$ is contained in the (α, β) -core $C_{\alpha, \beta}(G[A'_U, \mathcal{A}_L])$.*

According to Lemma 3.2, we can easily infer that, for any vertex $u \in C_{\alpha, \beta}(G[A_U, \mathcal{A}_L])$, if $C_{\alpha, \beta}(G[A_U, \mathcal{A}_L])$ is contained in $C_{\alpha, \beta}(G[A'_U, \mathcal{A}_L])$, then u must be in $C_{\alpha, \beta}(G[A'_U, \mathcal{A}_L])$. For convenience, we call A an attribute-constrained set for u and (α, β) if the vertex u is in the (α, β) -core of $G[A, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, A]$). Since there may exist multiple attribute-constrained sets for u and (α, β) , we also use $AS(u, (\alpha, \beta))$ to denote the set of attribute-constrained set over u and (α, β) . From this perspective, we propose a novel concept of minimal attribute-constrained set as follows.

DEFINITION 3.1 (MINIMAL ATTRIBUTE-CONSTRAINED SET). *Given an attributed bipartite graph G , for an attribute-constrained set $A \in AS(u, (\alpha, \beta))$, we define A as a Minimal Attribute-constrained Set (MAS, for short) over u and (α, β) , if there does not exist any other attribute-constrained set $A' \in AS(u, (\alpha, \beta))$ such that $A' \subset A$. Considering that there may be multiple minimal attribute-constrained sets, we use $MAS(u, (\alpha, \beta))$ to denote the set of the minimal attribute-constrained set over u and (α, β) .*

Based on Definition 3.1, we can infer that, for an attribute set $A_U \subseteq \mathcal{A}_U$ (resp. \mathcal{A}_L), if there does not exist any MAS $A'_U \in \text{MAS}(u, (\alpha, \beta))$ such that $A'_U \subset A_U$, then the vertex u would not be in the (α, β) -core of any subgraph $G[A_U, A_L]$ with $A_L \subseteq \mathcal{A}_L$. Therefore, computing the MASs for each vertex and each possible (α, β) pair can help filter out unpromising vertices which are not in the (α, β) -core of $G[Q_U, Q_L]$ over two query attribute sets Q_U, Q_L .

Example 3.1. Given an attributed bipartite graph G in Figure 2, for the vertex u_3 and an integer pair $(\alpha, \beta) = (3, 3)$, there are two MASs for u_3 and $(3, 3)$, i.e., $\{c\}$ and $\{a, b\}$.

3.2 The VC-Index Structure

In this subsection, we propose a Vertex-based Core Index (VC-Index, for short), which maintains some pivotal attribute sets for each vertex and each possible (α, β) pair. On one hand, VC-Index exploits attribute constraints to guarantee that each added attribute set is a MAS, thereby improving pruning efficiency. On the other hand, VC-Index makes use of the partial nested relationship of (α, β) -core to further compress storage redundancy. We formally define the VC-Index as follows.

DEFINITION 3.2 (VC-INDEX). Given an attributed bipartite graph $G = (V = (U, L), E, \mathcal{A}_U, \mathcal{A}_L)$, the VC-Index of G , formally denoted by I_{VC} , is essentially a two-dimensional matrix structure with two parts for all vertices in $U(G)$ and $L(G)$, respectively, denoted by I_{VC}^U and I_{VC}^L . Each cell $I_{VC}^U(u, (\alpha, \beta))$ in I_{VC}^U is determined by both a vertex u and an integer pair (α, β) , which maintains some attribute sets for u and (α, β) , adhering to the following criteria (I_{VC}^L is symmetric to I_{VC}^U):

- **Attribute constraint.** For each cell $I_{VC}^U(u, (\alpha, \beta))$ in I_{VC}^U and an attribute set $A \in I_{VC}^U(u, (\alpha, \beta))$, A must be a MAS for u and (α, β) .
- **Partial constraint.** For each cell $I_{VC}^U(u, (\alpha, \beta))$ in I_{VC}^U and an attribute set $A \in I_{VC}^U(u, (\alpha, \beta))$, there does not exist any MAS $(u, (\alpha, \beta'))$ with $\beta' > \beta$ that contains A .

For the sake of presentation only, we use $I_{VC}^{U, \mathcal{A}_U}$ to denote such a two-dimensional matrix structure with each cell $I_{VC}^{U, \mathcal{A}_U}(u, (\alpha, \beta))$ maintaining a set of attribute sets $A \in \mathcal{A}_U$ for u and (α, β) . We have a symmetrical definition for $I_{VC}^{U, \mathcal{A}_L}$. Note that, I_{VC}^U can be essentially viewed as a combination of $I_{VC}^{U, \mathcal{A}_U}$ and $I_{VC}^{U, \mathcal{A}_L}$. Moreover, we do not consider the case of $(\alpha, \beta) = (1, 1)$ for any vertex in I_{VC} , since the query can be efficiently done by checking if the vertices have neighbors over query attribute sets. The space complexity of VC-Index is shown in Theorem 3.1.

Theorem 3.1. Given an attributed bipartite graph G , the size of I_{VC} is bounded by $O(n \cdot D_{max}^2 \cdot \mu)$, where μ denotes the average number of MASs for each vertex and each possible (α, β) pair in I_{VC} .

PROOF. The possible (α, β) pair of all vertices in I_{VC} takes up no more than $O(\sum_{u \in V(G)} (deg(u, G) \cdot D_{max}))$ space. Since $deg(u, G) < D_{max}$, we have $\sum_{u \in V(G)} (deg(u, G) \cdot D_{max}) \leq \sum_{u \in V(G)} D_{max}^2 = n \cdot D_{max}^2$. Hence, the space is bounded by $O(n \cdot D_{max}^2 \cdot \mu)$. \square

Algorithm 2: Query based on VC-Index

Input: An attributed bipartite graph G , a VC-Index I_{VC} of G , two integers α, β and two query attribute sets Q_U, Q_L

Output: The result of SQAC

```

1  $R \leftarrow \emptyset$ ;
2 foreach vertex  $u \in U(G)$  do
3   foreach  $\beta' : \text{from } \beta \text{ to } \beta_{max}(u, \alpha)$  do
4     if there exist two MASs  $A_U, A_L \in I_{VC}^U(u, (\alpha, \beta'))$  such
       that  $A_U \subseteq Q_U$  and  $A_L \subseteq Q_L$  then
5        $R \leftarrow R \cup \{u\}$ ;
6     break;
7 foreach vertex  $v \in L(G)$  do
8   Do the operations that are symmetrical to those in Lines 3-6 in
     Algorithm 2;
9  $G' \leftarrow$  the subgraph of  $G$  induced by  $R$ ;
10 Do the same operation as Lines 2-13 in Algorithm 1 to guarantee
    the correctness of the final results;
11 return  $G'$ ;

```

3.3 Query Processing

Let's discuss the VC-Index based query algorithm, which is supported by the following lemma. Note that we only present I_{VC}^U since I_{VC}^L is symmetrical to I_{VC}^U . For convenience, we use $\beta_{max}(u, \alpha)$ to denote the maximum β value of the vertex $u \in U(G)$ for α in I_{VC}^U .

Lemma 3.3. Given a VC-Index I_{VC} of an attribute bipartite graph G , two integers α, β and two query attribute sets Q_U, Q_L . For a vertex $u \in U(G)$, if there exist two MASs $A_U, A_L \in I_{VC}^U(u, (\alpha, \beta'))$ ($\beta \leq \beta' \leq \beta_{max}(u, \alpha)$) in I_{VC} such that $A_U \subseteq Q_U$ and $A_L \subseteq Q_L$, then u may be in the (α, β) -core of $G[Q_U, Q_L]$. Otherwise, u would not be in the (α, β) -core of $G[Q_U, Q_L]$.

PROOF. This lemma can be easily derive from Definition 3.2. \square

Example 3.2. Given a part of the I_{VC} in Figure 5, two integers $\alpha = 3, \beta = 3$ and two query attribute sets $Q_U = \{a, c\}, Q_L = \{w, x\}$. For the vertex u_2 , the maximum value of β is 3 for $\alpha = 3$. We can infer that u_2 is not in the $(3, 3)$ -core of $G[\{a, c\}, \{w, x\}]$, since there does not exist any MAS in $I_{VC}^{U, \mathcal{A}_U}(u_2, (3, 3))$ that is a subset of Q_U .

Query Algorithm. The pseudocode for VC-Index based query is presented in Algorithm 2. Given an index I_{VC} of G , two integers α, β and two query attribute sets Q_U, Q_L , we first traverse all vertices in $U(G)$ and push the candidate vertices that may be included in $C_{\alpha, \beta}(G[Q_U, Q_L])$ into the temporary result set R (Lines 2-6) according to Lemma 3.3. We then perform the symmetrical operation for the vertices in $L(G)$ and obtain the subgraph G' of G induced by R (Lines 7-9). For the sake of presentation, we also refer to G' as the **intermediate subgraph**, which is essentially the intersection of the (α, β) -core in $G[Q_U, \mathcal{A}_L]$ and $G[\mathcal{A}_U, Q_L]$. Finally, to guarantee the correctness of the final result, we iteratively remove the vertices in G' that do not satisfy the degree constraints (Line 10).

Theorem 3.2. Given a I_{VC} of G , the time complexity of the I_{VC} based query algorithm is bounded by $O(n \cdot \mu \cdot D_{max} + m')$, where m' denotes the number of edges in the intermediate subgraph.

		Vertex					
		u_1	u_2	u_3	u_4	u_5	u_6
(α, β)	(4,2)	{b}, {c}				{b}, {c}	
	(4,1)	{a}					
	(3,3)	{b}, {c}	{b}, {c, d}	{c}, {a, b}		{b}, {c}	
	(3,2)	{a}	{a, d}	{a}			
	(3,1)		{d}				
	(2,4)	{a, b}, {a, c}	{a, b}, {b, c}	{a, b}, {a, c}	{a, b}, {a, c}	{b, c}, {b, c}	
		{b, c}, {b, d}	{b, d}, {c, d}	{b, c}, {c, d}	{b, c}, {b, d}	{c, d}	
	(2,3)	{a}, {b}, {c}	{b}, {a, d}	{a}, {c}	{a}	{c}, {b}	
	(2,2)		{d}		{d}		{d}
	(2,1)						
	(1,5)	{a, b}, {a, c}	{a, b}, {c, d}	{a, b}, {a, c}	{a, b}, {a, c}	{a, b}, {a, c}	{a, c}, {c, d}
		{c, d}		{c, d}		{c, d}	
	(1,4)	{c}, {a, d}	{a, d}, {b, c}	{c}, {a, d}	{a, d}, {b, d}	{c}, {b, d}	{c}, {a, d}
		{b, d}	{b, d}				{b, d}
	(1,3)	{a}, {b}	{b}	{a}		{b}	
	(1,2)		{d}		{d}		{d}

(a) $\mathcal{I}_{VC}^{U, \mathcal{A}_U}$

		Vertex					
		u_1	u_2	u_3	u_4	u_5	u_6
(α, β)	(4,2)	{x, z}, {x, y}				{x, z}, {x, y}	
		{y, w}				{y, w}	
	(4,1)						
	(3,3)	{x}, {y, w}	{x, z}, {x, y}	{x}, {y, w}		{x}, {y, w}	
			{y, w}				
	(3,2)	{y, z}, {z, w}				{y, z}, {z, w}	
	(3,1)						
	(2,4)	{x}, {y, z}	{x}, {y, z}	{x}, {y, z}	{x}, {y, z}	{x}, {y, z}	{x, z}, {x, y}
		{z, w}, {y, w}	{z, w}, {y, w}	{z, w}, {y, w}	{z, w}, {y, w}	{z, w}, {y, w}	{y, z}, {z, w}
	(2,3)	{z, y, w}	{y}	{w}		{z, y, w}	{z}
	(2,2)						
	(2,1)						
	(1,5)	{x, y, z, w}	{x}, {y}	{x, y, z, w}	{x, y, z, w}	{x, y, z, w}	{x, z, w}
	(1,4)		{z}, {w}				{y}
	(1,3)						
	(1,2)						

(b) $\mathcal{I}_{VC}^{U, \mathcal{A}_L}$

Figure 5: The VC-Index \mathcal{I}_{VC}^U for G in Figure 2

PROOF. In the Lines 2-8 of Algorithm 2, checking all vertices of G takes $O(n \cdot \mu \cdot D_{max})$ time. Also, the running time of performing an online query on the intermediate subgraph G' (Line 10 in Algorithm 2) is bounded by $O(m')$. The theorem holds. \square

3.4 Index Construction

According to Definition 3.2, computing the MASs for each vertex and each possible (α, β) pair is an important step for the VC-Index construction. A basic approach is first to compute all bi-core numbers (i.e., Definition 2.3) for each attribute-induced subgraph over any possible attribute set, and then derive the MASs of VC-Index directly based on Definition 3.2. However, this approach would cost $O(\delta \cdot m \cdot (2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|}))$ time, making it difficult to apply in large-scale datasets because of its poor scalability. Actually, in the real application, for two attribute sets $A_U, A'_U \in \mathcal{A}_U$ (resp. \mathcal{A}_L) and their corresponding attribute-induced subgraphs $G[A_U, \mathcal{A}_L]$, $G[A'_U, \mathcal{A}_L]$, if $A_U = \{A'_U \cup \{a\}\}$, where a denotes an individual attribute, we observe that the bi-core numbers of most vertices in $G[A_U, \mathcal{A}_L]$ are usually equivalent to that in $G[A'_U, \mathcal{A}_L]$. Consequently, it is natural for us to consider employing the (α, β) -core maintenance technology of dynamic graphs to optimize the computation of bi-core number, thereby enhancing the construction efficiency. There are two feasible maintenance-oriented frameworks for index construction: ① bottom-up: From the subgraph induced by an individual attribute, we gradually expand it to a larger attribute-induced subgraph through insertion algorithm of maintenance; ② top-down: from the original graph, we progressively shrink it to a smaller attribute-induced subgraph through deletion algorithm of maintenance. Note that, these two approaches, in the process of computing the bi-core numbers, also need to derive the MAS of VC-Index simultaneously according to Definition 3.2. Although the time complexity for edge insertion and edge deletion is the same according to the SOTA method of (α, β) -core maintenance [31], the top-down approach evidently incurs extra time overhead. For example, given two attribute sets A_U and $A'_U = \{A_U \cup \{a\}\}$, when the bottom-up approach is applied, we would obtain the bi-core numbers of $G[A'_U, \mathcal{A}_L]$ by inserting edges that contain the attribute a into $G[A_U, \mathcal{A}_L]$ (assume the bi-core numbers of $G[A_U, \mathcal{A}_L]$ have

been computed). During the computation process, if the bi-core number of a vertex u differs between $G[A'_U, \mathcal{A}_L]$ and $G[A_U, \mathcal{A}_L]$, then A'_U should be examined to determine whether it can be added to the temporary cell in \mathcal{I}_{VC} . Otherwise, we need to do nothing according to Definition 3.1. While, when we adopt the top-down approach, we would obtain the bi-core numbers of $G[A_U, \mathcal{A}_L]$ by removing edges that contain s from $G[A'_U, \mathcal{A}_L]$ (assume the bi-core numbers of $G[A'_U, \mathcal{A}_L]$ have been computed). If a vertex u is contained within the (α, β) -core of $G[A_U, \mathcal{A}_L]$ and $G[A'_U, \mathcal{A}_L]$, and A'_U is exactly in the temporary cell $\mathcal{I}_{VC}^U(u, (\alpha, \beta))$, then we need to replace A'_U with A_U since A'_U cannot be a MAS for u and (α, β) pair. Otherwise, A_U still need to check whether it can be added to the temporary cell in \mathcal{I}_{VC} , since A_U may be a MAS according to Definition 3.2. Compared to the former approach, which only focuses on the vertices whose bi-core numbers change, the latter requires attention to all vertices, evidently increasing the additional time cost. Consequently, we will apply the bottom-up approach to construct the VC-Index in this paper.

Construction Algorithm. The pseudocode of constructing VC-Index \mathcal{I}_{VC} is presented in Algorithm 3. We first perform the (α, β) -core decomposition for the subgraph $G[\{a\}, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, \{a\}]$) induced by each individual attribute $a \in \mathcal{A}_U$ (resp. \mathcal{A}_L), and then update the temporary \mathcal{I}_{VC} (Lines 2-9). In particular, for each vertex u and its associated (α, β) -core in $G[\{a\}, \mathcal{A}_L]$, the attribute set $\{a\}$ is obviously a MAS over u and (α, β) according to Definition 3.1. We can directly add $\{a\}$ into $\mathcal{I}_{VC}^U(u, (\alpha, \beta))$ if $\{a\}$ also satisfies partial constraint in \mathcal{I}_{VC} (Lines 4-6). After that, we adopt the bottom-up approach to iteratively derive the bi-core numbers of subgraphs induced by larger attribute sets (Lines 10-19). During the computation process, the temporary index also needs to be updated simultaneously based on Definition 3.2 (Lines 14-16). Finally, we can construct the \mathcal{I}_{VC} when all possible attribute sets are verified.

Theorem 3.3. *Given an attributed bipartite graph G , the time cost of constructing the VC-Index \mathcal{I}_{VC} is bounded by $O(n \cdot D_{max}^2 \cdot \mu)$.*

PROOF. A key step in constructing \mathcal{I}_{VC} is computing all possible bi-core numbers of each vertex u in G , which would take $O(n \cdot \deg(u, G) \cdot D_{max})$ time. Given that $\deg(u, G) \leq D_{max}$, we have

Algorithm 3: Construction for VC-Index

Input: An attributed bipartite graph $G = (V, E, \mathcal{A}_U, \mathcal{A}_L)$
Output: The VC-Index \mathcal{I}_{VC} of G

```

1  $C_U, C_L, \mathcal{I}_{VC} \leftarrow \emptyset$ ;
2 foreach individual attribute  $a \in \mathcal{A}_U$  do
3    $CD[\{a\}, \mathcal{A}_L] \leftarrow$  the bi-core numbers for  $G[\{a\}, \mathcal{A}_L]$ ;
4   foreach vertex  $u \in G[\{a\}, \mathcal{A}_L]$  do
5      $B \leftarrow$  the set of bi-core numbers of  $u$  in  $G[\{a\}, \mathcal{A}_L]$ ;
6     Add  $\{a\}$  into  $\mathcal{I}_{VC}^U(u, (\alpha, \beta))$  for  $u$  and its associated  $(\alpha, \beta)$ -core if  $\{a\}$  also satisfies partial constraint in  $\mathcal{I}_{VC}$ ;
7    $C_U.push(CD[\{a\}, \mathcal{A}_L])$ ;
8 foreach individual attribute  $a' \in \mathcal{A}_L$  do
9   Do the operations that are symmetrical to those in Lines 3-7 in Algorithm 3;
10 while  $C_U \neq \emptyset$  do
11    $CD[A_U, \mathcal{A}_L] \leftarrow C_U.pop()$ ;
12   foreach individual attribute  $a \in \mathcal{A}_U$  do
13      $A'_U \leftarrow A_U \cup \{a\}$ ;
14     if the bi-core numbers of  $G[A'_U, \mathcal{A}_L]$  have not been computed then
15        $CD[A'_U, \mathcal{A}_L] \leftarrow$  compute the bi-core numbers of  $G[A'_U, \mathcal{A}_L]$  by performing an updating procedure for inserting all edges containing  $a$  into  $G[A_U, \mathcal{A}_L]$ ;
16       Update the temporary  $\mathcal{I}_{VC}$  based on Definition 3.2;
17        $C_U.push(CD[A'_U, \mathcal{A}_L])$ ;
18 while  $C_L \neq \emptyset$  do
19   Do the operations that are symmetrical to those in Lines 11-17 in Algorithm 3;
20 return  $\mathcal{I}_{VC}$ ;
  
```

$n \cdot \deg(u, G) \cdot D_{max} \leq n \cdot D_{max}^2$. During the computation of MAS, checking whether an attribute set can be added to the index requires a comparison with the temporary attribute sets in \mathcal{I}_{VC} , which takes $O(|\mathcal{I}_{VC}(u, (\alpha, \beta))|)$ time. Hence, the total time complexity is bounded by $O(n \cdot D_{max}^2 \cdot \mu)$. \square

4 ATTRIBUTE-BASED CORE INDEX

While the VC-Index can correctly answer the query of the SQAC problem, its query efficiency is limited due to the factor n in Theorem 3.2. Therefore, it is essential to eliminate the impact of graph size on query efficiency. In this section, we design a query-optimized index that reduces the query time cost from $O(n \cdot \mu \cdot D_{max} + m')$ to $O(n' + m')$, where n' and m' denote the number of vertices and edges in the intermediate subgraph, respectively.

4.1 Overview of the AC-Index

To eliminate the influence of factor n on the query efficiency of VC-Index, a natural idea is to optimize the query performance to be related to each query rather than the graph size. Considering the existing Bicore-Index [29] for the query of (α, β) -core over non-attributed bipartite graphs, its query efficiency depends solely on the size of each query result. Drawing inspiration from this, we propose a new index structure, called Atttribute-based Core Index (AC-Index, for short), which computes the Bicore-Index for the subgraphs induced by each possible attribute set and can help directly extract the intermediate subgraph.

Note that, although both AC-Index and ABi-Index compute the Bicore-Index for the attribute-induced subgraphs, ABi-Index considers the subgraph induced by two types of attribute set, leading to extremely poor scalability; in contrast, the AC-Index focuses only on the subgraph induced by one type of attribute set. Therefore, the AC-Index can essentially be viewed as a one-dimensional version of the ABi-Index. We formally define the AC-Index as follows.

DEFINITION 4.1 (AC-INDEX). Given an attributed bipartite graph $G = (V = (U, L), E, \mathcal{A}_U, \mathcal{A}_L)$, the AC-Index of G , formally denoted by \mathcal{I}_{AC} , where each sub-index $\mathcal{I}_{AC}(A)$ maintains a Bicore-Index for a subgraph $G[A, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, A]$) induced by an attribute set $A \in \mathcal{A}_U$ (resp. \mathcal{A}_L) and is composed of $\mathcal{I}_{AC}^U(A)$ and $\mathcal{I}_{AC}^L(A)$.

Theorem 4.1. Given an AC-Index \mathcal{I}_{AC} of G , the index size is bounded by $O(m \cdot (2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|}))$.

PROOF. For the subgraph induced by each possible attribute set, the size of corresponding sub-index of \mathcal{I}_{AC} is bounded by $O(m)$ according to the existing work Bicore-Index [29]. The theorem holds since there are $(2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|})$ such possible graphs. \square

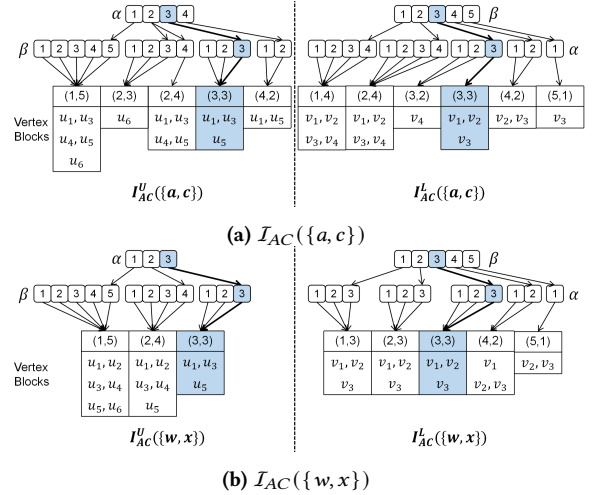


Figure 6: The partial AC-Index \mathcal{I}_{AC} of G in Figure 2

Query Algorithm. Given an AC-Index \mathcal{I}_{AC} of G , two integers α, β and two query attribute sets Q_U, Q_L . We first locate the sub-index $\mathcal{I}_{AC}(Q_U)$ and $\mathcal{I}_{AC}(Q_L)$ in \mathcal{I}_{AC} . For $\mathcal{I}_{AC}(Q_U)$, we obtain partial candidate vertices by collecting the vertices in vertex block $\mathcal{I}_{AC}(Q_U, (\alpha, \beta'))$ with $\beta' \geq \beta$. We then perform a symmetrical operation for $\mathcal{I}_{AC}(Q_L)$ and obtain the final candidate vertex set R (i.e., $C_{\alpha, \beta}(G[Q_U, \mathcal{A}_L]) \cap C_{\alpha, \beta}(G[\mathcal{A}_U, Q_L])$). Finally, for the intermediate subgraph G' induced by R , we iteratively remove the vertices that violate the degree constraints to obtain the final result.

Theorem 4.2. Given an AC-Index \mathcal{I}_{AC} of G , the time complexity of the AC-Index based query algorithm is bounded by $O(n' + m')$, where n' and m' are the number of vertices and edges in the intermediate subgraph, respectively.

PROOF. Firstly, the algorithm needs to take $O(\lambda \cdot n')$ time to extract the intermediate subgraph G' , where λ denotes the average occurrence frequency of all vertices in G' . Due to the fact that λ is a small integer in the large-scale datasets, the time cost can be

loosely bounded by $O(n')$. Additionally, the algorithm will take $O(m')$ time to peel the intermediate subgraph. We can easily infer that, the intermediate subgraph G' for both \mathcal{I}_{AC} and \mathcal{I}_{VC} over the same query attribute set (Q_U, Q_L) is identical, since G' is essentially the intersection of the (α, β) -core in $G[Q_U, \mathcal{A}_L]$ and $G[\mathcal{A}_U, Q_L]$. Thus, the theorem holds. \square

Example 4.1. Considering the attributed bipartite graph in Figure 2, a part of the \mathcal{I}_{AC} in Figure 6, two integers $\alpha = 3, \beta = 3$ and two query attribute sets $Q_U = \{a, c\}, Q_L = \{w, x\}$, the intermediate graph G' retrieved by $\mathcal{I}_{AC}(\{a, c\})$ and $\mathcal{I}_{AC}(\{w, x\})$, is $G' = \{U = \{u_1, u_3, u_5\}, L = \{v_1, v_2, v_3\}\}$. The corresponding result is marked in blue. Note that G' is also the answer, since all vertices in G' satisfy the degree constraints.

Construction Algorithm. According to Definition 4.1, a basic approach for AC-Index construction is to invoke the Bicore-Index construction procedure for each subgraph induced by any possible attribute set. However, this approach would cost $O(\delta \cdot m \cdot (2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|}))$ time, making it difficult to apply to large graphs due to its poor scalability. Given that the computation of the bi-core number is an essential step for deriving the Bicore-Index, re-computing it for each possible subgraph is a significant waste of time. Therefore, we also consider maintenance-oriented technology (mentioned in Section 3.4) to optimize the computation of bi-core number. Due to the factor that the time complexity for edge insertion and edge deletion in the maintenance algorithm is the same, we employ the bottom-up framework to construct AC-Index in this paper.

The pseudocode for AC-Index \mathcal{I}_{AC} construction is presented in Algorithm 4. We first construct the Bicore-Index for the subgraph $G[\{a\}, \mathcal{A}_L]$ induced by each individual attribute $a \in \mathcal{A}_U$ (Lines 3-5). Based on these smaller subgraphs that have been computed, we then perform insertion algorithm of (α, β) -core maintenance to compute the bi-core numbers for larger subgraphs and derive the corresponding Bicore-Index (Lines 6-13). Finally, to obtain a full \mathcal{I}_{AC} , we perform symmetrical operation for the subgraph induced by the attribute set $A_L \in \mathcal{A}_L$ (Lines 14-15).

Theorem 4.3. Given an attributed bipartite graph G , the time cost of constructing AC-Index \mathcal{I}_{AC} is bounded by $O(n \cdot D_{max}^2 \cdot \bar{d} \cdot \log \gamma)$, where \bar{d} denotes the average degree of each vertex in G , and γ denotes the average number of vertices for each vertex block in AC-Index.

PROOF. A key step in constructing \mathcal{I}_{AC} is computing all possible bi-core numbers of each vertex u in G , which would take $O(n \cdot \deg(u, G) \cdot D_{max})$ time. Given that $\deg(u, G) \leq D_{max}$, we have $n \cdot \deg(u, G) \cdot D_{max} \leq n \cdot D_{max}^2$. During the computation of bi-core numbers, if the bi-core number of a vertex changes, it would take $O(\bar{d} \cdot \log \gamma)$ time to update the temporary \mathcal{I}_{AC} . Hence, the total time cost is bounded by $O(n \cdot D_{max}^2 \cdot \bar{d} \cdot \log \gamma)$. \square

5 MINIMIZED ATTRIBUTE-BASED CORE INDEX

Although the AC-Index achieves a great improvement of query performance from $O(n \cdot \mu \cdot D_{max} + m')$ to $O(n' + m')$, it still suffers from poor scalability in terms of index size, which limits its applicability to large-scale datasets. Consequently, it is necessary to further compress redundant information in the AC-Index. In this section, we

Algorithm 4: Construction for AC-Index

Input: An attributed bipartite graph $G = (V, E, \mathcal{A}_U, \mathcal{A}_L)$
Output: The AC-Index \mathcal{I}_{AC} of G

```

1  $C_U, C_L, \mathcal{I}_{AC} \leftarrow \emptyset$ ;
2 foreach individual attribute  $a \in \mathcal{A}_U$  do
3    $CD[\{a\}, \mathcal{A}_L] \leftarrow$  the bi-core numbers for  $G[\{a\}, \mathcal{A}_L]$ ;
4   Compute sub-index  $\mathcal{I}_{AC}(\{a\})$  for the subgraph  $G[\{a\}, \mathcal{A}_L]$ ;
5    $C_U.push(CD[\{a\}, \mathcal{A}_L])$ ;
6   while  $C_U \neq \emptyset$  do
7      $CD[A_U, \mathcal{A}_L] \leftarrow C_U.pop$ ;
8     foreach individual attribute  $a' \in \mathcal{A}_U \wedge a' \neq a$  do
9        $A'_U \leftarrow A_U \cup \{a\}$ ;
10      if the bi-core numbers of  $G[A'_U, \mathcal{A}_L]$  have not been
          computed then
11         $CD[A'_U, \mathcal{A}_L] \leftarrow$  perform updating procedure by
            inserting edges containing  $a'$  into  $G[A_U, \mathcal{A}_L]$ ;
12        Compute  $\mathcal{I}_{AC}(A'_U)$  for subgraph  $G[A'_U, \mathcal{A}_L]$ ;
13         $C_U.push(CD[A'_U, \mathcal{A}_L])$ ;
14 foreach individual attribute  $a \in \mathcal{A}_L$  do
15   Do the operations that are symmetrical to those in Lines 3-13 in
      Algorithm 4;
16 return  $\mathcal{I}_{AC}$ ;
```

demonstrate the compressibility of vertex blocks in the AC-Index, proposing the minimized AC-Index, which not only significantly optimizes the space cost of AC-Index from $O(m \cdot (2^{|\mathcal{A}_U|} + 2^{|\mathcal{A}_L|}))$ to $O(\rho \cdot n)$, but also ensures query efficiency comparable to the AC-Index. Here, ρ denotes the average occurrence frequency of all vertices in G (ρ is a small integer and $\rho \ll n$ in practice).

5.1 Overview of the MAC-Index

In the AC-Index \mathcal{I}_{AC} , we observe that a vertex u may be stored in different sub-indexes of \mathcal{I}_{AC} for the same (α, β) , which would result in significant storage redundancy in \mathcal{I}_{AC} according to the containment relationship of the attribute set for the (α, β) -core. Next, we will introduce the concept of the attribute-dominant vertex and illustrate how such a vertex utilizes the dominant property to compress redundant information in \mathcal{I}_{AC} . Finally, we propose the minimized AC-Index based on these properties.

DEFINITION 5.1 (ATTRIBUTE-DOMINANT VERTEX). Given an AC-Index \mathcal{I}_{AC} of G , for any vertex $u \in \mathcal{I}_{AC}^U(A, (\alpha, \beta))$, if there does not exist any other vertex block $\mathcal{I}_{AC}^U(A', (\alpha, \beta))$ containing u with $A' \subset A$, then the vertex u is an attribute-dominant vertex over the attribute set A and the integer pair (α, β) . The symmetric case for the vertices in any $\mathcal{I}_{AC}^L(A, (\alpha, \beta))$.

Example 5.1. Considering the AC-Index \mathcal{I}_{AC} of the attributed bipartite graph in Figure 2. The vertex u_3 is contained in $\mathcal{I}_{AC}^U(\{c\}, (3, 3))$, $\mathcal{I}_{AC}^U(\{a, c\}, (3, 3))$, and so on, then u_3 is an attribute-dominant vertex over $\{c\}$ and $(3, 3)$.

According to Definition 5.1, we can easily infer that, for a vertex $u \in U(G)$, if u is an attribute-dominant vertex in the vertex block $\mathcal{I}_{AC}^U(A, (\alpha, \beta))$, then u must also be contained in the (α, β) -core of any subgraph $G[A', \mathcal{A}_L]$ with $A \subseteq A'$. Consequently, when querying the (α, β) -core of the subgraph $G[A, \mathcal{A}_L]$, we can also

retrieve the result vertices $U(G')$ (resp. $L(G')$) from the attribute-dominant vertices in any vertex block $I_{AC}^U(A', (\alpha, \beta'))$ with $A' \subseteq A$ and $\beta \leq \beta'$ (resp. $I_{AC}^L(A', (\alpha', \beta))$ with $A' \subseteq A$ and $\alpha \leq \alpha'$).

With the dominant property among the vertices of different vertex blocks in I_{AC} , we construct a superior-optimized index, named Minimized Attribute-based Core Index (MAC-Index, for short), which significantly compresses the space of I_{AC} , while efficient query performance is guaranteed. We formally define the MAC-Index as follows.

DEFINITION 5.2 (MAC-INDEX). *Given an attributed bipartite graph $G = (V = (U, L), E, \mathcal{A}_U, \mathcal{A}_L)$, the MAC-Index of G , formally denoted by I_{MAC} , which just maintains the attribute-dominant vertices for each sub-index in the AC-Index, i.e., for each vertex u in any vertex block $I_{MAC}^U(A, (\alpha, \beta))$ (resp. $I_{MAC}^L(A, (\alpha, \beta))$), u is an attribute-dominant vertex over u and (α, β) .*

For convenience, we use $I_{MAC}(A)$, $I_{MAC}^U(A)$, and $I_{MAC}^L(A)$ to correspond to $I_{AC}(A)$, $I_{AC}^U(A)$ and $I_{AC}^L(A)$ of the AC-Index, respectively. We also use $I_{MAC}^{\mathcal{A}_U}$ (resp. $I_{MAC}^{\mathcal{A}_L}$) to denote the set of $I_{MAC}^U(A)$ (resp. $I_{MAC}^L(A)$) with $A \subseteq \mathcal{A}$. Additionally, we use $I_{MAC}^{\mathcal{A}}$ to denote the set of $I_{MAC}(A)$ with $A \in \mathcal{A}$.

Example 5.2. *Considering the attributed bipartite graph in Figure 2, a part of the MAC-Index I_{MAC} is presented in Figure 7 due to the space limit. Compared to the I_{AC} in Figure 6, we can see that, for the same attribute set and (α, β) pair, the corresponding vertex block in I_{MAC} contains fewer vertices. For example, $I_{AC}^U(\{a, c\}, (3, 3)) = \{u_1, u_3, u_5\}$ and $I_{MAC}^U(\{a, c\}, (3, 3)) = \emptyset$.*

Theorem 5.1. *Given a MAC-Index I_{MAC} of G , the index size is bounded by $O(\rho \cdot n)$, where ρ denotes the average occurrence frequency of all vertices in G within I_{MAC} .*

PROOF. The primary space cost of the I_{MAC} stems from storing the vertices of the vertex blocks. For convenience, we use $|I_{MAC}^U(A)|$ to denote the total number of vertices in all vertex blocks of $I_{MAC}^U(A)$. The total number of vertices in I_{MAC} can be represented as $\hat{n} = \sum_{A \in \{\mathcal{A}_U \cup \mathcal{A}_L\}} (|I_{MAC}^U(A)| + |I_{MAC}^L(A)|)$. Thus, we have $\rho = \frac{\hat{n}}{n}$. The theorem holds. \square

Query Algorithm. Let's discuss the MAC-Index based query algorithm, which is supported by the following lemma.

Lemma 5.1. *Given a MAC-Index I_{MAC} of an attributed bipartite graph G , two integers α, β and two query attribute sets Q_U, Q_L , for any vertex $u \in I_{MAC}^U(A, (\alpha, \beta'))$ (resp. $I_{MAC}^L(A, (\alpha', \beta))$) with $A \subseteq Q_U$ and $\beta' \geq \beta$ (resp. $\alpha' \geq \alpha$), the vertex u must be contained in the (α, β) -core of $G[Q_U, \mathcal{A}_L]$. The symmetric case for the vertices in $I_{MAC}^U(A', (\alpha, \beta'))$ with $A' \subseteq Q_L$.*

PROOF. This can be directly derived from Definition 5.1. \square

The pseudocode for MAC-Index based query algorithm is presented in Algorithm 5. Given a MAC-Index I_{AC} of G , two integers α, β and two query attribute sets Q_U, Q_L . In Line 4-8, we first collect the candidate vertices from the vertex blocks $I_{MAC}^U(A, (\alpha, \beta'))$ (resp. $I_{MAC}^L(A, (\alpha', \beta))$) with $A \in Q_U$ and $\beta' \geq \beta$ (resp. $\alpha' \geq \alpha$) since these vertices may be the part of final result according to

Algorithm 5: Query based on MAC-Index

Input: An attributed bipartite graph G , a MAC-Index I_{MAC} of G , two integers α, β and two query attribute sets Q_U, Q_L
Output: The result of SQAC

- 1 $R \leftarrow \emptyset$;
- 2 $Q_U \leftarrow$ All possible subsets of Q_U ;
- 3 $Q_L \leftarrow$ All possible subsets of Q_L ;
- 4 **foreach** attribute set $A \in Q_U$ **do**
- 5 **foreach** integer $\beta' \geq \beta$ **do**
- 6 **if** $I_{MAC}^U(A, (\alpha, \beta'))$ exists **then**
- 7 $R \leftarrow R \cup I_{MAC}^U(A, (\alpha, \beta'))$;
- 8 For any $I_{MAC}^L(A, (\alpha', \beta))$ with $\alpha' \geq \alpha$, we do the operations that are symmetrical to those in Line 5-7 in Algorithm 5;
- 9 **foreach** attribute set $A \subseteq Q_L$ **do**
- 10 Do the same operations as Lines 5-8 in Algorithm 5.
- 11 $G' \leftarrow$ the subgraph of G induced by R ;
- 12 Do the same operation as Lines 2-13 in Algorithm 1 to guarantee the correctness of the final results;
- 13 **return** G' ;

Lemma 5.1. We then collect the remaining candidate vertices from $I_{MAC}^{\mathcal{A}_L}$ in Line 8-10. Finally, we remove the vertices that do not satisfy degree constraints to obtain the final result. The overall time cost is $O(\lambda' \cdot n' + m')$, where λ' denotes the average occurrence frequency of all vertices in G' . In fact, λ' is a small integer and $\lambda' \ll n'$. Thus, the time cost can be loosely bounded by $O(n' + m')$.

5.2 The Construction for MAC-Index

Computing the attribute-dominant vertices for a specific attribute set and an integer pair (α, β) is a crucial step in constructing the MAC-Index, significantly reducing the size of the index. In this subsection, we will illustrate the construction framework of the MAC-Index and the computation of the attribute-dominant vertex.

Algorithm 6: Construction for MAC-Index

Input: An attributed bipartite graph $G = (V, E, \mathcal{A}_U, \mathcal{A}_L)$
Output: The MAC-Index I_{MAC} of G

- 1 $C_U, C_L, I_{MAC} \leftarrow \emptyset$;
- 2 **foreach** individual attribute $a \in \mathcal{A}_U$ **do**
- 3 $CD[\{a\}, \mathcal{A}_L] \leftarrow$ compute bi-core numbers for $G[\{a\}, \mathcal{A}_L]$;
- 4 Construct the $I_{MAC}^U(\{a\}, (\alpha, \beta))$ (resp. $I_{MAC}^L(\{a\}, (\alpha, \beta))$) of the I_{MAC} according to Definition 5.1;
- 5 $C_U.push(CD[\{a\}, \mathcal{A}_L])$;
- 6 **foreach** individual attribute $a' \in \mathcal{A}_L$ **do**
- 7 Do the operations that are symmetrical to those in Lines 3-5 in Algorithm 6;
- 8 Call the procedure *ComputeADVertex* (G, C_U, C_L, I_{MAC}) to compute all attribute-dominant vertices;

The Framework. The pseudocode for MAC-Index construction is presented in Algorithm 6. Firstly, we compute the bi-core numbers for each subgraph $G[\{a\}, \mathcal{A}_L]$ induced by the individual attribute $a \in \mathcal{A}_U$ and construct the vertex block $I_{MAC}^U(\{a\}, (\alpha, \beta))$ (resp. $I_{MAC}^L(\{a\}, (\alpha, \beta))$) over a and each possible (α, β) pair according to Definition 5.1 (Lines 2-5). Actually, the vertex blocks in $I_{MAC}(\{a\})$ are identical to those in $I_{AC}(\{a\})$ since $\{a\}$ can only be a subset of other sets and each vertex in $I_{AC}(\{a\})$ is an attribute-dominant

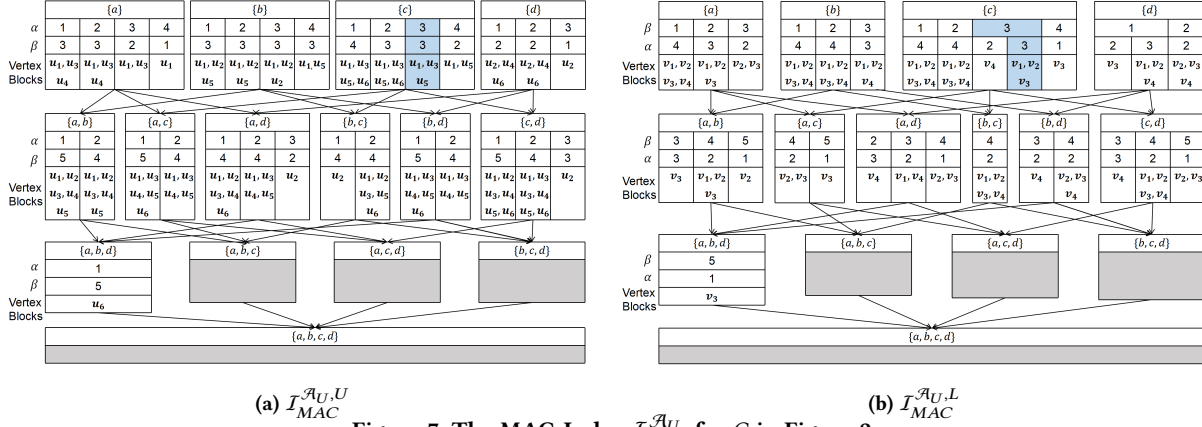


Figure 7: The MAC-Index $\mathcal{I}_{MAC}^{\mathcal{A}_U}$ for G in Figure 2

Algorithm 7: Attribute-dominant Vertex Computation

```

1 Procedure ComputeADVertex( $G, C_U, C_L, \mathcal{I}_{MAC}$ )
2   while  $C_U \neq \emptyset$  do
3      $CD[A_U, \mathcal{A}_L] \leftarrow C_U.pop()$ ;
4     foreach individual attribute  $a \in \mathcal{A}_U$  do
5        $A'_U \leftarrow A_U \cup \{a\}$ ;
6       if the bi-core numbers of  $G[A'_U, \mathcal{A}_L]$  have not been
          computed then
7          $V^* \leftarrow$  Vertices with bi-core number changed by
          inserting edges containing  $a$  into  $G[A_U, \mathcal{A}_L]$ ;
8         foreach vertex  $u \in V^* \wedge u \in U(G)$  do
9            $(\alpha, \beta) \leftarrow$  the current bi-core number of  $u$ ;
10          foreach integer  $\alpha' \leq \alpha$  do
11            if  $u$  is an attribute-dominant vertex for
               $A'_U$  and  $(\alpha', \beta)$  based on Definition 5.1
              then
12              Add  $u$  to  $\mathcal{I}_{MAC}^U(A'_U, (\alpha', \beta))$ ;
13          foreach vertex  $v \in V^* \wedge v \in L(G)$  do
14            Do the operation that are symmetrical to
              those in Lines 6-12 in Algorithm 7;
15           $C_U.push(CD[A'_U, \mathcal{A}_L])$ ;
16   while  $C_L \neq \emptyset$  do
17     Do the operations that are symmetrical to those in
        Lines 3-15 of Algorithm 7;
18   return  $\mathcal{I}_{MAC}$ ;

```

vertex. We then perform symmetrical operations for each individual attribute $a' \in \mathcal{A}_L$ in Lines 6-7. Finally, we compute the remaining attribute-dominant vertices by invoking the procedure *ComputeADVertex* in Algorithm 7.

The Computation of Attribute-dominant Vertex. Similar to the AC-Index construction, we also employ the bottom-up approach to effectively derive the bi-core numbers of subgraph induced by larger attribute sets and simultaneously update the temporary MAC-Index. Specifically, during the process of computing the bi-core numbers of an attribute-induced subgraph $G[A'_U, \mathcal{A}_L]$, for a vertex $u \in U(G)$ (resp. $L(G)$) with a bi-core number changed (assume the current bi-core number is (α, β)), if u is an attribute-dominant vertex for A'_U and (α', β) pair with $\alpha' \leq \alpha$ based on Definition 5.1, then u would be added to the vertex block $\mathcal{I}_{MAC}^U(A'_U, (\alpha', \beta))$ (Lines 6-12 in Algorithm 7), which costs $O(\bar{d} \cdot \log n)$ time. Thus, the time

complexity of constructing \mathcal{I}_{MAC} is bounded by $O(n \cdot D_{max}^2 \cdot \bar{d} \cdot \log n)$ due to the fact that there are $|deg(u, G) \cdot D_{max}|$ possible bi-core numbers for each vertex u in G .

5.3 Discussions on Maintenance of MAC-Index

Considering that the bipartite graph-structured data is constantly changing in the real world, forming the dynamic bipartite graphs. In this subsection, we propose the maintenance algorithm for the MAC-Index including edge insertion and edge deletion.

When edge insertion happens, assume the incoming edge is $e_i = (u_i, v_i)$, with $A(u_i)$ and $A(v_i)$ denoting the set of attributes associated with u_i and v_i , respectively. Obviously, we only need to consider the affected subgraph $G[A, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, A]$) induced by the attribute set A containing at least one attribute from $A(u_i)$ (resp. $A(v_i)$), i.e., $A \cap A(u_i) \neq \emptyset$. This is due to the bi-core numbers for the subgraphs induced by other attribute sets remaining unchanged. Next, with the properties of attribute-dominant vertices, we will still adopt a bottom-up strategy to update the vertex blocks in each affected subgraph. From the smaller affected subgraph $G[A, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, A]$), we first call the insertion algorithm of (α, β) -core maintenance [31] to obtain the vertices V^* whose bi-core number changes. Then, for a vertex $u \in \{V^* \cap U(G)\}$ (resp. $v \in \{V^* \cap V(G)\}$) with the current bi-core number (α, β) after edge insertion, if u is an attribute-dominant vertex for A and (α', β) pair with $\alpha' \leq \alpha$, then we would add u to the vertex block $\mathcal{I}_{MAC}^U(A, (\alpha', \beta))$. Meanwhile, if u is exactly included in $\mathcal{I}_{MAC}^U(A, (\alpha', \beta'))$ with $\beta' < \beta$ before, then u needs to be removed based on Definition 5.2. For edge deletion, we have operations similar to those for edge insertion. When deleting an edge $e_j = (u_j, v_j)$, we also need to focus on the affected subgraph as in the case of edge insertion and employ the bottom-up approach to update the MAC-Index. From the smaller affected subgraph $G[A, \mathcal{A}_L]$ (resp. $G[\mathcal{A}_U, A]$), we will first call the deletion algorithm of (α, β) -core maintenance to obtain the vertices V^* whose bi-core number changes. Then, for a vertex $u \in \{V^* \cap U(G)\}$ (resp. $v \in \{V^* \cap V(G)\}$) with a previous bi-core number (α', β') and a current bi-core number (α, β) after edge deletion, we would remove u from previous $\mathcal{I}_{MAC}^U(A, (\alpha', \beta'))$ according to the definition of Bicore-Index. After that, if u is an attribute-dominant vertex for A and (α'', β) pair with $\alpha'' \leq \alpha$, then we would add u to the vertex block $\mathcal{I}_{MAC}^U(A, (\alpha'', \beta))$. Meanwhile, if u is exactly

included in $\mathcal{I}_{MAC}^U(A, (\alpha'', \beta''))$ with $\beta'' > \beta$ before, then u should be removed based on Definition 5.2.

Table 2: Statistics of Datasets

Datasets	$ E $	$ U $	$ L $	$ \mathcal{A}_U $	$ \mathcal{A}_L $	d_{max}^U	d_{max}^L	Real Labels
IMDB-5M	5.46M	1.05M	2.13M	6	6	359	27	✓
IMDB-7M	7.66M	1.02M	2.34M	6	6	341	18	✓
IMDB-10M	10.53M	1.05M	2.57M	5	7	611	45	✓
Tmall-3M	3.56M	1.07M	0.31M	4	8	128	158	✓
Tmall-10M	10.47M	2.06M	0.69M	4	9	101	68	✓
DBLP	12.28M	1.95M	5.62M	6	6	1386	287	✗

6 EXPERIMENTAL EVALUATION

In this section, we evaluate the performance of the proposed algorithms. We also conduct two case studies for our proposed SQAC. All algorithms are implemented in C++ and run on a CentOS machine of 1TB memory and an Intel(R) Xeon(R) Silver-4210R 2.40GHz CPU. We terminate an algorithm if it runs more than 72 hours, which is denoted as INF. Note that the ABi-Index cannot be constructed within 72 hours over all datasets.

6.1 Setup

6.1.1 Datasets. Table 2 presents the statistics of datasets used in the experiments. We extract three real graphs from IMDB (<https://www.imdb.com>) and two real graphs from Tmall (<https://www.tmall.com/>). The **IMDB** dataset consists of a bipartite person-movie network, where an edge denotes that a person has contributed to a movie. The attributes of a person indicate occupation (e.g., "actor", "producer"), and the attributes associated with a movie represent its genre. The **Tmall** dataset is an E-commerce website that contains the customer-product bipartite network, where the attribute of a customer indicates behavior (e.g., "click", "purchase") and the attributes associated with a product represent its category. The **DBLP** dataset is obtained from KONECT (<http://konect.cc/networks/>) and lacks vertex attributes; we generate synthetic attributes as in previous work [46], following a specific uniform distribution. All codes and datasets are available on GitHub [1].

6.1.2 Parameters setting. For the default query, we set the size of the query attribute set $|Q_U| = 0.4 \cdot |\mathcal{A}_U|$ and $Q_L = 0.4 \cdot |\mathcal{A}_L|$. Also, we set the query integer $\alpha = 0.1 \cdot d_{max}^U$ and $\beta = 0.1 \cdot d_{max}^L$. To evaluate the effect of α and β , we set α from $0.05 \cdot d_{max}^U$ to $0.9 \cdot d_{max}^U$ and β from $0.05 \cdot d_{max}^L$ to $0.9 \cdot d_{max}^L$. For evaluating the effect of the size of query attribute set, we vary $|Q_U|$ as 20%, 40%, 60%, 80% of $|\mathcal{A}_U|$ and $|Q_L|$ as 20%, 40%, 60%, 80% of $|\mathcal{A}_L|$. The reported time and space under a given group of settings are obtained by averaging those from the corresponding generated queries.

6.2 Query Efficiency

6.2.1 Varying α and β . We evaluate the query performance by varying the value of α and β over all datasets. Figure 8 shows that all index-based query algorithms significantly outperform the online algorithm in all settings. Also, we can observe that both AC-Index and MAC-Index demonstrate superior query performance compared to the VC-Index. Additionally, the query efficiency of AC-Index and MAC-Index is nearly identical. With the growth of α and β , the performance of the online algorithm is relatively stable, since the online algorithm always scans the entire original graph

for each query, regardless of the value of α and β . For the VC-Index, it exhibits a downward trend in running time as the values of α and β increase, since the number of MASs for a vertex u and (α, β) will evidently decrease, resulting in lower time costs to check whether a vertex is in the (α, β) -core. Regarding the AC-Index and MAC-Index, they also show better query efficiency when the values of α and β get larger, which is primarily attributed to the size of the query result decreasing significantly as α and β grow, allowing them to quickly retrieve relevant vertices.

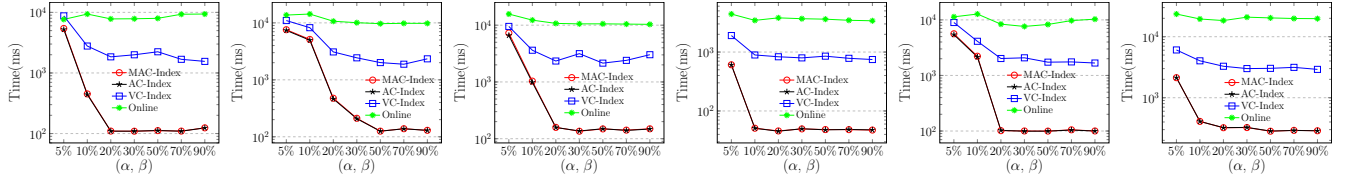
6.2.2 Varying the size of query attribute set. Figure 9 presents the query performance by varying the size of the query attribute set. We can see that, our all index-based query algorithms consistently outperform online algorithm by 1~2 orders of magnitude. In addition, both MAC-Index and AC-Index show better query efficiency than the VC-Index. When the size of the query attribute set increases, all index-based query algorithms show a slight upward trend on time cost, which is primarily due to the fact that the size of the intermediate subgraph increases, allowing more time overhead for the peeling process. In contrast, online query is not sensitive to the size of the query attribute set due to the peeling of the entire original graph from scratch.

6.3 Index Construction & Maintenance

6.3.1 Construction performance on all datasets. In Figure 12, we evaluate the construction performance for all indexes. We can observe that the construction efficiency among these indexes is very marginal, primarily because they all employ the bottom-up construction framework. Moreover, the construction time of the MAC-Index is slightly higher than that of AC-Index, which is also reasonable due to the computation of attribute-dominant vertex.

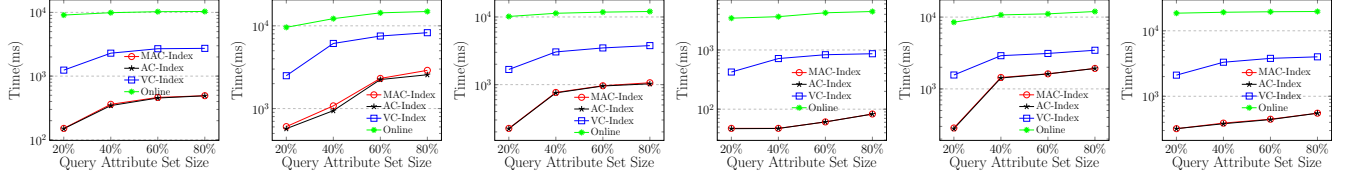
6.3.2 Construction performance varying graph size. We evaluate the scalability of our construction methods in Figure 10. For each dataset, we randomly sample 20%, 40%, 60%, 80%, 100% of the edges from the original graph to perform the construction algorithm. We can see that the running time of all construction methods shows an increasing trend as the graph size grows, since these methods rely on the graph size and the number of attribute types in the datasets. In addition, the construction time of all methods is comparable, which corresponds to the results shown in Figure 12.

6.3.3 Maintenance performance on all datasets. We report the maintenance performance of MAC-Index in Figure 14. Note that there are three algorithms for the maintenance of MAC-Index. The first is our proposed maintenance algorithm for edge insertion, denoted as **MAC-Ins**. The second is our proposed maintenance algorithm for edge deletion, denoted as **MAC-Del**. The last is the index construction algorithm adapted for the maintenance of MAC-Index, denoted as **MAC-Constr**. When evaluating the efficiency of edge deletion, we randomly remove 1000 edges from the input graph and report the average time required for each edge deletion. Then, for testing the edge insertion, we insert these removed edges into the graph and report the average time required for each edge insertion. From the result shown in Figure 14, we can see that the running time of our proposed maintenance algorithm is significantly faster than the basic solution. The primary reason is that we only perform



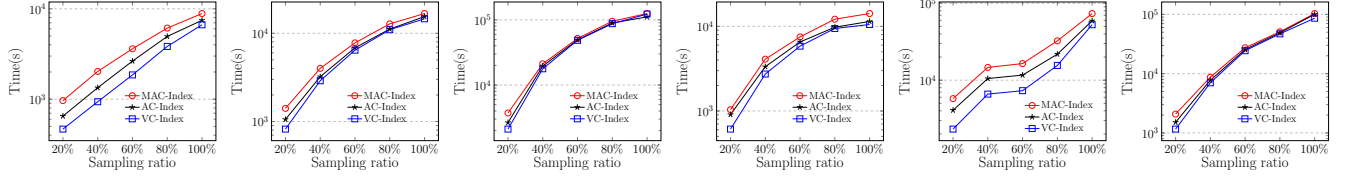
(a) IMDB-5M (b) IMDB-7M (c) IMDB-10M (d) Tmall-3M (e) Tmall-10M (f) DBLP

Figure 8: Query time varying α and β



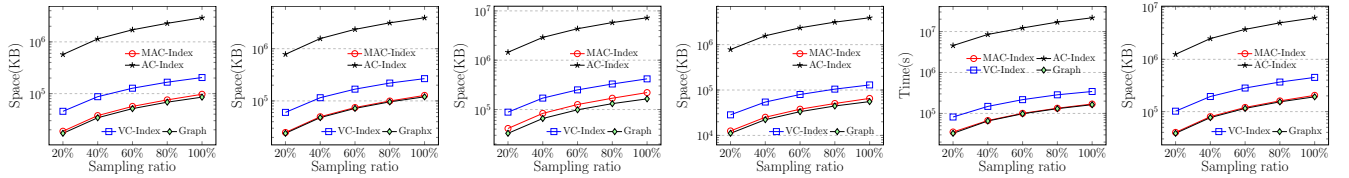
(a) IMDB-5M (b) IMDB-7M (c) IMDB-10M (d) Tmall-3M (e) Tmall-10M (f) DBLP

Figure 9: Query time varying query attribute set size



(a) IMDB-5M (b) IMDB-7M (c) IMDB-10M (d) Tmall-3M (e) Tmall-10M (f) DBLP

Figure 10: Construction performance for different sampling ratio



(a) IMDB-5M (b) IMDB-7M (c) IMDB-10M (d) Tmall-3M (e) Tmall-10M (f) DBLP

Figure 11: Index size for different sampling ratio

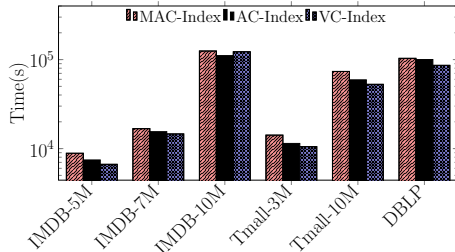


Figure 12: Construction performance on all datasets

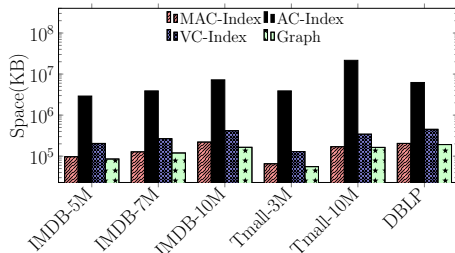


Figure 13: Index size on all datasets

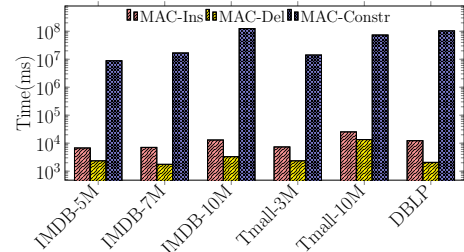


Figure 14: Maintenance performance on all datasets

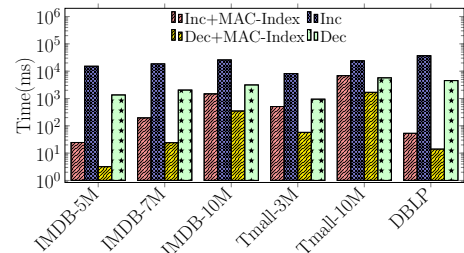


Figure 15: Query performance of attributed communities

a localized update on the index during the edge insertion/deletion happens, thus avoiding unnecessary time overhead.

6.4 Index Size

6.4.1 Index size on all datasets. We compare the space overhead of different indexes and choose the size of the original graph as baseline. We can see from the result in Figure 13 that the space efficiency of MAC-Index significantly outperforms the comparative indexes and is comparable to that of the original graph. Apparently, the size of AC-Index is the largest among all indexes since it almost stores all query results. In addition, the space efficiency of VC-Index is better than that of AC-Index, as its size primarily depends on the graph size and the average number of MASs.

6.4.2 Index size varying graph size. We also evaluate the scalability of the index size by varying the graph size. For each dataset, we adopt the same sampling strategy as presented in Figure 10. We can see from Figure 11 that the space efficiency of MAC-Index significantly outperforms other indexes in all settings, demonstrating superior scalability. As the graph size grows, the size of all indexes becomes larger, since they are dependent on the graph size and the number of attribute types.



Figure 16: Case study on IMDB graph ((7, 11)-core)

6.5 Case Studies

To demonstrate the effectiveness of the SQAC problem in real-world applications, we conduct two case studies as follows:

- (1) Personalized recommendation. We conduct a case study on the IMDB graph of the famous actor "Evgeniy Stychkin". According to the acting style of "Evgeniy Stychkin", we choose two groups of query attribute set: {"History", "Biography"} and {"Drama", "Romance"}. Figure 16 presents (7, 11)-cores of "Evgeniy Stychkin"s ego network on the IMDB graph under two different query attribute sets. We can see that the resulting subgraphs show significant differences between the two query attribute sets.
- (2) Accelerate the querying of the attributed communities. We evaluate the query performance of the attributed (α, β) -community based on our MAC-Index. Xu et al. [46] proposed two SOTA algorithms, namely Inc and Dec, for the computation of the attributed (α, β) -community over attributed bipartite graphs. The Inc method has been thoroughly discussed in Section 1. For the Dec method, the main idea is to start from a larger attribute-induced subgraph and check whether the final community can be found. If not, it progressively narrows down to smaller attribute-induced graphs until the answer is discovered. The optimized algorithm employs MAC-Index to prune all vertices that cannot be in the result of the SQAC under the query attribute sets,

which can greatly narrow the search scope before computing the community. Regarding the query parameters, we set the default values of α and β to 8. Also, a query attribute set is configured as the set of attributes contained in the query vertex q , and another query attribute set is configured as the set of attribute types among the neighbors of q . The experiment result is presented in Figure 15. We can see that the time efficiency of the two algorithms has improved significantly when incorporating our proposed SQAC, demonstrating the effectiveness of pruning.

7 RELATED WORK

In this section, we review several closely related works to our SQAC problem as follows.

Dense Subgraph models Over Bipartite Graphs. There are many dense subgraph models in bipartite graphs, including (α, β) -core [6, 9, 23, 29–31], k -bitruss [36, 41, 49], and biclique [14, 16, 20, 32, 33, 44], which have attracted much attention in recent years. The (α, β) -core is the extension of the k -core [17, 35, 37, 47] in general graph. In [3], Ahmed et al. first introduce the concept of (α, β) -core. Ding et al. [21] propose a linear-time online approach for the computation of (α, β) -core, which iteratively removes the vertices that violate the degree constraints over bipartite graphs. The SOTA work is proposed by [29], which utilizes the Bicore-Index to efficiently retrieve the (α, β) -core. In addition, Huang et al. [26] propose a parallel version and Liu et al. [30] design an (α, β) -core decomposition algorithm in distributed settings. Luo et al. [31] propose the SOTA algorithm for the maintenance of (α, β) -core over dynamic bipartite graphs. As another significant dense subgraph structure on bipartite graphs, the k -bitruss is derived by extending the k -truss [4, 13, 25, 39] to the bipartite graph. Note that, for the truss-like structure, it focuses on edge-centric cohesiveness, while the core-like structure is used to measure the cohesiveness of vertex. Zou et al. [49] are the first to formulate the k -bitruss model in bipartite graphs. Wang et al. [41] propose the SOTA algorithm for the computation of k -bitruss, which utilizes the advance index to accelerate the query efficiency. An distributed algorithm for k -bitruss decomposition is proposed by [45]. Sariyuce et al. [36] define a novel dense subgraph model based on k -bitruss. Like clique-like structure [10–12, 42] in general graphs, biclique is a complete subgraph in a bipartite graph. Chen et al. [14] propose a SOTA algorithm for maximal biclique enumeration over bipartite graphs, which has been proven to be NP-Hard [33]. Chen et al. [15] design an exact algorithm to compute the maximum balanced biclique, which is also proven to be NP-Hard [44]. Deng et al. [20] propose an efficient algorithm for the maintenance of top k (p, q) -bicliques over streaming bipartite graphs. However, all of these existing works do not consider the attributes of vertices and cannot be effectively applied to our problem.

Attributed Core-like Subgraph Query. Recently, many different types of attributed core-like subgraph models have been extensively studied over the past decade. In general graph, Fang et al. [22] propose a community model based on k -core over vertex-attributed graphs, which aims to find a connected subgraph that meets both structure cohesiveness and attribute cohesiveness (i.e., its vertices share the most attributes of given query attribute set). Deng et al. [19] propose an efficient algorithm to identify a k -core subgraph

that incorporates edge-labeled constraints. Kang et al. [27] define a k -core-based community model over edge-attributed graphs. While, these studies over k -core structure are inherently difficult to adapt to our problem of (α, β) -core. In bipartite graph, there are several works that focus on the weights of edges. Wang et al. [43] propose a significant community model based on (α, β) -core structure that adapts minimum edge weight to measure the significance of community over edge-weighted bipartite graphs. Li et al. [28] study a more densely (α, β) -attributed weight community model, which considers both weight constraints on the edges and attribute constraints on the vertices in bipartite graphs. These works are fundamentally different from our focus on vertex-attributed bipartite graphs. There are two existing studies that also consider vertex-attributed bipartite graphs. Xu et al. [46] study the problem of attributed (α, β) -community search over vertex-attributed bipartite graphs, which is a simple extension of [22] from general graphs to bipartite graphs. Zhang et al. [48] formulate a novel community model named Pareto-optimal (α, β) -community, which is the first to consider both structure cohesiveness and vertex importance on vertex-weighted bipartite graphs. In contrast to these works, our problem just concentrates on (α, β) -core itself and relaxes the strict attribute constraint on the vertices. Consequently, these existing methods can hardly contribute to our problem.

We can see that, we are the first to propose the SQAC problem, and existing methods are far from solving this problem.

8 CONCLUSIONS

In this paper, we investigated the problem of querying the (α, β) -core with attribute constraints over attributed bipartite graphs. We are the first to propose the problem and existing studies over (α, β) -core disregard the vertex attributes, making it difficult for previous methods to adapt well to our problem. We proposed two index-based algorithms to support efficient querying for our SQAC. To further improve performance, we designed the MAC-Index, which not only significantly reduces the index size, but also guarantees efficient query efficiency. Additionally, efficient construction and maintenance algorithms were also proposed for MAC-Index. Extensive experiments over real-world datasets confirmed the efficiency and effectiveness of our index-based algorithms.

REFERENCES

- [1] 2025. Codes. <https://github.com/XueQingSBQ/SQAC>.
- [2] 2025. SQAC: Scalable Querying of Attribute-Constrained (α, β) -Cores over Large Bipartite Graphs [technical report]. <https://github.com/XueQingSBQ/SQAC/blob/master/full.pdf>.
- [3] Adel Ahmed, Vladimir Batagelj, Xiaoyan Fu, Seok-Hee Hong, Damian Merrick, and Andrej Mrvar. 2007. Visualisation and analysis of the internet movie database. In *2007 6th International Asia-Pacific Symposium on Visualization*. IEEE, 17–24.
- [4] Esra Akbas and Peixiang Zhao. 2017. Truss-based community search: a truss-equivalence based indexing approach. *Proceedings of the VLDB Endowment* 10, 11 (2017), 1298–1309.
- [5] Mohammad Allahbakhsh, Aleksandar Ignjatovic, Boualem Benatallah, Seyed-Mehdi-Reza Beheshti, Elisa Bertino, and Norman Foo. 2013. Collusion detection in online rating systems. In *Web Technologies and Applications: 15th Asia-Pacific Web Conference, APWeb 2013, Sydney, Australia, April 4-6, 2013. Proceedings 15*. Springer, 196–207.
- [6] Wen Bai, Yadi Chen, Di Wu, Zhichuan Huang, Yipeng Zhou, and Chuan Xu. 2022. Generalized core maintenance of dynamic bipartite graphs. *Data Mining and Knowledge Discovery* (2022), 1–31.
- [7] Vladimir Batagelj and Matjaz Zaversnik. 2003. An $o(m)$ algorithm for cores decomposition of networks. CoRR. *arXiv preprint cs.DS/0310049* 37 (2003).
- [8] Alex Beutel, Wanhong Xu, Venkatesan Guruswami, Christopher Palow, and Christos Faloutsos. 2013. Copycatch: stopping group attacks by spotting lockstep behavior in social networks. In *Proceedings of the 22nd international conference on World Wide Web*. 119–130.
- [9] Monika Cerinšek and Vladimir Batagelj. 2015. Generalized two-mode cores. *Social Networks* 42 (2015), 80–87.
- [10] Lijun Chang. 2020. Efficient maximum clique computation and enumeration over large sparse graphs. *The VLDB Journal* 29, 5 (2020), 999–1022.
- [11] Lijun Chang. 2023. Efficient maximum k -defective clique computation with improved time complexity. *Proceedings of the ACM on Management of Data* 1, 3 (2023), 1–26.
- [12] Lijun Chang, Rashmika Gamage, and Jeffrey Xu Yu. 2024. Efficient k -Clique count estimation with accuracy guarantee. *Proceedings of the VLDB Endowment* 17, 11 (2024), 3707–3719.
- [13] Huipeng Chen, Alessio Conte, Roberto Grossi, Grigorios Loukides, Solon P Pissis, and Michelle Sweering. 2021. On breaking truss-based communities. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 117–126.
- [14] Jiujian Chen, Kai Wang, Rong-Hua Li, Hongchao Qin, Xuemin Lin, and Guoren Wang. 2024. Maximal Biclique Enumeration: A Prefix Tree Based Approach. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*. IEEE, 2544–2556.
- [15] Lu Chen, Chengfei Liu, Rui Zhou, Jiajie Xu, and Jianxin Li. 2021. Efficient exact algorithms for maximum balanced biclique search in bipartite graphs. In *Proceedings of the 2021 International Conference on Management of Data*. 248–260.
- [16] Lu Chen, Chengfei Liu, Rui Zhou, Jiajie Xu, and Jianxin Li. 2022. Efficient maximal biclique enumeration for large sparse bipartite graphs. *Proceedings of the VLDB Endowment* 15, 8 (2022), 1559–1571.
- [17] James Cheng, Yiping Ke, Shumo Chu, and M Tamer Özsu. 2011. Efficient core decomposition in massive networks. In *2011 IEEE 27th International Conference on Data Engineering*. IEEE, 51–62.
- [18] Hongbo Deng, Michael R Lyu, and Irwin King. 2009. A generalized co-hits algorithm and its application to bipartite graphs. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. 239–248.
- [19] Xin Deng, Peng Peng, Chuanyu Liu, Xianyan Xie, Hui Zhou, and Zheng Qin. 2025. Efficient Indexing for Label-Constrained Cohesive Subgraph Queries over Large Graphs. In *2025 IEEE 41st International Conference on Data Engineering (ICDE)*. IEEE Computer Society, 2135–2147.
- [20] Xin Deng, Zheng Qin, Peng Peng, and Hui Zhou. 2025. TopK-BC: Efficient Maintenance of Top k (p, q) -Bicliques over Streaming Bipartite Graphs. In *2025 IEEE 41st International Conference on Data Engineering (ICDE)*. IEEE Computer Society, 3696–3709.
- [21] Danhao Ding, Hui Li, Zhipeng Huang, and Nikos Mamoulis. 2017. Efficient fault-tolerant group recommendation using alpha-beta-core. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. 2047–2050.
- [22] Yixiang Fang, CK Cheng, Siqiang Luo, and Jiafeng Hu. 2016. Effective community search for large attributed graphs. *Proceedings of the VLDB Endowment* (2016).
- [23] Yunguo Guan, Rongxing Lu, Yandong Zheng, Songnian Zhang, Jun Shao, and Guiyi Wei. 2022. Achieving Efficient and Privacy-Preserving (α, β) -Core Query Over Bipartite Graphs in Cloud. *IEEE Transactions on Dependable and Secure Computing* 20, 3 (2022), 1979–1993.
- [24] Stephan Gunnemann, Emmanuel Muller, Sebastian Raubach, and Thomas Seidl. 2011. Flexible fault tolerant subspace clustering for data with missing values. In *2011 IEEE 11th International Conference on Data Mining*. IEEE, 231–240.
- [25] Xin Huang, Hong Cheng, Lu Qin, Wentao Tian, and Jeffrey Xu Yu. 2014. Querying k -truss community in large and dynamic graphs. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. 1311–1322.
- [26] Yihao Huang, Claire Wang, Jessica Shi, and Julian Shun. 2023. Efficient algorithms for parallel bi-core decomposition. In *2023 Symposium on Algorithmic Principles of Computer Systems (APOCS)*. SIAM, 17–32.
- [27] Qinma Kang, Yunfan Kang, and Hanzhang Kong. 2018. Edge-attributed community search for large graphs. In *Proceedings of the 2nd International Conference on Big Data Research*. 114–118.
- [28] Dengshi Li, Xiaocong Liang, Ruimin Hu, Lu Zeng, and Xiaochen Wang. 2022. (α, β) -AWCS: (α, β) -Attributed Weighted Community Search on Bipartite Graphs. In *IJCNN*. IEEE, 1–8.
- [29] Boge Liu, Long Yuan, Xuemin Lin, Lu Qin, Wenjie Zhang, and Jingren Zhou. 2019. Efficient (α, β) -core computation: An index-based approach. In *The World Wide Web Conference*. 1130–1141.
- [30] Qing Liu, Xuankun Liao, Xin Huang, Jianliang Xu, and Yunjun Gao. 2023. Distributed (α, β) -core decomposition over bipartite graphs. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 909–921.
- [31] Wensheng Luo, Qiaoyuan Yang, Yixiang Fang, and Xu Zhou. 2023. Efficient Core Maintenance in Large Bipartite Graphs. *Proceedings of the ACM on Management of Data* 1, 3 (2023), 1–26.
- [32] Bingqing Lyu, Lu Qin, Xuemin Lin, Ying Zhang, Zhengping Qian, and Jingren Zhou. 2020. Maximum biclique search at billion scale. *Proceedings of the VLDB Endowment* (2020).
- [33] René Peeters. 2003. The maximum edge biclique problem is NP-complete. *Discrete Applied Mathematics* 131, 3 (2003), 651–654.
- [34] Ardian Kristanto Poernomo and Vivekanand Gopalkrishnan. 2009. Towards efficient mining of proportional fault-tolerant frequent itemsets. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. 697–706.
- [35] Ahmet Erdem Sariyüce, Buğra Gedik, Gabriela Jacques-Silva, Kun-Lung Wu, and Umit V Çatalyürek. 2016. Incremental k -core decomposition: algorithms and evaluation. *The VLDB Journal* 25 (2016), 425–447.
- [36] Ahmet Erdem Sariyüce and Ali Pinar. 2018. Peeling bipartite networks for dense subgraph discovery. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*. 504–512.
- [37] Bintao Sun, T-H Hubert Chan, and Mauro Sozio. 2020. Fully dynamic approximate k -core decomposition in hypergraphs. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 14, 4 (2020), 1–21.
- [38] Haibo Wang, Chuan Zhou, Jia Wu, Weizhen Dang, Xingquan Zhu, and Jilong Wang. 2018. Deep structure learning for fraud detection. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 567–576.
- [39] Jia Wang and James Cheng. 2012. Truss Decomposition in Massive Networks. *Proceedings of the VLDB Endowment* 5, 9 (2012).
- [40] Jun Wang, Arjen P De Vries, and Marcel JT Reinders. 2006. Unifying user-based and item-based collaborative filtering approaches by similarity fusion. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*. 501–508.
- [41] Kai Wang, Xuemin Lin, Lu Qin, Wenjie Zhang, and Ying Zhang. 2020. Efficient bitruss decomposition for large-scale bipartite graphs. In *2020 IEEE 36th International Conference on Data Engineering (ICDE)*. IEEE, 661–672.
- [42] Kaixin Wang, Kaiqiang Yu, and Cheng Long. 2024. Efficient k -Clique listing: an edge-oriented branching strategy. *Proceedings of the ACM on Management of Data* 2, 1 (2024), 1–26.
- [43] Kai Wang, Wenjie Zhang, Xuemin Lin, Ying Zhang, Lu Qin, and Yuting Zhang. 2021. Efficient and effective community search on large-scale bipartite graphs. In *2021 IEEE 37th International Conference on Data Engineering (ICDE)*. IEEE, 85–96.
- [44] Yiyuan Wang, Shaowei Cai, and Minghao Yin. 2018. New heuristic approaches for maximum balanced biclique problem. *Information Sciences* 432 (2018), 362–375.
- [45] Yue Wang, Ruiqi Xu, Xun Jian, Alexander Zhou, and Lei Chen. 2022. Towards distributed bitruss decomposition on bipartite graphs. *Proceedings of the VLDB Endowment* 15, 9 (2022), 1889–1901.
- [46] Zongyu Xu, Yihao Zhang, Long Yuan, Yuwen Qian, Zi Chen, Mingliang Zhou, Qin Mao, and Weibin Pan. 2023. Effective community search on large attributed bipartite graphs. *International Journal of Pattern Recognition and Artificial Intelligence* 37, 02 (2023), 2359002.
- [47] Yang Zhang and Srinivasan Parthasarathy. 2012. Extracting analyzing and visualizing triangle k -core motifs within networks. In *2012 IEEE 28th international conference on data engineering*. IEEE, 1049–1060.
- [48] Yuting Zhang, Kai Wang, Wenjie Zhang, Xuemin Lin, and Ying Zhang. 2021. Pareto-optimal community search on large bipartite graphs. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 2647–2656.
- [49] Zhaonian Zou. 2016. Bitruss decomposition of bipartite graphs. In *International conference on database systems for advanced applications*. Springer, 218–233.