# Exploring Cohesive Subgraphs with Vertex Engagement and Tie Strength in Bipartite Graphs

[1]Yizhang He, [1]Kai Wang*, [1]Wenjie Zhang, [1]Xuemin Lin, [2]Ying Zhang

[1]*School of Computer Science and Engineering, University of New South Wales, NSW 2033, Australia*
[2]*Centre for AI, University of Technology Sydney, NSW 2007, Australia*

**Abstract**

We propose a novel cohesive subgraph model called $\tau$-strengthened $(\alpha, \beta)$-core (denoted as $(\alpha, \beta)_\tau$-core), which is the first to consider both tie strength and vertex engagement on bipartite graphs. An edge is a strong tie if contained in at least $\tau$ butterflies ($2 \times 2$-bicliques). $(\alpha, \beta)_\tau$-core requires each vertex on the upper or lower level to have at least $\alpha$ or $\beta$ strong ties, given strength level $\tau$. To retrieve the vertices of $(\alpha, \beta)_\tau$-core optimally, we construct index $I_{\alpha,\beta,\tau}$ to store all $(\alpha, \beta)_\tau$-cores. Effective optimization techniques are proposed to improve index construction. To make our idea practical on large graphs, we propose 2D-indexes $I_{\alpha,\beta}, I_{\beta,\tau}$, and $I_{\alpha,\tau}$ that selectively store the vertices of $(\alpha, \beta)_\tau$-core for some $\alpha, \beta$, and $\tau$. The 2D-indexes are more space-efficient and require less construction time, each of which can support $(\alpha, \beta)_\tau$-core queries. As query efficiency depends on input parameters and the choice of 2D-index, we propose a learning-based hybrid computation paradigm by training a feed-forward neural network to predict the optimal choice of 2D-index that minimizes the query time. Extensive experiments show that (1) $(\alpha, \beta)_\tau$-core is an effective model capturing unique and important cohesive subgraphs; (2) the proposed techniques significantly improve the efficiency of index construction and query processing.

*Keywords:* Bipartite graph; Cohesive subgraph; Classification; Vertex engagement; Tie strength

## 1. Introduction

Bipartite graphs are widely used to represent networks with two different groups of entities such as user-item networks [1], author-paper networks [2], and member-activity networks [3]. In bipartite graphs, cohesive subgraph mining has numerous applications including fraudsters detection [4, 5, 6], group recommendation [7, 8] and discovering inter-corporate relations [9, 10].

$(\alpha, \beta)$-core and bitruss are two representative cohesive subgraph models in bipartite graphs extended from the unipartite $k$-core [11] and $k$-truss [12] models. $(\alpha, \beta)$-core is the maximal subgraph of a bipartite graph $G$ such that the vertices on upper or lower layer have at least $\alpha$ or $\beta$ neighbors respectively. $(\alpha, \beta)$-core models vertex engagement as degrees and treats each edge equally, but ties (edges) in real networks have different strengths. $k$-bitruss is the maximal subgraph where each edge is contained in at least $k$ butterflies (i.e. 2x2-biclique), which can model the tie strength [13, 14].

In the author-paper network as shown in Figure 1, the graph is the $(\alpha, \beta)$-core ($\alpha=2$, $\beta=2$) and the light blue region is the $k$-bitruss ($k=2$). Without considering tie strength, $(\alpha, \beta)$-core blindly includes research groups of different levels of cohesiveness. We can see that $v_0$ and $v_1$ are not as closely connected as the rest
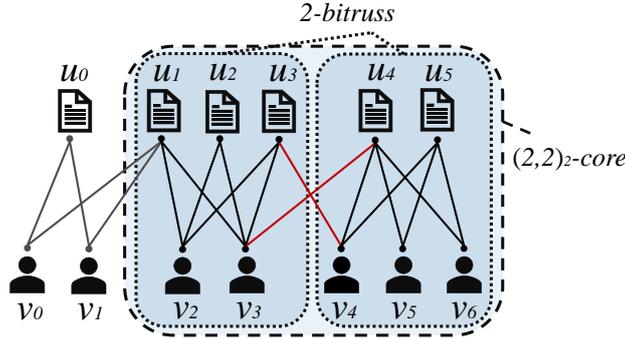
---

Figure 1: Motivation example

authors. The $k$-bitruss model can exclude the relatively sparse subgraph containing $v_0$ and $v_1$, but it also deletes edges $(u_3, v_4)$ and $(u_4, v_3)$ when their incident vertices are present. This exposes the drawbacks of the $k$-bitruss model: (1) As $k$-bitruss only keeps strong ties, the weak ties between important vertices are missed. In Fig 1, it fails to recognize the contributions of authors $v_3, v_4$ in papers $u_3, u_4$. (2) After removing weak ties, the tie strengths are modeled inaccurately. Edges $(u_3, v_3)$ and $(u_4, v_4)$ have more supporting butterflies ($u_3, u_4, v_3, v_4$ form a butterfly) than $(u_1, v_2)$, but their tie strengths are modeled as equal.

In this paper, we study the efficient and scalable computation of $\tau$-strengthened $(\alpha, \beta)$-core, which is the first cohesive subgraph model on bipartite graphs to consider both tie strength and vertex engagement. Given a bipartite graph $G$, we model the tie strength of each edge as the number of butterflies containing it. With a strength level $\tau$, we consider the edges with tie strength no less than $\tau$ to be *strong ties*. The *engagement* of a vertex is modeled as the number of strong ties it is incident to. Given engagement constraints $\alpha, \beta$ and strength level $\tau$, $(\alpha, \beta)_\tau$-core is the maximal subgraph of $G$ such that each upper or lower vertex in the subgraph has at least $\alpha$ or $\beta$ strong ties. The $(\alpha, \beta)_\tau$-core model is highly flexible and is able to capture unique structures. For instance, in Figure 1, the subgraph induced by vertices $\{u_1, u_2, u_3, u_4, u_5, v_2, v_3, v_4, v_5, v_6\}$ is the $(2, 2)_2$-core which cannot be found by $(\alpha, \beta)$-core or $k$-bitruss for any $\alpha, \beta$ or $k$. Also, as shown in Figure 1, $(\alpha, \beta)$-core can preserve the weak ties if the incident vertices are present (e.g., the red edges are preserved due to $u_3, u_4, v_3$ and $v_4$), which better resembles reality. The flexibility of the $(\alpha, \beta)_\tau$-core model is also evaluated in another experiment conducted on dataset `DBpedia-producer`. Figure 2 shows the subgraphs of different densities found by $(\alpha, \beta)_\tau$-core and $(\alpha, \beta)$-core, where density is the ratio between the number of existing edges and the number of all possible edges [13]. 165 subgraphs with a density greater than 0.2 are found by $(\alpha, \beta)_\tau$-core while only 9 such subgraphs are found by $(\alpha, \beta)$-core.
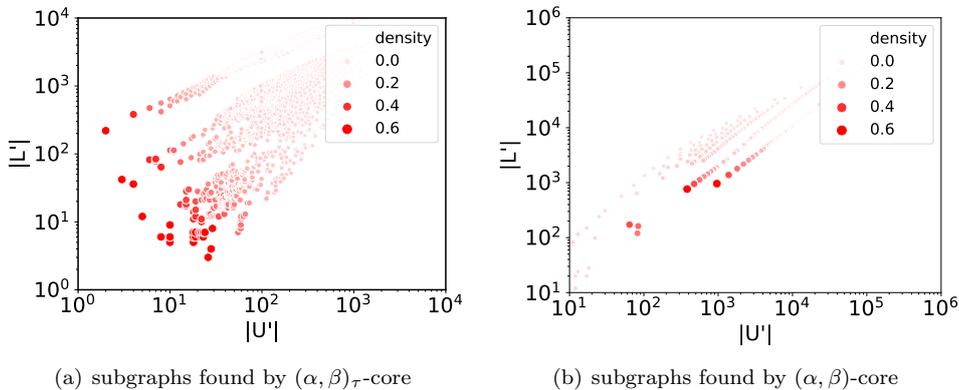


(a) subgraphs found by $(\alpha, \beta)_\tau$-core

(b) subgraphs found by $(\alpha, \beta)$-core

Figure 2: Dense subgraphs in `DBpedia-producer`.

2

**Applications.** The $\tau$-strengthened $(\alpha, \beta)$-core model has many applications. We list some of them below.

● *Identify nested communities.* On Internet forums like Reddit, Quora, and StackOverflow, users hold conversations on topics that interest them. The users and the topics form a bipartite network. In these networks, communities naturally exist and are nested. For instance, Reddit displays a list of top communities like "News", "Gaming" and "Sports" on the front page. "Sports" community contains many sub-communities including "Cricket", "Bicycling" and "Golf". The edges in sub-communities have higher tie strength because users and topics within them are more closely connected. By increasing strength level $\tau$, $(\alpha, \beta)_\tau$-core captures the subgraphs forming a hierarchy, which can model nested communities on bipartite networks.

● *Group similar users and items.* In online shopping platforms like Amazon, eBay and Alibaba, users and items form a bipartite graph, where each edge indicates a purchasing record. Such a network consists of many closely connected communities, where some items are repeatedly bought by the same group of users (i.e., the target market). Examples of such communities include gym-goers and gym attires, students and stationery, diabetic patients and no-sugar foods, etc. Within one community, items are considered more similar and users tend to be alike due to their common shopping habits. As the edges between these users and items have high tie strength (butterfly support), we can use $(\alpha, \beta)_\tau$-core to find these communities and group similar users or items together.

**Challenges.** To obtain the $(\alpha, \beta)_\tau$-core from the input graph, we can first compute the support of edges and the engagement of vertices and then iteratively delete the vertices not meeting the engagement constraints. When $\alpha$, $\beta$, $\tau$ are large, $(\alpha, \beta)_\tau$-core is small and computing $(\alpha, \beta)_\tau$-core from the input graph is time-consuming. Thus, the online computation method cannot support a large number of $(\alpha, \beta)_\tau$-core queries.

In this paper, we resort to index-based approaches. A straightforward solution is to compute all possible $(\alpha, \beta)_\tau$-cores and build a total index $I_{\alpha,\beta,\tau}$ based on them. Instead of computing all $(\alpha, \beta)_\tau$-cores from the input graph, we take advantage of the nested property of the $(\alpha, \beta)_\tau$-core, which means that if $\alpha \geq \alpha^*$, $\beta \geq \beta^*$ and $\tau \geq \tau^*$, $(\alpha, \beta)_\tau$-core is a subgraph of $(\alpha^*, \beta^*)_{\tau^*}$-core. Specifically, for all possible $\alpha$ and $\beta$, we first find $(\alpha, \beta)_1$-core and then compute $(\alpha, \beta)_\tau$-core while gradually increasing strength level $\tau$. In this manner, we can compute all $(\alpha, \beta)_\tau$-cores and construct the index $I_{\alpha,\beta,\tau}$. Although $I_{\alpha,\beta,\tau}$ supports optimal retrieval of the vertex set of any $(\alpha, \beta)_\tau$-core, it still suffers from long construction time on large graphs. To devise more practical index-based approaches, we face the following challenges.

1. When building index $I_{\alpha,\beta,\tau}$, it is time-consuming to enumerate all butterflies containing the deleted edges. Also, the $I_{\alpha,\beta,\tau}$ index construction algorithm is prone to visit the same $(\alpha, \beta)_\tau$-core subgraph repeatedly as it can correspond to different combinations of $\alpha, \beta$, and $\tau$. It is a challenge to speed up butterfly enumeration and avoid repeatedly visiting the same subgraphs during the construction of the total index $I_{\alpha,\beta,\tau}$.

2. Due to the flexibility of the $(\alpha, \beta)_\tau$-core model, there are a large number of $(\alpha, \beta)_\tau$-cores corresponding to different combinations of $\alpha$, $\beta$, and $\tau$. The time cost of indexing all $(\alpha, \beta)_\tau$-cores becomes not affordable on large graphs. It is also a challenge to strike a balance between building space-efficient indexes and supporting efficient and scalable query processing.

**Our approaches.** To address the first challenge, we extend the butterfly enumeration techniques in [15] and propose novel computation sharing optimizations to speed up the index construction process of $I_{\alpha,\beta,\tau}$. Specifically, we build a `Bloom-Edge-Index` (hereafter denoted by `BE-Index`) proposed in [15] to quickly fetch the butterflies containing an edge. The `BE-Index` captures the relationships between edges and $(2 \times k)$-bicliques (also called *blooms*). When an edge is deleted, we can quickly locate the blooms containing this edge in the `BE-Index` and update the support of affected edges in these blooms accordingly. In addition, computation-sharing optimization is based on the fact that the same $(\alpha, \beta)_\tau$-core subgraph corresponds to various parameter combinations. If we realize the vertices in a subgraph have already been recorded, we can choose to skip the current parameter combination.

To address the second challenge, we introduce space-efficient 2D-indexes including $I_{\alpha,\beta}$, $I_{\beta,\tau}$, and $I_{\alpha,\tau}$, and train a feed-forward neural network to predict the most promising index to handle an $(\alpha, \beta)_\tau$-core query. Instead of indexing all $(\alpha, \beta)_\tau$-cores, the 2D-indexes $I_{\alpha,\beta}$, $I_{\beta,\tau}$, and $I_{\alpha,\tau}$ store the vertex sets of all $(\alpha, \beta)$-core, $(1, \beta)_\tau$-core, and $(\alpha, 1)_\tau$-core respectively. These 2D-indexes are much smaller in size and require significantly less build time, each of which can be used to handle $(\alpha, \beta)_\tau$-core queries. For example,

to compute $(\alpha, \beta)_\tau$-core using $I_{\beta, \tau}$, we fetch the vertices in $(1, \beta)_\tau$-core and recover the edges of $(1, \beta)_\tau$-core. Then, we iteratively remove the vertices not having enough engagement from $(1, \beta)_\tau$-core until we find $(\alpha, \beta)_\tau$-core. However, the query processing performance based on each 2D-index is highly sensitive to parameters $\alpha$, $\beta$, and $\tau$. This is because the 2D-indexes only store the vertices in $(\alpha, \beta)$-core, $(1, \beta)_\tau$-core, and $(\alpha, 1)_\tau$-core and the size difference between $(\alpha, \beta)_\tau$-core and each of these subgraphs is uncertain. We also observe that there are no simple rules to partition the parameter space so that queries from each partition can be efficiently handled by one type of index. This motivates us to resort to machine learning techniques and train a feed-forward neural network as the classifier to predict the optimal choice of the index for each incoming query of $(\alpha, \beta)_\tau$-core. Since we aim to minimize the query time instead of accuracy, we propose a scoring function, *time-sensitive-error*, to tune the hyper-parameters of the classifier. The experiment results show that the resulting hybrid computation algorithm significantly outperforms the query processing algorithms based on $I_{\alpha, \beta}, I_{\beta, \tau}$, and $I_{\alpha, \tau}$, and it is less sensitive to varying parameters.

**Contribution.** Our major contributions are summarized here:

• We propose the first cohesive subgraph model $\tau$-strengthened $(\alpha, \beta)$-core on bipartite graphs which considers both tie strength and vertex engagement. The flexibility of our model allows it to capture unique and useful structures on bipartite graphs.

• We construct index $I_{\alpha, \beta, \tau}$ to support optimal retrieval of the vertex set of any $(\alpha, \beta)_\tau$-core. We also devise computation sharing and `BE-Index` based optimizations to effectively reduce its construction time.

• We build 2D-indexes that are more space-efficient and require significantly less build time. Also, we propose a learning-based hybrid computation paradigm to predict which index to choose to minimize the response time for an incoming $(\alpha, \beta)_\tau$-core query.

• We validate the efficiency of proposed algorithms and the effectiveness of our model through extensive experiments on real-world datasets. Results show that the 2D-indexes are scalable and the hybrid computation algorithm on a well trained neural network can outperform the algorithms based on each 2D-index alone.

**Organization.** The rest of the paper is organized as follows. Section 2 reviews the related work. Section 3 summarizes important notations and definitions and introduces $(\alpha, \beta)$-core and $\tau$-strengthened $(\alpha, \beta)$-core. Section 4 presents the online computation algorithm. Section 5 and 6 presents the total index $I_{\alpha, \beta, \tau}$ and optimizations of the index construction process. Section 7 presents the learning-based hybrid computation paradigm. Section 8 shows the experimental results and Section 9 concludes the paper.

## 2. Related work

In the literature, there are many recent studies on cohesive subgraph models on both unipartite graphs and bipartite graphs.

*Unipartite graphs.* $k$-core [11, 16, 17, 18] and $k$-truss [12, 19, 20] are two of the most well-known cohesive subgraph models on general, unipartite graphs. Given a unipartite graph, $k$-core is the maximal subgraph such that each vertex in the subgraph has at least $k$ neighbors. $k$-core models vertex engagement as degrees and assumes the importance of each tie to be equal. However, on real networks, ties (edges) have different strengths and are not of equal importance [21]. As triangles are considered as the smallest cohesive units, the number of triangles containing an edge is used to model tie strength on unipartite graphs. Thus, $k$-truss is proposed to better model tie strength, which is the maximal subgraph such that each edge in the subgraph is contained in at least $(k-2)$ triangles. The issue with $k$-truss is that it does not tolerate the existence of weak ties, which is inflexible for modeling real networks. To consider both vertex engagement and tie strength, the $(k,s)$-core model is proposed in [22]. In addition, recent works studied problems related to variants of $k$-core such as radius-bounded $k$-core on geo-social networks [23], core maintenance on dynamic graphs [24], core decomposition on uncertain graphs [25, 26], and anchored $k$-core problem [27, 28]. Variants of $k$-truss are also studied including $k$-truss communities on dynamic graphs [19], $k$-truss decomposition on uncertain graphs [29], and anchored $k$-truss problem [30]. However, these algorithms do not apply to bipartite graphs. Attempts to project the bipartite graph to general graphs will incur information loss and size inflation [13].

*Bipartite graphs.* In correspondence to $k$-core and $k$-truss, $(\alpha, \beta)$-core [7, 6] and $k$-bitruss [14, 15] are proposed on bipartite graphs. $(\alpha, \beta)$-core is the maximal subgraph such that each vertex on the upper or

lower level in the subgraph has at least $\alpha$ or $\beta$ neighbors. Just like $k$-core on unipartite graphs, $(\alpha, \beta)$-core cannot distinguish weak ties from strong ties. On bipartite graphs, tie strength is often modeled as the number of butterflies (i.e., $(2 \times 2)$-bicliques) containing an edge because butterflies are viewed as analogs of triangles [31, 32, 33, 15]. $k$-bitruss is the maximal subgraph such that each edge in the subgraph is contained in at least $k$ butterflies, which can model tie strength. $k$-bitruss suffers from the same issue as its counterpart $k$-truss does: it forcefully deletes all weak ties even if the incident vertices are strongly-engaged. Other works for bipartite graph analysis using cohesive structures such as $(p,q)$-core [34], fractional $k$-core [35] cannot be used to address these issues. In contrast to the above studies, we propose the first cohesive subgraph model $\tau$-strengthened $(\alpha, \beta)$-core that considers both vertex engagement and tie strength on bipartite graphs.

## 3. Problem Definition

Table 1: Summary of Notations

| Notation | Definition |
|---|---|
| $G$ | a bipartite graph |
| $\alpha, \beta$ | the engagement constraints |
| $\tau$ | the strength level |
| $nb(u, G)$ | the set of adjacent vertices of $u$ in $G$ |
| $deg(u, G)$ | the number of adjacent vertices of $u$ in $G$ |
| $\bowtie_G$ | the number of butterflies in $G$ |
| $sup(e)$ | the number of butterflies containing $e$ |
| $eng(u)$ | the number of strong ties adjacent to $u$ |
| $(\alpha, \beta)_\tau$-core | the $\tau$-strengthened $(\alpha, \beta)$-core |
| $I_{\alpha, \beta, \tau}$ | the decomposition-based index |
| $I_{\alpha, \beta}, I_{\beta, \tau}, I_{\alpha, \tau}$ | the 2D-indexes |

In this section, we formally define our cohesive subgraph model $\tau$-strengthened $(\alpha, \beta)$-core. We consider an unweighted, undirected bipartite graph $G(V, E)$. $V(G) = U(G) \cup L(G)$ denotes the set of vertices in $G$ where $U(G)$ and $L(G)$ represent the upper and lower layer, respectively. $E(G) \subseteq U(G) \times L(G)$ denotes the set of edges in $G$. We use $n = |V(G)|$ to denote the number of vertices and $m = |E(G)|$ to denote the number of edges. The maximum degree in the upper and lower layer is denoted as $d_{max}(U)$ and $d_{max}(L)$ respectively. The set of neighbors of a vertex $u$ in $G$ is denoted as $nb(u, G)$. The degree of a vertex is $deg(u, G) = |nb(u, G)|$. When the context is clear, we omit the input graph $G$ in notations.

**Definition 1.** $(\alpha, \beta)$**-core.** *Given a bipartite graph $G$ and degree constraints $\alpha$ and $\beta$, a subgraph $G'$ is the $(\alpha, \beta)$-core, denoted by $C_{\alpha, \beta}(G)$, if (1) all vertices in $G'$ satisfy degree constraints, i.e. $deg(u, G') \geq \alpha$ for each $u \in U(G')$ and $deg(v, G') \geq \beta$ for each $v \in L(G')$; and (2) $G'$ is maximal, i.e. any subgraph $G'' \supseteq G'$ is not an $(\alpha, \beta)$-core.*

**Definition 2. Butterfly.** *In a bipartite graph $G$, given vertices $u, w \in U(G)$ and $v, x \in L(G)$, a butterfly $\bowtie$ is the complete subgraph induced by $u, v, w, x$, which means both $u$ and $w$ are connected to $v$ and $x$ by edges. The total number of butterflies in $G$ is denoted as $\bowtie_G$.*

$(\alpha, \beta)$-core is a vertex-induced subgraph model, which assumes that the edges are of equal importance. To better model the strength of an edge $e$, we define the support $sup(e)$ to be the number of butterflies containing $e$.

**Definition 3. Strong Tie.** *Given an integer $\tau$, an edge $e \in E(G)$ is called a strong tie if $sup(e) \geq \tau$, where $\tau$ is called the strength level. Weak ties are the edges $e$ such that $sup(e) < \tau$.*

**Definition 4. Vertex Engagement.** *Given a strength level $\tau$ and $u \in V(G)$, the engagement $eng(u)$ is the number of strong ties incident to $u$. At strength level 0, $eng(u) = deg(u, G)$.*

5

If the engagement of an upper or lower vertex is at least $\alpha$ or $\beta$, we call it a **strongly-engaged** vertex. Otherwise, it is a **weakly-engaged** vertex.

**Definition 5. $\tau$-strengthened $(\alpha, \beta)$-core.** *Given a bipartite graph $G$ and engagement constraints $\alpha$ and $\beta$, and strength level $\tau$, a subgraph $G'$ is the $\tau$-strengthened $(\alpha, \beta)$-core, denoted by $(\alpha, \beta)_\tau$-core, if (1) $eng(u) \geq \alpha$ for each $u \in U(G')$ and $eng(v) \geq \beta$ for each $v \in L(G')$; and (2) $G'$ is maximal, i.e. any subgraph $G'' \supseteq G'$ is not a $\tau$-strengthened $(\alpha, \beta)$-core.*

**Problem Statement.** Given a bipartite graph $G$ and parameters $\alpha, \beta$ and $\tau$, we study the problem of scalable and efficient computation of $(\alpha, \beta)_\tau$-core in $G$.

---

**Algorithm 1:** OnlineComputation

**Input:** $G, \alpha, \beta, \tau$
**Output:** $(\alpha, \beta)_\tau$-core
1 Compute $sup(e)$ foreach $e \in E(G)$
2 Compute $eng(u)$ foreach $u \in V(G)$
3 *Peeling($G, \alpha, \beta, \tau, sup, eng$)*
4 **return** $G$

---

**Algorithm 2:** Peeling

**Input:** $G, \alpha, \beta, \tau, sup, eng$
**Output:** $(\alpha, \beta)_\tau$-core
1 **while** *exists $u \in V(G)$ without enough engagement* **do**
2     **foreach** $v \in nb(u)$ **do**
3         **if** $sup((u,v)) \geq \tau$ **then**
4             $eng(v) \leftarrow eng(v) - 1$
5         **foreach** $\boxtimes$ *containing* $(u,v)$ **do**
6             **foreach** *edge $e' = (u', v') \in \boxtimes$ s.t. $e' \neq e$ and $sup(e') \geq \tau$* **do**
7                 $sup(e') \leftarrow sup(e') - 1$
8                 **if** $sup(e') = \tau - 1$ **then**
9                     decrease $eng(u')$ and $eng(v')$ by 1
10         remove $(u,v)$ from $G$
11     remove $u$ from $G$
12 **return** $G$

---

## 4. The Online Computation Algorithm

Given engagement constraints $\alpha$, $\beta$ and strength level $\tau$, the online algorithm to compute the $(\alpha, \beta)_\tau$-core is outlined in Algorithm 1. First, we compute the support of each edge $e$ using the algorithm in [33] and count how many strong ties each vertex $u$ has. Then, Algorithm 2 is invoked to iteratively remove the vertices without enough engagement along with their incident edges. The vertices in $U(G)$ and $L(G)$ are sorted by engagement and the edges are sorted by support. In this manner, we can always delete the vertices with the smallest engagement first and quickly identify which edges are strong ties. When an edge $e$ is removed due to lack of support (i.e., $sup(e) < \tau$), we go through all the butterflies containing $e$ and update the supports of the edges in these butterflies (lines 5-9). Specifically, we do not need to update the support of the edges connected to weakly-engaged vertices, because they will be removed (line 10). Neither do we update the support of weak ties because they do not contribute to any vertex engagement. In other words, we only update the support of strong ties between strongly-engaged vertices. When a strong tie

6

becomes a weak tie, we decrease the engagement of its incident vertices (lines 4,9). The order of vertices and edges are maintained in linear heaps [36] after their engagement and support are updated. Here we evaluate the time and space complexity of Algorithm 1.

**Lemma 1.** *The time complexity of Algorithm 1 is* $O(\sum_{(u,v)\in E(G)} \sum_{w\in nb(v)} min(deg(u), deg(w)))$ *and the space complexity is* $O(m)$.

*Proof.* The butterfly counting process takes $O(\sum_{(u,v)\in E(G)} min(deg(u), deg(v)))$ time [33]. After the support of each edge is computed, it takes $O(m)$ time to compute the engagement for each vertex. Then, we need to run Algorithm 2. For each weakly engaged vertex $u$, we need to delete all its incident edges, which dominates the time cost of Algorithm 2.

For each deleted edge $(u, v)$, we need to enumerate the butterflies containing it. Let $w$ be a vertex in $nb(v, G) \setminus \{u\}$. For each vertex $x$ in $nb(u, G) \cap nb(w, G)$, the induced subgraph of $\{u, v, w, x\}$ is a butterfly. The set intersection (computing $nb(u, G) \cap nb(w, G)$) can be implemented in $O(min(deg(u), deg(w))$ time by using a $O(m)$ hash table to store the neighbor set of each vertex.

Thus, the butterfly enumeration for each delete edge $(u, v)$ takes $O(\sum_{w\in nb(v)} min(deg(u), deg(w)))$ time. As each edge can only be deleted once, the total time complexity of the butterfly enumeration for all deleted edges takes $O(\sum_{(u,v)\in E(G)} \sum_{w\in nb(v)} min(deg(u), deg(w)))$ time. Thus, the time complexity of Algorithm 1 is $O(\sum_{(u,v)\in E(G)} \sum_{w\in nb(v)} min(deg(u), deg(w)))$ which is denoted as $T_{peel}(G)$ hereafter.

We store the neighbors of each vertex as adjacency lists as well as the support of edges and engagement of vertices, which in total takes $O(m)$ space. Therefore, the space complexity of Algorithm 1 is $O(m)$. □
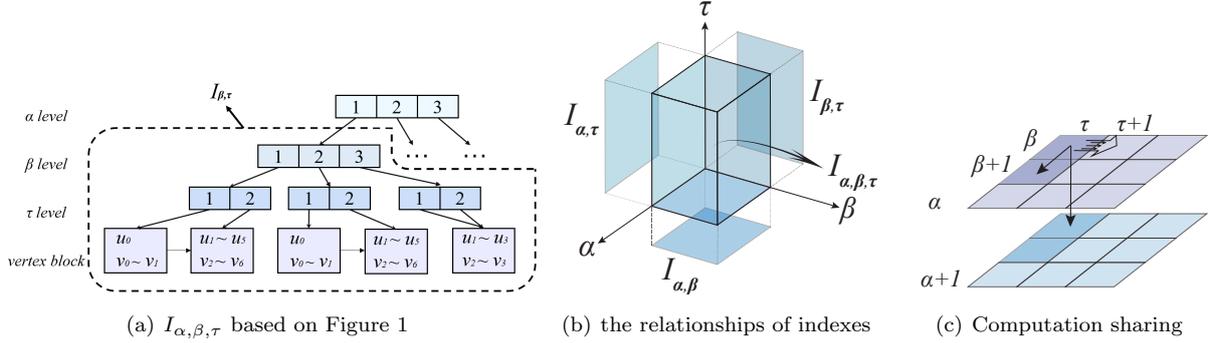


(a) $I_{\alpha,\beta,\tau}$ based on Figure 1     (b) the relationships of indexes     (c) Computation sharing

Figure 3: Illustrating our ideas

## 5. The Decomposition Based Total Index

Given $\alpha$, $\beta$, and $\tau$, Algorithm 1 computes the $(\alpha, \beta)_\tau$-core from the input graph, which is slow and cannot handle a large number of queries. In this section, we present a decomposition algorithm that retrieves all $(\alpha, \beta)_\tau$-cores and build a total index based on the decomposition output to support efficient query processing.

**The decomposition algorithm.** The following lemma is immediate based on Definition 5, which depicts the nested relationships among $(\alpha, \beta)_\tau$-cores.

**Lemma 2.** $(\alpha, \beta)_\tau$-core $\subseteq (\alpha', \beta')_{\tau'}$-core if $\alpha \geq \alpha'$, $\beta \geq \beta'$, and $\tau \geq \tau'$.

Based on Lemma 2, if a vertex $u$ is in $(\alpha, \beta)_{\tau'}$-core, $u$ is also in $(\alpha, \beta)_\tau$-core if $\tau < \tau'$. Thus, in the decomposition, for given vertex $u$ and $\alpha$, $\beta$ values, we aim to retrieve the maximum $\tau$ value such that $u$ is in the corresponding $(\alpha, \beta)_\tau$-core, namely, $\tau_{max}(\alpha, \beta, u) = \max\{\tau | u \in (\alpha, \beta)_\tau\text{-}core\}$. For each vertex $u$ and all possible combinations of $\alpha$ and $\beta$, it is only necessary to store $u$ in $(\alpha, \beta)_{\tau'}$-core to build a

---
**Algorithm 3:** Decomposition

**Input:** $G(V = (U, L), E)$
**Output:** $\tau_{max}(\alpha, \beta, u)$, for all $\alpha, \beta, \forall u \in V(G)$

**1** $\alpha \leftarrow 1$, $\beta \leftarrow 1$, $\tau \leftarrow 1$
**2** Compute $sup(e)$, $\forall e \in E(G)$
**3** Compute $eng(u)$, $\forall u \in V(G)$
**4** **while** $(\alpha, 1)_1$-core in $G$ is not empty **do**
**5**   $Peeling((\alpha, 1)_1\text{-core}, \alpha, 1, 1, sup, eng)$; $\beta \leftarrow 1$
**6**   **while** $(\alpha, \beta)_1$-core in $G$ is not empty **do**
**7**    $sup' \leftarrow sup$; $eng' \leftarrow eng$; $\tau \leftarrow 1$
**8**    $Peeling((\alpha, \beta)_1\text{-core}, \alpha, \beta, 1, sup', eng')$
**9**    **while** $(\alpha, \beta)_\tau$-core in $G$ is not empty **do**
**10**     $sup'' \leftarrow sup'$; $eng'' \leftarrow eng'$
**11**     $Peeling((\alpha, \beta)_\tau\text{-core}, \alpha, \beta, \tau, sup'', eng'')$,   add $\tau_{max}(\alpha, \beta, u) \leftarrow \tau - 1$ before line 2
**12**     $\tau \leftarrow \tau + 1$
**13**     **foreach** $e = (u', v') \in E(G)$, $sup(e) = \tau$ **do**
**14**      $eng''(u') \leftarrow eng''(u') - 1$
**15**      $eng''(v') \leftarrow eng''(v') - 1$
**16**    $\beta \leftarrow \beta + 1$
**17**   $\alpha \leftarrow \alpha + 1$
**18** **return** $\tau_{max}(\alpha, \beta, u)$, for all $\alpha, \beta, \forall u \in V(G)$

---

space compact index, where $\tau' = \tau_{max}(\alpha, \beta, u)$ and $u \in (\alpha, \beta)_\tau$-core can be implied if $\tau < \tau'$. Algorithm 3 is devised for $(\alpha, \beta)_\tau$-core decomposition which applies three nested loops to go through all possible $\alpha, \beta, \tau$ combinations. Note that when computing the $(\alpha, \beta)_\tau$-core, we record $\tau_{max}(\alpha, \beta, u)$ for each vertex $u$. Specifically, for a vertex $u$ contained in $(\alpha, \beta)_{\tau_0}$-core but is removed when computing $(\alpha, \beta)_{\tau_0+1}$-core, we assign $\tau_0$ to $\tau_{max}(\alpha, \beta, u)$ (line 11). Here we evaluate the time and space complexity of Algorithm 3.

**Lemma 3.** *The time complexity of Algorithm 3 is $O(\sum_{\alpha=1}^{\alpha_{max}} \sum_{\beta=1}^{\beta_{max}(\alpha)} T_{peel}((\alpha, \beta)_1\text{-core}))$, where $\alpha_{max}$ is the maximal $\alpha$ such that $(\alpha, 1)_1$-core exists and $\beta_{max}(\alpha)$ is the maximal $\beta$ such that $(\alpha, \beta)_1$-core exists. Algorithm 3 takes up $O(m)$ space.*

*Proof.* In the outer while-loop, we start from the input graph $G$ and compute $(\alpha, 1)_1$-core, which takes $T_{peel}(G)$. In the middle while-loop, we calculate $(\alpha, \beta)_1$-core from $(\alpha, 1)_1$-core for all possible $\alpha$, which takes $\sum_{\alpha=1}^{\alpha_{max}} T_{peel}((\alpha, 1)_1\text{-core})$. The dominant cost occurs in the innermost while-loop when we run the Peeling algorithm on $(\alpha, \beta)_1$-core for all possible $\alpha$ and $\beta$. As iterative removing edges and vertices in $(\alpha, \beta)_1$-core until it is empty takes $T_{peel}((\alpha, \beta)_1\text{-core})$, the overall time complexity is $O(\sum_{\alpha=1}^{\alpha_{max}} \sum_{\beta=1}^{\beta_{max}(\alpha)} T_{peel}((\alpha, \beta)_1\text{-core}))$.

At any time during the execution of this algorithm, we always store $G$, $(\alpha, 1)_1$-core, $(\alpha, \beta)_1$-core, and $(\alpha, \beta)_\tau$-core in memory, which takes $O(m)$ space. We also store the support of edges and engagement of vertices in these graphs, which also takes $O(m)$ space. Thus, the total space complexity is $O(m)$. □

**Decomposition-based index.** Based on the decomposition results, a four level index $I_{\alpha, \beta, \tau}$ can be constructed to support query processing as shown in Figure 3(a).

- $\alpha$ *level.* The $\alpha$ level of $I_{\alpha, \beta, \tau}$ is an array of pointers, each of which points to an array in the $\beta$ level. The length of the array in the $\alpha$ level is $\alpha_{max}$. The $k_{th}$ element is denoted as $I_{\alpha, \beta, \tau}[k]$.
- $\beta$ *level.* The $\beta$ level has $\alpha_{max}$ arrays of pointers. The array pointed by $I_{\alpha, \beta, \tau}[k]$ has length $\beta_{max}(\alpha)$. The $j_{th}$ pointer in the $k_{th}$ array is denoted as $I_{\alpha, \beta, \tau}[k][j]$, which points to an array in the $\tau$ level.
- $\tau$ *level.* The $\tau$ level has $\sum_{i=1}^{\alpha_{max}} \beta_{max}(i)$ arrays of pointers to vertex blocks, corresponding to all pairs of $\alpha, \beta$. The array pointed by $I_{\alpha, \beta, \tau}[k][j]$ has length $\tau_{max}(\alpha, \beta) = \max\{\tau | (\alpha, \beta)_\tau\text{-core exists}\}$.

---

**Algorithm 4:** DecompQuery

> **Input:** $I_{\alpha,\beta,\tau}, \alpha, \beta, \tau, G$
> **Output:** $(\alpha,\beta)_\tau$-core

**1** $U', V', E' \leftarrow \emptyset$
**2** **if** $I_{\alpha,\beta,\tau}.size < \alpha$ *or* $I_{\alpha,\beta,\tau}[\alpha].size < \beta$ *or* $I_{\alpha,\beta,\tau}[\alpha][\beta].size < \tau$ **then**
**3** $\quad$ return $\emptyset$
**4** $ptr \leftarrow I_{\alpha,\beta,\tau}[\alpha][\beta][\tau]$
**5** **while** *ptr is not null* **do**
**6** $\quad$ $v \leftarrow$ vertices in vertex block pointed by $ptr$
**7** $\quad$ **if** $v \in U(G)$ **then**
**8** $\quad\quad$ $U' \leftarrow U' \cup v$
**9** $\quad$ **else**
**10** $\quad\quad$ $V' \leftarrow V' \cup v$
**11** $\quad$ $ptr \leftarrow$ the address of the next vertex block
**12** $E' \leftarrow E(G) \cap (U' \times V')$
**13** **return** $G' = (U', V', E')$

---

- *vertex blocks.* The fourth level of $I_{\alpha,\beta,\tau}$ is a singly linked list of vertex blocks. Each vertex block corresponds to a set of $\alpha, \beta, \tau$ value, which contains all vertex $u$ such that $\tau_{max}(\alpha, \beta, u) = \tau$ along with a pointer to the next vertex block. The vertex blocks with the same $\alpha$ and $\beta$ are sorted by the associated $\tau$ values and each of them has a pointer to the next. Among them, the vertex block with the largest $\tau$ has its pointer pointing to null.

We can construct index $I_{\alpha,\beta,\tau}$ from the output of Algorithm 4. Given $\alpha$ and $\beta$, we store all vertices $u$ in the same vertex block if they have the same $\tau_{max}(\alpha, \beta, u)$, so each vertex block has an associated $(\alpha, \beta, \tau)$ value. In each $I_{\alpha,\beta,\tau}[\alpha][\beta][\tau]$, we store the address of the vertex block associated with $(\alpha, \beta, \tau')$, where $\tau'$ is the smallest integer such that $\tau \leq \tau'$. The index construction time is linear to the size of decomposition results, which is bounded by the time complexity of Algorithm 4.

**Lemma 4.** *The space complexity of index $I_{\alpha,\beta,\tau}$ is $O(\sum_{i=1}^{\alpha_{max}} \sum_{j=1}^{\beta_{max}(i)} (\tau_{max}(i,j) + n))$, where $\tau_{max}(i,j)$ be the maximal $\tau$ in all $(i,j)_\tau$-cores.*

*Proof.* By construction, the space complexity of the first two levels of pointers is bounded by that of $\tau$ level. The space complexity of $\tau$ level is $\sum_{i=1}^{\alpha_{max}} \sum_{j=1}^{\beta_{max}(i)} \tau_{max}(i,j)$. Given vertex $u$, let $\alpha_{max}(u)$ be the maximal $\alpha$ such that $u \in (\alpha, \beta)_\tau$-*core*. The space complexity of the vertex blocks is $\sum_{u \in V(G)} \sum_{i=1}^{\alpha_{max}(u)} \max\{\beta | u \in (i, \beta)_1\text{-}core\} \leq \sum_{i=1}^{\alpha_{max}} \sum_{j=1}^{\beta_{max}(i)} n$. Adding it to the space complexity of level $\tau$ completes the proof. $\square$

**Index based query processing.** When the index $I_{\alpha,\beta,\tau}$ is built on a bipartite graph $G$, Algorithm 4 outlines how to restore $(\alpha, \beta)_\tau$-core given $\alpha, \beta$, $\tau$ and $I_{\alpha,\beta,\tau}$. First, it checks the validity of the input parameters. If the queried $(\alpha, \beta)_\tau$-core does not exist, it terminates immediately (lines 2,3). Otherwise, it collects the vertices of $(\alpha, \beta)_\tau$-core from $I_{\alpha,\beta,\tau}$ and restores the edges in the queried subgraph.

**Lemma 5.** *Given a graph $G$ and parameters $\alpha, \beta$ and $\tau$, Algorithm 4 retrieves $V((\alpha, \beta)_\tau$-core$)$ from index $I_{\alpha,\beta,\tau}$ in $O(|V((\alpha, \beta)_\tau$-core$)|)$ time. The edges in $(\alpha, \beta)_\tau$-core can be retrieved in $O(\sum_{v \in V((\alpha,\beta)_\tau\text{-}core)} deg(v))$ time after obtaining the vertex set.*

*Proof.* As each vertex block only stores the vertices with one given $\tau$ value, the vertex blocks pointed by $I_{\alpha,\beta,\tau}[\alpha][\beta][\tau']$ where $\tau' \geq \tau$ give us all the vertices in $(\alpha, \beta)_\tau$-core, which takes $O(|V((\alpha, \beta)_\tau$-core$)|)$ time. To restore the edges in $(\alpha, \beta)_\tau$-core, for each vertex $v$ in $(\alpha, \beta)_\tau$-core, we go through each of its neighbors in $G$ and check if it is in $V((\alpha, \beta)_\tau$-core$)$. This takes $O(\sum_{v \in V((\alpha,\beta)_\tau\text{-}core)} (deg(v)))$ time. $\square$

According to Lemma 5, given $\alpha, \beta$ and $\tau$, the vertex set of the $(\alpha, \beta)_\tau$-core is retrieved in optimal time.

**Example 1.** *Figure 3(a) illustrates the $I_{\alpha,\beta,\tau}$ index for the bipartite graph in Figure 1. When querying $(1,2)_1$-core, we start with the vertex block pointed by $I_{\alpha,\beta,\tau}[1][2][1]$ ($u_0, v_0$ and $v_1$). Then we keep collecting the vertices until we have fetched the vertices pointed by $I_{\alpha,\beta,\tau}[1][2][2]$ ($u_1$ to $u_5$ and $v_2$ to $v_6$) . All the collected vertices provides the final solution to vertex set of $(1,2)_1$-core, which are $u_0$ to $u_5$ and $v_0$ to $v_6$.*

## 6. Optimizations of Index Construction

The above decomposition algorithm has these issues: (1) the same subgraph can be computed repeatedly for different $\alpha$ and $\beta$ values. For example, if $(1,1)_\tau$-core is the same subgraph as $(1,2)_\tau$-core, then we will compute it twice when $\beta$=1 and $\beta$=2. (2) when removing an edge $e$, we need to enumerate all the butterflies containing $e$. The basic implementation of butterfly enumeration is inefficient, which involves finding three connected vertices first and then check if a fourth vertex can form a butterfly with the existing ones. We devise computation-sharing optimizations to address the first issue, and adopt the `Bloom-Edge-Index` proposed in [33] to speed up the butterfly enumeration process.

**Computation sharing optimizations.** In this part, we reduce the times of visiting the same $(\alpha,\beta)_\tau$-core subgraphs by skipping some combinations of $\alpha$, $\beta$, and $\tau$, while yielding the same decomposition results.

• *Skip computation for $\tau$.* In Algorithm 3, if vertex $u$ is removed when computing $(\alpha,\beta)_{\tau+1}$-core from $(\alpha,\beta)_\tau$-core, we conclude that $\tau_{max}(\alpha,\beta,u)=\tau$. However, if both $(\alpha,\beta)_\tau$-core and $(\alpha,\beta)_{\tau+1}$-core are already visited for other $\beta$, this process is redundant and the current $\tau$ value can be skipped. Specifically, in the innermost while loop (lines 9-15), we can use an array $\beta_{min}[\tau]$ to store the minimal lower engagement of $(\alpha,\beta)_\tau$-core for each $\tau$. If $\beta_{min}[\tau] \geq \beta$, then the current $(\alpha,\beta)_\tau$-core has already been visited. We only compute $\tau_{max}(\alpha,\beta,u)$ values when one of $(\alpha,\beta)_\tau$-core and $(\alpha,\beta)_{\tau+1}$-core is not visited. Otherwise, we skip the current $\tau$ value. The following lemma explains how to correctly obtain multiple $\tau_{max}(\alpha,\beta,u)$ values when removing one vertex.

**Lemma 6.** *Given $\alpha, \beta, \tau$ and graph $G$, let $u$ be a vertex in $(\alpha,\beta)_\tau$-core but not in $(\alpha,\beta)_{\tau+1}$-core. If there exists an integer $k$ such that $(\alpha,\beta)_\tau$-core $= (\alpha,\beta+k)_\tau$-core, then $u \notin (\alpha,\beta+k)_{\tau+1}$-core.*

This lemma is immediate from Lemma 2, because if vertex $u$ is in $(\alpha,\beta+k)_{\tau+1}$-core, then it must also be contained in $(\alpha,\beta)_{\tau+1}$-core, which contradicts our assumption. Therefore, for vertices like $u$, we can conclude that $\tau_{max}(\alpha,\beta',u)=\tau$, for all $\beta \leq \beta' \leq \beta + k$. In this way, we fully preserve the decomposition outputs of Algorithm 3.

• *Skip computation for $\alpha$ and $\beta$.* To skip some $\beta$ values, we keep track of the minimal engagement of lower level vertices of the visited $(\alpha,\beta)_\tau$-cores in the middle while-loop (lines 7-16) in Algorithm 3 as $\beta^*$. At line 16, if $\beta^* > \beta$, then $\beta$ should be directly increased to $\beta^* + 1$ (the first value which is not computed yet) and values from $\beta$ to $\beta^*$ are skipped. This is because for all $\beta \leq \beta' \leq \beta^*$, the decomposition process are exactly the same. Likewise, to skip some $\alpha$ values, we record the minimal engagement of upper-level vertices of the visited $(\alpha,\beta)_\tau$-cores in the outermost while-loop (lines 5-17) of Algorithm 3 as $\alpha^*$. At line 17, if $\alpha^* > \alpha$, then $\alpha$ should be directly increased to $\alpha^* + 1$ and values from $\alpha$ to $\alpha^*$ are skipped.

**Example 2.** *As shown in Figure 3(a), when $\beta$=1, we remove $u_0, v_0$ and $v_1$ when computing $(1,1)_2$-core from $(1,1)_1$-core. The minimal engagement of lower vertices in $(1,1)_1$-core and $(1,1)_2$-core are 2, so the array $\beta_{min}$ is $[2,2]$. This means that $\tau_{max}(1,\beta',u) = 1$ and $\tau_{max}(1,\beta',u') = 2$, where $u \in \{u_0,v_0,v_1\}$ and $u' \in \{u_1,u_2,u_3,u_4,u_5,v_2,v_3,v_4,v_5,v_6\}$ and $\beta' \in \{1,2\}$. When $\beta$=2, we infer that $(1,2)_1$-core and $(1,2)_2$-core are already visited based on array $\beta_{min}$, so we can skip $\beta$=2. When $\beta$=3, the current $\beta$ value exceeds the values in $\beta_{min}$, so it cannot be skipped.*

**Bloom-Edge-Index-based optimization.** During edge deletions of Algorithm 3, we need to repeatedly retrieve the butterflies containing the deleted edges. To efficiently address this, We deploy a `Bloom-Edge-Index` (hereafter denoted as `BE-Index`) proposed in [15] to facilitate butterfly enumeration. Specifically, a bloom is a $2 \times k$-biclique, which contains $\frac{k \times (k-1)}{2}$ butterflies. Each edge in the bloom is contained in $k-1$ butterflies. The `BE-Index` compresses butterflies into blooms and keeps track of the edges they contain. The space complexity of `BE-Index` and the time complexity to build it are both $O(\sum_{e=(u,v)\in E(G)} min(deg(u), deg(v))$

[15]. Hereafter we also use $T_{BE}$ to represent $\sum_{e=(u,v)\in E(G)} min(deg(u), deg(v))$. Deleting an edge $e$ based on BE-Index takes $O(sup(e))$ time, where $sup(e)$ is the number of butterflies containing $e$. This is because when $e$ is deleted, BE-Index fetches the associated blooms and updates the support number of the affected edges in these blooms. In total, there are $O(sup(e))$ affected edges if edge $e$ is deleted. Here we evaluate the BE-Index's impact on the overall time and space complexity of Algorithm 3.

**Lemma 7.** *By adopting the BE-Index for edge deletions, the time complexity of Algorithm 3 becomes* $O(T_{BE}) + O(\sum_{\alpha=1}^{\alpha_{max}} \sum_{\beta=1}^{\beta_{max}(\alpha)} \boxtimes_{(\alpha,\beta)_1\text{-}core})$, *where* $\boxtimes_{(\alpha,\beta)_1\text{-}core}$ *is the number of butterflies in* $(\alpha,\beta)_1\text{-}core$.

*Proof.* In the innermost loop of Algorithm 4, we remove the edges from $(\alpha,\beta)_1$-core to get $(\alpha,\beta)_\tau$-core. As each edge deletion operation takes $sup(e)$ [15], it takes $\sum_{e\in E((\alpha,\beta)_1\text{-}core)} = O(\boxtimes_{(\alpha,\beta)_1\text{-}core})$ to compute $(\alpha,\beta)_\tau$-core from $(\alpha,\beta)_1$-core. As we are doing this for all possible $\alpha$ and $\beta$, the time complexity within the while-loops becomes $O(\sum_{\alpha=1}^{\alpha_{max}} \sum_{\beta=1}^{\beta_{max}(\alpha)} \boxtimes_{(\alpha,\beta)_1\text{-}core})$. Adding the BE-Index construction time to it completes the proof. $\square$

## 7. A Learning-based Hybrid Computation Paradigm

Although the index $I_{\alpha,\beta,\tau}$ supports the optimal retrieval of the vertices in the queried $(\alpha,\beta)_\tau$-core, it does not scale well to large graphs due to its long build time and large space complexity even with the related optimizations. For instance, on datasets Team, Wiki-en, Amazon, and DBLP, the index $I_{\alpha,\beta,\tau}$ cannot be built within two hours as evaluated in our experiments. In this section, we present 2D-indexes that selectively store the vertices of $(\alpha,\beta)_\tau$-core for some combinations of $\alpha, \beta,$ and $\tau$. We also train a feed-forward neural network on a small portion of the queries to predict the choice of 2D-index that minimizes the running time for each new incoming query.

| Index | Space Complexity | Build Time | Query Time |
|---|---|---|---|
| $I_{\alpha,\beta,\tau}$ | $O(\sum_{i=1}^{\alpha_{max}} \sum_{j=1}^{\beta_{max}(i)} (\tau_{max}(i,j) + n))$ | $O(T_{BE} + \sum_{\alpha=1}^{\alpha_{max}} \sum_{\beta=1}^{\beta_{max}(\alpha)} \boxtimes_{(\alpha,\beta)_1\text{-}core})$ | $O(\sum_{v\in V((\alpha,\beta)_\tau\text{-}core)} (deg(v)))$ |
| $I_{\alpha,\beta}$ | $O(m)$ | $O(\delta \cdot m)$ | $O(T_{peel}((\alpha,\beta)\text{-}core))$ |
| $I_{\beta,\tau}$ | $O(\sum_{j=1}^{\beta_{max}} (\tau_{max}(1,j) + n))$ | $O(T_{BE} + \sum_{\beta=1}^{\beta_{max}} \boxtimes_{(1,\beta)_1\text{-}core})$ | $O(T_{peel}((1,\beta)_\tau\text{-}core))$ |
| $I_{\alpha,\tau}$ | $O(\sum_{i=1}^{\alpha_{max}} (\tau_{max}(i,1) + n))$ | $O(T_{BE} + \sum_{\alpha=1}^{\alpha_{max}} \boxtimes_{(\alpha,1)_1\text{-}core})$ | $O(T_{peel}((\alpha,1)_\tau\text{-}core))$ |

Table 2: Space complexity, index construction time and query processing time of different indexes

**2D-indexes.** We introduce three 2D-indexes $I_{\alpha,\beta}$, $I_{\beta,\tau}$, and $I_{\alpha,\tau}$ in this part. Each of them is a three-level index with two levels of pointers and one level of vertex blocks. The main structures of them are presented as follows.
• $I_{\beta,\tau}$. For all $\beta$, $\tau$, $I_{\beta,\tau}[\beta][\tau]$ points to the vertices $u \in V(G)$ s.t. $\tau_{max}(1,\beta,u) = \tau$. It is the component of $I_{\alpha,\beta,\tau}$ where $\alpha=1$, which can fetch the vertices of all subgraphs of the form $(1,\beta)_\tau$-core in optimal time.
• $I_{\alpha,\tau}$. For all $\alpha$, $\tau$, $I_{\alpha,\tau}[\alpha][\tau]$ points to the vertices $u \in V(G)$ s.t. $\tau_{max}(\alpha,1,u) = \tau$. It is the component of $I_{\alpha,\beta,\tau}$ where $\beta=1$, which can fetch the vertices of all subgraphs of the form $(\alpha,1)_\tau$-core in optimal time.
• $I_{\alpha,\beta}$ consists of $I_{\alpha,\beta}U$ and $I_{\alpha,\beta}V$ to store the vertices in $U(G)$ and $L(G)$ separately. $I_{\alpha,\beta}U[\alpha][\beta]$ points to the vertices $u \in U(G)$ s.t. $\beta = \max\{\beta'|u \in (\alpha,\beta')\text{-}core\}$ and $I_{\alpha,\beta}V[\beta][\alpha]$ points to the vertices $v \in L(G)$ such that $\alpha = \max\{\alpha'|v \in (\alpha',\beta)\text{-}core\}$. $I_{\alpha,\beta}$ can fetch the vertices of any $(\alpha,\beta)$-core in optimal time.

Note that, $I_{\alpha,\beta}$ is proposed to support efficient $(\alpha,\beta)$-core computation as introduced in [6] while $I_{\beta,\tau}$ and $I_{\alpha,\tau}$ are essentially parts of $I_{\alpha,\beta,\tau}$. For each type of 2D-indexes, we analyze its construction time, space complexity, and the query time to compute $(\alpha,\beta)_\tau$-core based on it.

**Lemma 8.** *The time complexity to build $I_{\beta,\tau}$ is $O(T_{BE} + \sum_{\beta=1}^{\beta_{max}} \boxtimes_{(1,\beta)_1\text{-}core})$ and the space complexity of $I_{\beta,\tau}$ is $O(\sum_{j=1}^{\beta_{max}} (\tau_{max}(1,j) + n))$. It takes $T_{peel}((1,\beta)_\tau\text{-}core)$ time to compute $(\alpha,\beta)_\tau$-core using $I_{\beta,\tau}$.*

*Proof.* As discussed in Lemma 7, the BE-Index can significantly speed up butterfly enumeration during edge deletions, which takes $O(T_{BE})$ to construct. Then, we fix $\alpha$ to one and run lines 6-16 of Algorithm

3 to compute all $(1,\beta)_\tau$-core. For each possible $\beta$, this process takes $O(\bigtimes_{(1,\beta)_1\text{-}core})$ time, so in total it takes $O(\sum_{\beta=1}^{\beta_{max}}\bigtimes_{(1,\beta)_1\text{-}core}))$ time. Adding it to the `BE-Index` construction time $(O(T_{BE}))$ gives the time complexity of $I_{\beta,\tau}$ construction.

As $I_{\beta,\tau}$ is the part of $I_{\alpha,\beta,\tau}$ with $\alpha=1$, its space is equal to the part of $I_{\alpha,\beta,\tau}$ that is pointed by $I_{\alpha,\beta,\tau}[1]$. The size of the arrays of pointers in $I_{\beta,\tau}$ is bounded by $O(\sum_{j=1}^{\beta_{max}}(\tau_{max}(1,j))))$. The size of the vertex blocks is bounded by the number of vertices in all $(1,\beta)_\tau$-core, which is $O(\sum_{j=1}^{\beta_{max}}(|V((1,\beta)_\tau\text{-}core)|)) = O(\sum_{j=1}^{\beta_{max}} n)$. Therefore, the space complexity of $I_{\beta,\tau}$ is $O(\sum_{j=1}^{\beta_{max}}(\tau_{max}(1,j)+n))$. $\qquad\square$

In order to query $(\alpha,\beta)_\tau$-core based on $I_{\beta,\tau}$, we first find the $(1,\beta)_\tau$-core from $I_{\beta,\tau}$ and then compute $(\alpha,\beta)_\tau$-core from $(1,\beta)_\tau$-core.

**Lemma 9.** *The query time of $(\alpha,\beta)_\tau$-core based on $I_{\beta,\tau}$ is $T_{peel}((1,\beta)_\tau\text{-}core)$.*

*Proof.* Given engagement constraints $\alpha,\beta$ and strength level $\tau$, let $G'$ be the $(1,\beta)_\tau$-core on bipartite graph $G$. First, it takes $O(|V(G')|)$ time to fetch the vertices in $(1,\beta)_\tau$-core from $I_{\beta,\tau}$. Restoring the edges of $(1,\beta)_\tau$-core from $G$ takes $O(\sum_{u\in V(G')} deg(u,G'))$ time. Then, we call the Peeling algorithm on $(1,\beta)_\tau$-core to compute $(\alpha,\beta)_\tau$-core, which takes $O(T_{peel}((1,\beta)_\tau\text{-}core))$ time. $\qquad\square$

**Example 3.** *In Figure 3(a), the component of $I_{\alpha,\beta,\tau}$ wrapped in dotted line is the $I_{\beta,\tau}$ of the graph in Figure 1. If $(2,2)_2$-core is queried, we first to obtain $(1,2)_2$-core from $I_{\beta,\tau}$ and compute $(2,2)_2$-core by calling the peeling algorithm.*

Note that index $I_{\alpha,\tau}$ is symmetric to index $I_{\beta,\tau}$, with $\alpha$ and $\beta$ switched. It is immediate that it takes $O(T_{BE}+\sum_{\alpha=1}^{\alpha_{max}}\bigtimes_{(\alpha,1)_1\text{-}core})$ time to construct $I_{\alpha,\tau}$ and its space complexity is $O(\sum_{i=1}^{\alpha_{max}}(\tau_{max}(i,1)+n))$. It takes $T_{peel}((\alpha,1)_\tau\text{-}core)$ time to compute $(\alpha,\beta)_\tau$-core using $I_{\alpha,\tau}$. As for $I_{\alpha,\beta}$, it takes $O(\delta\cdot m)$ time to construct and its space complexity is $O(m)$, where $\delta$ is the degeneracy of the graph [6]. To compute $(\alpha,\beta)_\tau$-core using $I_{\alpha,\beta}$, we fetch the vertices of $(\alpha,\beta)$-core and restore the edges in $O(|V(G')|) + O(\sum_{u\in V(G')} deg(u,G'))$ time $(G' = (\alpha,\beta)$-core). Then, we call the peeling algorithm on $(\alpha,\beta)$-core to compute $(\alpha,\beta)_\tau$-core, which takes $O(T_{peel}((\alpha,\beta)\text{-}core))$ time.

The sizes of 2D-indexes can be considered as the projections of $I_{\alpha,\beta,\tau}$ onto 3 planes, as depicted in Figure 3(b). We also summarize the space complexity, build time, and query time of 2D-indexes in Table 2.
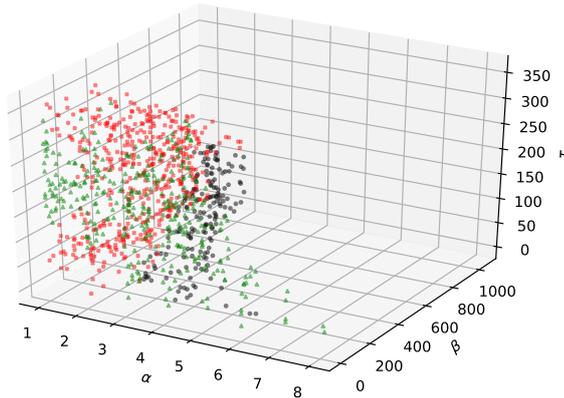


Figure 4: Motivation example for learning-based query processing (`DBpedia-Team`)

**Learning-based hybrid query processing.** As $I_{\alpha,\beta}, I_{\beta,\tau}$, and $I_{\alpha,\tau}$ do not store all the decomposition results like $I_{\alpha,\beta,\tau}$, the construction of these indexes is more time and space-efficient than $I_{\alpha,\beta,\tau}$. However, the reduced index computation inevitably compromises the query processing performance. This is because

the 2D-indexes only store the vertices in $(\alpha, \beta)$-core, $(1, \beta)_\tau$-core, and $(\alpha, 1)_\tau$-core. Clearly, computing the $(\alpha, \beta)_\tau$-core based on these 2D-indexes results in different response time. To better illustrate this point, we plot all parameter combinations on dataset `DBpedia-team` and give each combination a color based on which query processing algorithm performs the best in Figure 4. For ease of presentation, we denote the query processing algorithms based on index $I_{\alpha,\beta}$, $I_{\beta,\tau}$, and $I_{\alpha,\tau}$ as $Q_{\alpha,\beta}$, $Q_{\beta,\tau}$ and $Q_{\alpha,\tau}$ respectively. The red points indicate that $Q_{\beta,\tau}$ is the fastest among the three. The green ones represent the win cases for $Q_{\alpha,\tau}$ and the black points are the cases when $Q_{\alpha,\beta}$ is the best. Evidently, the points of different colors are mingled together and distributed across the parameter space. This suggests that finding simple rules to partition the parameter space is not promising in deciding which of $Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}$ is the fastest. Hence, we formulate it as a classification problem and resort to machine learning techniques to solve this problem.

---

**Algorithm 5:** Hybrid Computation Algorithm

---

    `// Offline training:`
    **Input:** $G$ : Input bipartite graph
    **Output:** Neural network $C : D \rightarrow \{Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}\}$
**1** Build $I_{\beta,\tau}, I_{\alpha,\tau}$ and $I_{\alpha,\beta}$ for $G$
**2** **foreach** $q \in \{N$ random queries$\}$ *on* $G$ **do**
**3**      $feature(q) \leftarrow [\alpha, \beta, \tau$ of $q]$
**4**      $label(q) \leftarrow$ the fastest algorithm in $\{Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}\}$
**5** $X = [feature(q)]$, $q \in N$ queries run on $G$
**6** $Y = [label(q)]$ , $q \in N$ queries run on $G$
**7** $C \leftarrow$ trained neural network on $X, Y$
    `// Online query processing:`
    **Input:** Query parameters: $\alpha, \beta, \tau$
    **Output:** $(\alpha, \beta)_\tau$-core in $G$
**1** $Q_{pred} \leftarrow C.predict(\alpha, \beta, \tau)$
**2** Run $Q_{pred}$ to compute $(\alpha, \beta)_\tau$-core
**3** **return** $(\alpha, \beta)_\tau$-core

---

We introduce a hybrid computation algorithm (Algorithm 5, denoted by $Q_{hb}$), which selects from $\{Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}\}$ based on the query parameters $\alpha, \beta$, and $\tau$. In the offline training phase, we build $I_{\beta,\tau}, I_{\alpha,\tau}$ and $I_{\alpha,\beta}$ on $G$ and obtain the runtime of $Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}$ on $N$ queries, where $N$ is chosen to be less than 5% of all possible queries. The label of a query is the algorithm that responds to it in the shortest time. Then, we train a feed-forward neural network $C$ on the $N$ labeled query instances. In the online query processing phase, given a new query of $(\alpha, \beta)_\tau$-core, the trained neural network makes a prediction based on $\alpha, \beta$, and $\tau$. Then, we use the predicted query processing algorithm to compute $(\alpha, \beta)_\tau$-core.

Here we detail how to train the feed-forward neural network. We impose only one hidden layer in $C$ to avoid over-fitting. The important hyper-parameters of $C$ include the number of hidden units $H$ and the type of optimizer. We use 5-fold cross-validation to evaluate the above hyper-parameters. Specifically, we split the $N$ labeled queries into 5 partitions and each time we take one partition as the validation set and the remainder as the training set. For each parameter setting, we build a classifier on the training sets for 5 times and calculate a performance metric on the validation set. In our model, we define a *time-sensitive error* on the validation set as the performance metric, which calculates a weighted mis-classification cost w.r.t the actual query time. Let $i_k, j_k \in \{1, 2, 3\}$ (encoding of $Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}$) be the predicted class and the actual class of the $k_{th}$ instance. The time-sensitive error is defined as

$$error(i_k, j_k) = e_{i_k}^T \begin{bmatrix} 0 & t_{1,k} - t_{2,k} & t_{1,k} - t_{3,k} \\ t_{2,k} - t_{1,k} & 0 & t_{2,k} - t_{3,k} \\ t_{3,k} - t_{1,k} & t_{3,k} - t_{2,k} & 0 \end{bmatrix} e_{j_k}$$

where $e_{i_k}$ and $e_{j_k}$ are one-hot vectors of length 3 with the $i_k, j_k$ position being 1. $t_{1,k}, t_{2,k}$ and $t_{3,k}$ are the

running time of $Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}$ on the $k_{th}$ instance respectively. The time-sensitive error measures the gap between the predicted query algorithm and the optimal query algorithm. It is averaged over all instances in the validation set and across 5 iterations of cross-validation. Then, the hyper-parameter setting with the lowest time-sensitive error should be chosen. In this way, we are more prone to find the parameter settings that allow us to minimize the query time instead of merely correctly classify each instance.

Note that, training a feed-forward neural network (lines 2 - 7) take significantly less time compared to the 2D-index construction process (line 1) since only $N$ ($N \le 5\%$ of the total number of possible queries) random queries are used.

**Example 4.** *On dataset* `DBpedia-starring`*, given $\alpha$=2, $\beta$=8, and $\tau$=9. $Q_{\alpha,\beta}$ takes 0.53 seconds to find the queried subgraph. $Q_{\beta,\tau}$ and $Q_{\alpha,\tau}$ takes 0.03 and 0.05 respectively. The optimal query processing algorithm on this instance is $Q_{\beta,\tau}$. Accuracy as a performance metric would give equal penalty to mis-classifying $Q_{\alpha,\beta}$ and $Q_{\alpha,\tau}$ as the best algorithm, which is clearly inappropriate. Instead, the time-sensitive error gives penalty of 0.02 if we predict $Q_{\beta,\tau}$ and 0.51 if we predict $Q_{\alpha,\tau}$.*

## 8. Experiments

In this section, we first validate the effectiveness of the $\tau$-strengthened $(\alpha, \beta)$-core model. Then, we evaluate the performance of the index construction algorithms as well as the query processing algorithms.

### 8.1. Experiments setting

**Algorithms.** Our empirical studies are conducted against the following algorithms:
• *Index construction algorithms.* We compare two $I_{\alpha,\beta,\tau}$ construction algorithms: the naive decomposition algorithm `decomp-naive` and the decomposition algorithm with optimizations `decomp-opt`. We also evaluate the index construction algorithms of $I_{\beta,\tau}$ and $I_{\alpha,\tau}$. As for $I_{\alpha,\beta}$, we report its size and build time by running the index construction algorithm in [6].
• *Query processing algorithms.* We use the online computation algorithm presented in Section 4 as the baseline method, denoted as $Q_{bs}$. We compare it to the index-based query processing algorithms $Q_{\alpha,\beta,\tau}, Q_{\alpha,\beta}, Q_{\beta,\tau}$, and $Q_{\alpha,\tau}$, which are based on $I_{\alpha,\beta,\tau}, I_{\alpha,\beta}, I_{\beta,\tau}$, and $I_{\alpha,\tau}$ respectively. We also evaluate the hybrid computation algorithm $Q_{hb}$, which depends on a well-trained classifier and the indexes $Q_{\alpha,\beta}, Q_{\beta,\tau}$, and $Q_{\alpha,\tau}$.

All algorithms are implemented in C++ and the experiments are run on a Linux server with Intel Xeon E3-1231 processors and 16GB main memory. *We end an algorithm if the running time exceeds two hours.*

| Dataset | $|E|$ | $|U|$ | $|L|$ | $\alpha_{max}$ | $\beta_{max}$ | $\tau_{max}$ | $\delta$ |
|---|---|---|---|---|---|---|---|
| Cond-mat (AC) | 58K | 38K | 16K | 37 | 13 | 63 | 8 |
| Writers (WR) | 144K | 135K | 89K | 11 | 82 | 99 | 6 |
| Producers (PR) | 207K | 187K | 48K | 220 | 18 | 219 | 6 |
| Movies (ST) | 281K | 157K | 76K | 19 | 215 | 222 | 7 |
| Location (LO) | 294K | 225K | 172K | 12 | 853 | 852 | 8 |
| BookCrossing (BX) | 434K | 264K | 78K | 376 | 100 | 375 | 13 |
| Teams (TM) | 1.4M | 935K | 901K | 11 | 1063 | 373 | 9 |
| Wiki-en (WC) | 3.80M | 2.04M | 1.85M | 39 | 7659 | 7658 | 18 |
| Amazon (AZ) | 5.74M | 3.38M | 2.15M | 659 | 294 | 658 | 26 |
| DBLP (DB) | 8.6M | 5.4M | 1.4M | 421 | 64 | 420 | 10 |

Table 3: This table reports the basic statistics of 10 real graph datasets.

**Datasets.** We use 10 real graphs in our experiments, which are obtained from the website KONECT [3]. Table 3 includes the statistics of these datasets, sorted by the number of edges in ascending order. The abbreviations of dataset names are listed in parentheses. $|E|$ is the number of edges in the graph. $|U|$ and $|L|$ are the number of vertices in the upper and lower levels. $\alpha_{max}$ is the largest $\alpha$ such that $(\alpha, 1)_1$-core exists. $\beta_{max}$ is the largest $\beta$ such that $(1, \beta)_1$-core exists. $\tau_{max}$ is the largest $\tau$ such that $(1, 1)_\tau$-core exists.

---

[3] http://konect.uni-koblenz.de/networks/
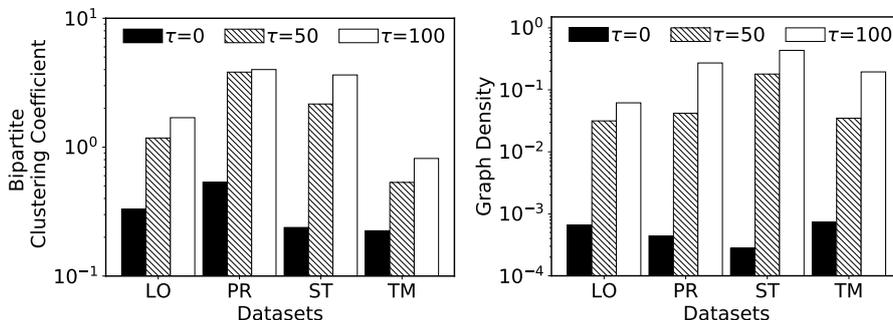
*8.2. Effectiveness Evaluation*



Figure 5: The cohesive metrics comparisons

In this section, we validate the effectiveness of the $\tau$-strengthened $(\alpha, \beta)$-core model. First, we compute some cohesive metrics for $(\alpha, \beta)$-core and $(\alpha, \beta)_\tau$-core. Then, we conduct a case study on dataset `DBLP-2019`.
**Compare $(\alpha, \beta)$-core with $\tau$-strengthened $(\alpha, \beta)$-core.** We compare the graph density and bipartite clustering coefficient for $(\alpha, \beta)$-core and $(\alpha, \beta)_\tau$-core. The graph density [13] of a bipartite graph is calculated as $|E|/(|U| \times |L|)$, where $|E|$ is the number of edges and $|U|$ and $|L|$ are the number of upper and lower vertices. The bipartite clustering coefficient [37] is a cohesive measurement of bipartite networks, which is calculated as $4 \times \bowtie_G / \asymp_G$ where $\asymp_G$ and $\bowtie_G$ are the number of caterpillars (three-path) and the number of butterflies in graph $G$ respectively. In Figure 5, the black bars with $\tau$=0, represents the $(\alpha, \beta)$-core. The shaded bars and the white bars represent the $(\alpha, \beta)_\tau$-core with $\tau$ being 50 and 100 respectively. The engagement constraints $\alpha$ and $\beta$ are set to $0.6\delta$ and $0.4\delta$ respectively, where $\delta$ is the graph degeneracy. As we can see, on all four datasets, the $(\alpha, \beta)_\tau$-core has a higher density and bipartite coefficient than $(\alpha, \beta)$-core. As $\tau$ increases, both of the metrics increase as well. This means that with higher values of $\tau$, we can find subgraphs within $(\alpha, \beta)$-core of higher density and cohesiveness.
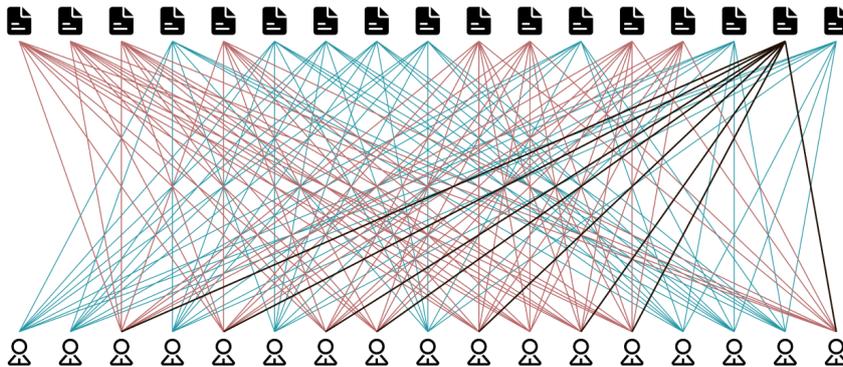


Figure 6: Case study on DBLP-2019

**Case study.** The effectiveness of our model is evaluated through a case study on the DBLP-2019 dataset. The graph in Fig 6 is an $(\alpha, \beta)$-core ($\alpha$=7, $\beta$=8). Given $\tau$=50, $(\alpha, \beta)_\tau$-core excludes the relatively sparse group represented by the light blue lines. The $k$-bitruss ($k$=56) represented by the red lines is in $(\alpha, \beta)_\tau$-core. The black lines are the edges included in $(\alpha, \beta)_\tau$-core but not in $k$-bitruss. The $(\alpha, \beta)_\tau$-core and $k$-bitruss involve the same authors, but $k$-bitruss removes the second last paper on the upper level to enforce the tie strength constraint. Figure 6 implies that: (1) Although $(\alpha, \beta)$-core models vertex engagement via degrees, it fails to distinguish between edges with different tie strength. (2) $k$-bitruss models tie strength via butterfly counting, but it forcefully excludes the weak ties between strongly engaged nodes, which leads to the imprecise estimation of tie strength and failure to include important nodes and their incident edges.

15

(3) $(\alpha, \beta)_\tau$-core considers both vertex engagement and tie strength. Its flexibility allows it to capture unique structures that better resemble the communities in reality.
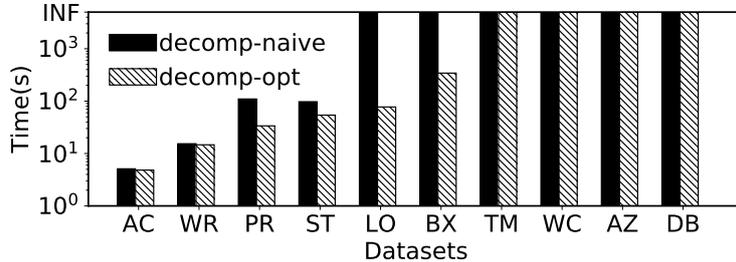


Figure 7: The $I_{\alpha,\beta,\tau}$ construction time

| Data | Index size (MB) | | | | Index construction time (sec) | | | |
|---|---|---|---|---|---|---|---|---|
| | $I_{\alpha,\beta,\tau}$ | $I_{\alpha,\beta}$ | $I_{\beta,\tau}$ | $I_{\alpha,\tau}$ | $I_{\alpha,\beta,\tau}$ | $I_{\alpha,\beta}$ | $I_{\beta,\tau}$ | $I_{\alpha,\tau}$ |
| Cond-mat | 1.29 | 0.78 | 0.26 | 0.50 | 5.11 | 0.11 | 0.68 | 1.37 |
| Writers | 2.14 | 2.24 | 0.93 | 0.24 | 18.81 | 0.24 | 5.01 | 1.31 |
| Producers | 5.55 | 3.16 | 0.37 | 2.43 | 79.43 | 0.38 | 4.99 | 28.37 |
| Movies | 5.26 | 3.51 | 2.01 | 0.53 | 66.18 | 0.46 | 34.29 | 6.26 |
| Location | 68.96 | 4.15 | 33.22 | 0.75 | 77.3978 | 0.36 | 49.5368 | 9.19316 |
| BookCrossing | 33.56 | 5.58 | 3.07 | 9.32 | 342.728 | 1.02 | 48.8196 | 75.2869 |
| Teams | – | 18.42 | 114.17 | 2.44 | time out | 1.94 | 944.102 | 127.265 |
| Wiki-en | – | 46.66 | 945.91 | 10.92 | time out | 9.079 | 3850.51 | 680.569 |
| Amazon | – | 72.96 | 74.47 | 129.21 | time out | 17.184 | 3598.13 | 4731.29 |
| DBLP | – | 112.60 | 29.13 | 159.74 | time out | 19.44 | 559.76 | 3000.08 |

Table 4: Evaluate the size of indexes and their build time.

### 8.3. Performance Evaluation

In this part, we evaluate the efficiency of the index construction algorithms and explore the appropriate hyperparameter settings for the feed-forward neural network that $Q_{hb}$ depends on. Then, we evaluate the efficiency of the query processing algorithms to retrieve $(\alpha, \beta)_\tau$-core.

**Index construction.** First, We compare the build time of $I_{\alpha,\beta,\tau}, I_{\alpha,\beta}, I_{\beta,\tau}$ and $I_{\alpha,\tau}$ on all datasets, as reported in Table 4. The reported build time corresponds to the index construction algorithms with the optimization techniques in Section 6. As shown in 7, although the computation-sharing and the `Bloom-Edge`-index based optimizations effectively reduce the running time, the $I_{\alpha,\beta,\tau}$ still cannot be built within time limit on `Teams`, `Wiki-en`, `Amazon` or `DBLP`. This is because $I_{\alpha,\beta,\tau}$ stores all the decomposition results and it takes the longest time to build, followed by $I_{\beta,\tau}, I_{\alpha,\tau}$, and $I_{\alpha,\beta}$. When the graph is denser on the upper level, $I_{\alpha,\tau}$ takes longer to construct than $I_{\beta,\tau}$. For example, in `DBLP`, where $d_{max}(L)$=119 $< d_{max}(U)$=951, $I_{\beta,\tau}$ is built within 10 minutes while $I_{\alpha,\tau}$ is built in 50 minutes. $I_{\alpha,\beta}$ construction is the fastest as it does not involve any butterfly counting or updating of support. In addition, we report the index sizes in Table 4. The size of $I_{\alpha,\beta,\tau}$ is larger than $I_{\alpha,\beta}, I_{\beta,\tau}$ and $I_{\alpha,\tau}$. In summary, the 2D-indexes that the hybrid computation algorithm depends on are space-efficient and can be built within a reasonable time.

**Tuning hyperparameters for the neural network.** When training the neural network for the hybrid computation algorithm, we choose the hyperparameters (the type of optimizer and size of the hidden layer) that minimizes the time-sensitive error from cross-validation. For each graph $G$, we set the size of hidden layer to 50 and test the *stochastic gradient descent, L-BFGS method* and *Adam* and compare the time-sensitive error. As shown in Figure 8, the L-BFGS method consistently outperforms the other methods and
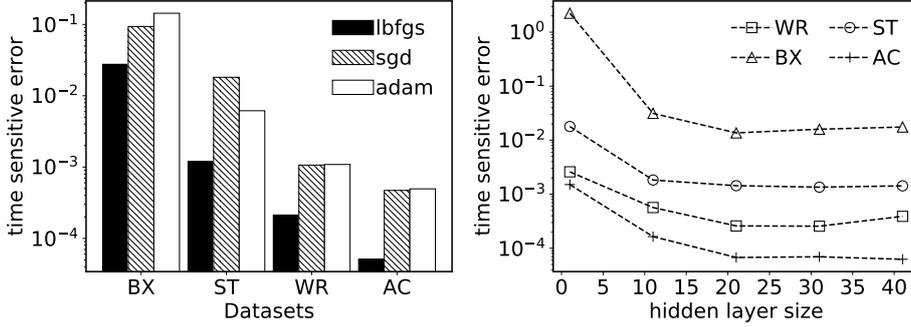
16

Figure 8: Effects of hyperparameters

| Dataset | $Q_{bs}$ | $Q_{\alpha,\beta,\tau}$ | $Q_{\alpha,\beta}$ | $Q_{\beta,\tau}$ | $Q_{\alpha,\tau}$ | $Q_{hb}$ |
|---|---|---|---|---|---|---|
| Cond-mat | 0.139 | 0.004 | 0.030 | 0.008 | 0.009 | 0.006 |
| Writers | 0.406 | 0.011 | 0.069 | 0.011 | 0.009 | 0.006 |
| Producers | 0.560 | 0.022 | 0.200 | 0.047 | 0.037 | 0.029 |
| Movies | 0.871 | 0.023 | 0.235 | 0.028 | 0.041 | 0.027 |
| Location | 11.782 | 0.285 | 7.005 | 1.755 | 2.895 | 0.234 |
| BookCrossing | 44.397 | 0.147 | 13.271 | 2.935 | 4.794 | 1.722 |
| Teams | 59.510 | − | 10.911 | 3.080 | 2.778 | 1.394 |
| Wiki-en | 128.536 | − | 28.589 | 2.638 | 13.613 | 1.775 |
| Amazon | 973.026 | − | 153.085 | 88.673 | 59.228 | 16.432 |
| DBLP | 61.101 | − | 1.843 | 1.321 | 1.055 | 0.269 |

Table 5: This table reports the average response time for all query processing algorithms.

thus is chosen in our model. Then, we explore the effect of the size of the hidden layer on our model. We report the change of time-sensitive error w.r.t varying hidden layer size on dataset `DBpedia-location` and the trends are similar on other datasets. As shown in the plot, 30 hidden units are enough for the classifier built on most tested datasets and beyond this point, more hidden units have little effect on the performance of the model.

**Average query time of** $Q_{\alpha,\beta,\tau}, Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}$ **and** $Q_{hb}$**.** For each algorithm, we report the average response time of 50 randomly generated queries on each dataset in Table 5. As expected, all index based algorithms outperform $Q_{bs}$, as $Q_{bs}$ computes $(\alpha, \beta)_\tau$-core from the input graph. Among them, $Q_{\alpha,\beta,\tau}$ performs the best on most datasets, as it fetches the vertices from $I_{\alpha,\beta,\tau}$ in optimal time and then restores the edges. However, the long build time of $I_{\alpha,\beta,\tau}$ makes $Q_{\alpha,\beta,\tau}$ not scalable to larger graphs like `Team,Wiki-en,Amazon` or `DBLP`. The performances of $Q_{\alpha,\beta}, Q_{\beta,\tau}$, and $Q_{\alpha,\tau}$ differ a lot from each other and across datasets. On average, $Q_{\alpha,\beta}$ is slower than $Q_{\beta,\tau}$ and $Q_{\alpha,\tau}$, because it needs to delete many edges from $(\alpha, \beta)$-core to get $(\alpha, \beta)_\tau$-core especially when $\tau$ is large. As $Q_{hb}$ is trained to pick the fastest from $\{Q_{\alpha,\beta}, Q_{\beta,\tau}, Q_{\alpha,\tau}\}$, it outperforms these 3 algorithms on average in all datasets. In addition, $Q_{hb}$ outperforms the online computation algorithm $Q_{bs}$ by up to two orders of magnitude.

**Evaluate the effect of** $\alpha, \beta$ **and** $\tau$**.** We investigate the effects of varying $\alpha, \beta, \tau$ on each query processing algorithm. We input three types of query streams for all algorithms, each of which increments one of $\alpha$, $\beta$, and $\tau$ while letting the other two parameters be generated randomly. As trends are similar, we only report the results on `Location` and `BookCrossing` in Figure 9. Each data point in the figure represents the average response time of 50 random queries. As expected, index-based query processing algorithms always outperform $Q_{bs}$. The performance of $Q_{\alpha,\beta,\tau}$ is not affected much by the varying parameters, while the 2D-index based algorithms are highly sensitive to them. Each of $Q_{\alpha,\beta}, Q_{\beta,\tau}$, and $Q_{\alpha,\tau}$ tends to perform better when the increased parameter results in a smaller subgraph in the index. For instance, $Q_{\beta,\tau}$ performs better as $\beta$ or $\tau$ increases. In contrast, the hybrid computation algorithm $Q_{hb}$ has very stable performance,
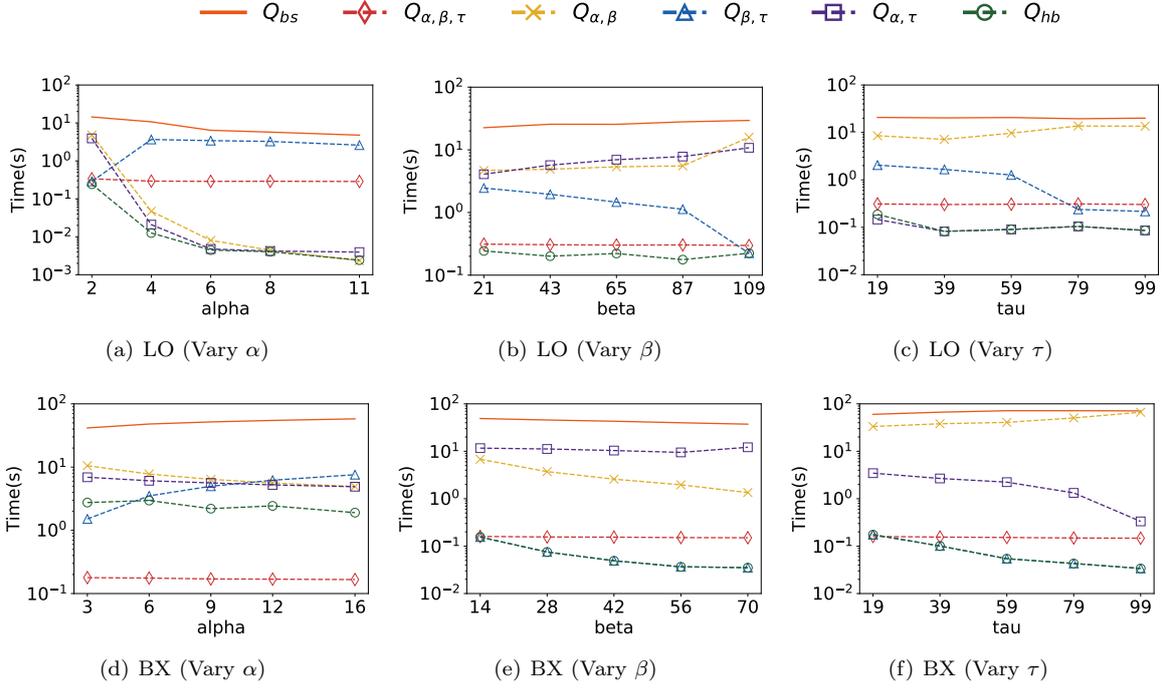
Figure 9: Effects of varying parameters on each query processing algorithm.

as it stays close to and in many cases outperforms the fastest of $Q_{\alpha,\beta}, Q_{\beta,\tau}$, and $Q_{\alpha,\tau}$. In summary, the hybrid computation algorithm $Q_{hb}$ with a well-trained classifier can adjust its querying processing algorithm to different parameters and datasets.

## 9. Conclusion

In this paper, we introduce a novel cohesive subgraph model, $\tau$-strengthened $(\alpha,\beta)$-core, which is the first to consider both tie strength and vertex engagement on bipartite graphs. We propose a decomposition-based index $I_{\alpha,\beta,\tau}$ that can retrieve vertices of any $(\alpha,\beta)_\tau$-core in optimal time. We also apply computation sharing and `BE-Index`-based optimizations to speed up the index construction process of $I_{\alpha,\beta,\tau}$. To balance space-efficient index construction and time-efficient query processing, we propose a learning-based hybrid computation paradigm. Under this paradigm, we introduce three 2D-indexes and train a feed-forward neural network to predict which index is the best choice to process an incoming $(\alpha,\beta)_\tau$-core query. The efficiency of the proposed algorithms and the effectiveness of our model are verified through extensive experiments.

## References

[1] J. Wang, A. P. De Vries, M. J. Reinders, 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, ACM, 2006, pp. 501–508.

[2] M. Ley, The DBLP computer science bibliography: Evolution, research issues, perspectives, in: Proc. Int. Symposium on String Processing and Information Retrieval, 2002, pp. 1–10.

[3] J. C. Brunson, Triadic analysis of affiliation networks, arXiv preprint arXiv:1502.07016.

[4] M. Allahbakhsh, A. Ignjatovic, B. Benatallah, E. Bertino, N. Foo, et al., Collusion detection in online rating systems, in: Asia-Pacific Web Conference, Springer, 2013, pp. 196–207.

[5] A. Beutel, W. Xu, V. Guruswami, C. Palow, C. Faloutsos, Copycatch: stopping group attacks by spotting lockstep behavior in social networks, in: Proceedings of the 22nd international conference on World Wide Web, ACM, 2013, pp. 119–130.

[6] B. Liu, L. Yuan, X. Lin, L. Qin, W. Zhang, J. Zhou, Efficient ($\alpha$,$\beta$)-core computation in bipartite graphs, The VLDB Journal (2020) 1–25.

[7] D. Ding, H. Li, Z. Huang, N. Mamoulis, Efficient fault-tolerant group recommendation using alpha-beta-core, in: Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, ACM, 2017, pp. 2047–2050.

[8] E. Ntoutsi, K. Stefanidis, K. Nørvåg, H.-P. Kriegel, Fast group recommendations by applying user clustering, in: International Conference on Conceptual Modeling, Springer, 2012, pp. 126–140.

[9] M. D. Ornstein, Interlocking directorates in canada: evidence from replacement patterns, Social Networks 4 (1) (1982) 3–25.

[10] D. Palmer, Interlocking directorates and intercorporate coordination, Social Networks: Critical Concepts in Sociology 3 (2002) 261.

[11] S. B. Seidman, Network structure and minimum degree, Social networks 5 (3) (1983) 269–287.

[12] J. Cohen, Trusses: Cohesive subgraphs for social network analysis, National security agency technical report 16 (2008) 3–1.

[13] A. E. Sarıyüce, A. Pinar, Peeling bipartite networks for dense subgraph discovery, in: Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining, ACM, 2018, pp. 504–512.

[14] Z. Zou, Bitruss decomposition of bipartite graphs, in: International Conference on Database Systems for Advanced Applications, Springer, 2016, pp. 218–233.

[15] K. Wang, X. Lin, L. Qin, W. Zhang, Y. Zhang, Efficient bitruss decomposition for large-scale bipartite graphs, in: 2020 IEEE 36th International Conference on Data Engineering (ICDE), IEEE, 2020, pp. 661–672.

[16] J. Cheng, Y. Ke, S. Chu, M. T. Özsu, Efficient core decomposition in massive networks, in: 2011 IEEE 27th International Conference on Data Engineering, IEEE, 2011, pp. 51–62.

[17] W. Khaouid, M. Barsky, V. Srinivasan, A. Thomo, K-core decomposition of large networks on a single pc, Proceedings of the VLDB Endowment 9 (1) (2015) 13–23.

[18] F. Zhang, C. Li, Y. Zhang, L. Qin, W. Zhang, Finding critical users in social communities: The collapsed core and truss problems, IEEE Transactions on Knowledge and Data Engineering.

[19] X. Huang, H. Cheng, L. Qin, W. Tian, J. X. Yu, Querying k-truss community in large and dynamic graphs, in: Proceedings of the 2014 ACM SIGMOD international conference on Management of data, ACM, 2014, pp. 1311–1322.

[20] Y. Shao, L. Chen, B. Cui, Efficient cohesive subgraphs detection in parallel, in: Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, 2014, pp. 613–624.

[21] M. S. Granovetter, The strength of weak ties, in: Social networks, Elsevier, 1977, pp. 347–367.

[22] F. Zhang, L. Yuan, Y. Zhang, L. Qin, X. Lin, A. Zhou, Discovering strong communities with user engagement and tie strength, in: International Conference on Database Systems for Advanced Applications, Springer, 2018, pp. 425–441.

[23] K. Wang, X. Cao, X. Lin, W. Zhang, L. Qin, Efficient computing of radius-bounded k-cores, in: 2018 IEEE 34th International Conference on Data Engineering (ICDE), IEEE, 2018, pp. 233–244.

[24] Y. Zhang, J. X. Yu, Y. Zhang, L. Qin, A fast order-based approach for core maintenance, in: 2017 IEEE 33rd International Conference on Data Engineering (ICDE), IEEE, 2017, pp. 337–348.

[25] F. Bonchi, F. Gullo, A. Kaltenbrunner, Y. Volkovich, Core decomposition of uncertain graphs, in: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2014, pp. 1316–1325.

[26] Y. Peng, Y. Zhang, W. Zhang, X. Lin, L. Qin, Efficient probabilistic k-core computation on uncertain graphs, in: 2018 IEEE 34th International Conference on Data Engineering (ICDE), IEEE, 2018, pp. 1192–1203.

[27] K. Bhawalkar, J. Kleinberg, K. Lewi, T. Roughgarden, A. Sharma, Preventing unraveling in social networks: the anchored k-core problem, SIAM Journal on Discrete Mathematics 29 (3) (2015) 1452–1475.

[28] F. Zhang, W. Zhang, Y. Zhang, L. Qin, X. Lin, Olak: an efficient algorithm to prevent unraveling in social networks, Proceedings of the VLDB Endowment 10 (6) (2017) 649–660.

[29] Z. Zou, R. Zhu, Truss decomposition of uncertain graphs, Knowledge and Information Systems 50 (1) (2017) 197–230.

[30] F. Zhang, Y. Zhang, L. Qin, W. Zhang, X. Lin, Efficiently reinforcing social networks over user engagement and tie strength, in: 2018 IEEE 34th International Conference on Data Engineering (ICDE), IEEE, 2018, pp. 557–568.

[31] S.-V. Sanei-Mehri, A. E. Sariyuce, S. Tirthapura, Butterfly counting in bipartite networks, in: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018, pp. 2150–2159.

[32] J. Wang, A. W.-C. Fu, J. Cheng, Rectangle counting in large bipartite graphs, in: 2014 IEEE International Congress on Big Data, IEEE, 2014, pp. 17–24.

[33] K. Wang, X. Lin, L. Qin, W. Zhang, Y. Zhang, Vertex priority based butterfly counting for large-scale bipartite networks, Proceedings of the VLDB Endowment 12 (10) (2019) 1139–1152.

[34] M. Cerinsek, V. Batagelj, Generalized two-mode cores, CoRR abs/1505.01817. `arXiv:1505.01817`.
URL `http://arxiv.org/abs/1505.01817`

[35] C. Giatsidis, D. M. Thilikos, M. Vazirgiannis, Evaluating cooperation in communities with the k-core structure, in: 2011 International conference on advances in social networks analysis and mining, IEEE, 2011, pp. 87–93.

[36] L. Chang, L. Qin, Cohesive subgraph computation over large sparse graphs, in: 2019 IEEE 35th International Conference on Data Engineering (ICDE), IEEE, 2019, pp. 2068–2071.

[37] S. G. Aksoy, T. G. Kolda, A. Pinar, Measuring and modeling bipartite graphs with community structure, Journal of Complex Networks 5 (4) (2017) 581–603.