



# Fast Local Subgraph Counting

Qiyang Li

The Chinese University of Hong Kong  
Hong Kong, China  
qyli@se.cuhk.edu.hk

Jeffrey Xu Yu

The Chinese University of Hong Kong  
Hong Kong, China  
yu@se.cuhk.edu.hk

## ABSTRACT

We study local subgraph counting queries,  $Q = (p, o)$ , to count how many times a given  $k$ -node pattern graph  $p$  appears around every node  $v$  in a data graph  $G$  when the given center node  $o$  in  $p$  maps to  $v$ . Such local subgraph counting becomes important in GNNs (Graph Neural Networks), where incorporating such counts for every node in  $G$  into the GNN architecture enhances the model’s ability to capture complex relationships within the graph  $G$ . It is challenging to count by subgraph isomorphism, which is known to be NP-hard. In this paper, we propose a novel approach by tree-decomposition-based counting. For a complex pattern graph  $p$  in  $Q$ , we find its best tree decomposition  $T$ , where a node in  $T$  represents a subgraph of  $p$ , and a node in  $p$  may appear in multiple nodes in  $T$ . Let  $p(T)$  be the pattern represented by  $T$ . Our approach is to count  $p(T)$  by homomorphism with a constraint to count the subgraph in every tree node by subgraph isomorphism. We apply symmetry-breaking rules to reduce the cost of counting by subgraph isomorphism for every node in  $T$ , and we develop a new multi-join algorithm to compute such counts. We confirm that our approach on a single machine using a single core can outperform the others significantly.

### PVLDB Reference Format:

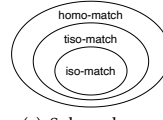
Qiyang Li and Jeffrey Xu Yu. Fast Local Subgraph Counting. PVLDB, 17(8): 1967 - 1980, 2024.  
doi:10.14778/3659437.3659451

### PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/magic62442/subgraph-counting>.

## 1 INTRODUCTION

A graph is a complex structure that has been widely used to support various real-world applications such as social network analysis. Given a large data graph  $G$ , a local subgraph counting query  $Q = (p, o)$  is to count how many times a given  $k$ -node pattern graph  $p$  that appears around every node  $v$  in  $G$  when the center node  $o$  in  $p$  maps to  $v$ . Local subgraph counting has many applications in data mining, such as community detection, graph clustering, network comparison and alignment, anomaly detection, and detecting strong ties in social networks [80, 101], and becomes important in GNNs (Graph Neural Networks), as it provides a large number of topological features. For example, when  $k = 5, 6$  and  $7$ , there are 58, 407 and 4, 306 such pattern graphs respectively. Such local



(a) Subgraph counts

	WS	BY	CA	SG	WW
$\sum_v  \text{HOM}_{p,o}(v) $	86.7	1,767	3,949	1,839	23,936
$\sum_v  \text{tISO}_{p,o,T}(v) $	78.3	1,620	3,817	1,584	16,461
$\sum_v  \text{ISO}_{p,o}(v) $	70.0	1,491	3,547	1,426	14,899
#enumerated match	0.6	6.7	8.1	12.9	232

(b) number of matches( $\times 10^8$ ) of the pattern  $p$  in Fig. 4

Figure 1: Comparing three types of matches

subgraph counts infuse GNNs with higher-order graph structural information. This is useful in tasks where the presence of specific subgraph patterns is indicative of certain properties or labels. In [9], Barceló et al. propose local graph parameter enabled GNNs, and in [76], Qian et al. further study subgraph-enhanced GNNs.

Local subgraph counting is challenging, particularly when it is to count by subgraph isomorphism, which is known to be NP-hard. To compute local subgraph counting queries, there are enumeration-based [72, 96], matrix-based [23, 41, 42, 63–65], and decomposition-based [73, 101] approaches. A survey can be found in [80]. Among the three categories, the state-of-the-art is the decomposition approach (e.g., EVOKE [73] and DISC [101]), which decomposes  $p$  into smaller pattern graphs to count. Here, EVOKE can only support pattern graphs up to 5 nodes [73], whereas DISC is an approach that can handle any  $k$ -node pattern graphs. In brief, DISC is to compute local subgraph counting queries under subgraph isomorphism by homomorphism counting [5, 13, 20, 22, 25, 28, 34, 40]. It computes a query by eliminating homomorphism counts for those that are not subgraph matches from the total homomorphism count for any possible homomorphism matches [101]. As shown in Fig. 1(a), to count by subgraph isomorphism or iso-match (the innermost circle), it counts by homomorphism or homo-matches (the outermost circle) and substrates the counts of the homo-matches between the innermost and outermost circles. The hardness of counting by homomorphism is the same as counting by subgraph isomorphism in general. But, homomorphism can be done efficiently as there are less constraints to check. However, DISC is inefficient in practice. It can only process some selected 6-node pattern graphs in a batch [101]. To the best of our knowledge, there is no reported result that can compute all 407 6-node pattern graphs in a batch together over real large graphs in a reasonable time.

In this paper, we propose a new decomposition approach, called *tree-decomposition-based counting*, to compute local subgraph counting queries for any  $k$ -node pattern graph  $p$  in a novel way. For a complex pattern  $p$ , we find its tree-shaped pattern,  $T$ , by tree decomposition [32], where a node in  $T$  represents a subgraph of  $p$ , and a node in  $p$  may appear in multiple nodes in  $T$ . By tree-decomposition-based counting, we count the entire tree  $T$  by homomorphism counting, where every node in  $p$  that appears in different nodes in  $T$  must map to the same node in the data graph  $G$ , and that every subgraph  $p'$  of  $p$  represented by a node in  $T$  is counted by subgraph isomorphism. We call it tiso-match and illustrate it

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](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 17, No. 8 ISSN 2150-8097.  
doi:10.14778/3659437.3659451

by the middle circle in Fig. 1(a). We subtract the counts of the tiso-matches between the innermost and the middle circle. To show the performance achievement, we show the number of matches for a local subgraph counting query  $Q = (p, o)$ , where  $p$  is a 6-node pattern  $p$  in Fig. 4 and  $o = u_1$  in  $p$ , using five datasets (Table 2) in Fig. 1(b). In Fig. 1(b), the first three rows are the total number of homo-matches, tiso-matches, and iso-matches, to be enumerated for the pattern  $p$ . We get such numbers using an enumeration-based approach. The number of enumerations by tiso-matches is significantly less than that by homo-matches since tiso-matches have more constraints than homo-matches. The actual number of enumerations by our approach is given in the last row in Fig. 1(b). We reduce the unnecessary enumerations significantly since we only enumerate subgraphs represented by tree nodes instead of  $p$ , and we apply symmetry-breaking rules to further reduce the enumeration of these subgraphs.

**Main Contributions:** First, we propose a novel tree decomposition-based counting and prove its correctness. Second, we explore automorphism orbits (structurally equivalent nodes) with symmetry-breaking rules to reduce the cost of finding iso-matches. Symmetry-breaking is an efficient technique for enumerating the iso-matches of an entire pattern graph. Different from existing works, we give a solution that can apply symmetry-breaking rules for tree nodes in a tree decomposition to compute aggregations. In addition, we propose an optimization technique to further reduce the cost. Third, we propose a new multi-join algorithm to compute each tree. Fourth, we have implemented our approach on a single machine, and we confirm that our approach SCOPE can outperform the state-of-the-art approach DISC significantly. We can compute the batch of all 407 6-node queries over real large graphs.

**Organizations.** We give preliminaries and the problem statement in Section 2, and discuss the existing homomorphism-based approach in Section 3. We propose a new decomposition-based approach in Section 4, discuss the new counting in Section 5, and introduce a new multi-join algorithm in Section 6. We discuss related works in Section 7, and confirm the efficiency of our approach in Section 8. We conclude our work in Section 9.

## 2 PRELIMINARIES

Following the notations in [101], we model a simple undirected graph as  $G = (V, E)$ , where  $V$  and  $E$  are the sets of nodes and edges in  $G$ , respectively. A graph  $G' = (V', E')$  is a subgraph of  $G$  if  $V' \subseteq V$  and  $E' \subseteq E$ , and is an induced subgraph of  $G$  if  $E'$  contains all the edges in  $G$  such that both endpoints belonging to  $V'$ . We call a graph a  $k$ -node graph (or  $k$ -graph) if it has  $k$  nodes. The number of nodes and the number of edges are denoted as  $n = |V|$  and  $m = |E|$ .

**Homomorphism & Subgraph Isomorphism:** Let  $G = (V, E)$  be a data graph and  $p = (V_p, E_p)$  be a pattern graph. A function  $f : V_p \mapsto V$  is called a homomorphism of  $p$  if for each edge  $(u, u') \in E_p$ , we have  $(f(u), f(u')) \in E$ , and a homomorphism  $f$  of  $p$  is called a subgraph isomorphism of  $p$  if  $f$  is injective. A subgraph  $G_f = (V_f, E_f)$  in a data graph  $G$  is an *iso-match* (*hom-match*) of a pattern graph  $p$  if  $f$  is a subgraph isomorphism (homomorphism).

**Automorphism Orbit:** For a given graph  $G = (V, E)$ , an automorphism is a bijective function  $\gamma : V \mapsto V$  such that  $(v, v') \in E$  iff  $(\gamma(v), \gamma(v')) \in E$ . Here, an automorphism  $\gamma$  maps  $G$  to itself in a structure-preserving manner. The set of automorphisms, denoted

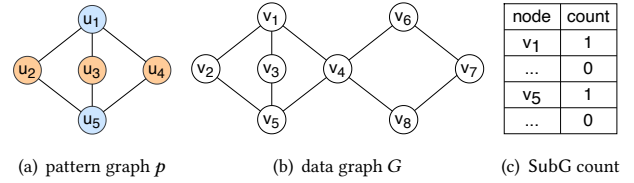


Figure 2: An Example of  $Q = (p, o)$  for  $o = u_1$

as  $\text{Aut}(G)$ , forms a group called the automorphism group of  $G$ . With  $\text{Aut}(G)$ , two nodes  $v, v' \in V$  are in an equivalence relation iff there exists an automorphism  $\gamma$  such that  $\gamma(v) = v'$  (or  $v \rightarrow v'$ ) for  $\gamma \in \text{Aut}(G)$ . Such equivalence relations partition nodes into equivalence classes, where an equivalence class is called an automorphism orbit, and is denoted as  $\vartheta$ . Below, we use  $o$  to denote a representative node in  $\vartheta$ , which we call an orbit.

**Problem Statement:** In this work, we study local subgraph counting queries by subgraph isomorphism. Here, a local subgraph counting query  $Q$  is defined as  $Q = (p, o)$ , where  $p$  is a connected pattern graph and  $o$  is an orbit of  $p$ . The local subgraph count for a specific node  $v$  in a data graph  $G$ , denoted as  $|\text{Sub}_{G,p,o}(v)|$ , is the count of all iso-matches of  $p$  that match  $v$  to  $o$  such that  $f(o) = v$ , where it only counts one for the iso-matches of  $p$  that are induced from the same edges. We compute  $|\text{Sub}_{G,p,o}(v)|$  for every node  $v$  in  $G$ .

It is known that  $|\text{Sub}_{G,p,o}(v)|$  can be processed by *local subgraph isomorphism counting* ( $|\text{ISO}_{p,o}(v)|$ ) or *local homomorphism counting* ( $|\text{HOM}_{p,o}(v)|$ ), where  $\text{ISO}_{p,o}(v)$  and  $\text{HOM}_{p,o}(v) = \{f \mid f \text{ are subgraph isomorphism and homomorphism of } p \text{ with } f(o) = v\}$ , respectively. It is important to mention that  $|\text{Sub}_{G,p,o}(v)| = |\text{ISO}_{p,o}(v)| \cdot |\vartheta|/|\text{Aut}(p)|$ , as shown in [101]. We explain it in Example 2.1.

**Example 2.1:** Fig. 2 shows an example with a pattern graph  $p = (V_p, E_p)$ , a data graph  $G = (V, E)$ , and the  $\text{Sub}_{G,p,o}(v_i)$  for every  $v_i$  in  $G$  for a given subgraph counting query  $Q = (p, o)$ , where  $o = u_1$ . The following two mappings,  $f_i$  and  $f_j$  from  $V_p$  to  $V$ ,  $f_i = \{u_1 \rightarrow v_1, u_2 \rightarrow v_2, u_3 \rightarrow v_3, u_4 \rightarrow v_4, u_5 \rightarrow v_5\}$  and  $f_j = \{u_1 \rightarrow v_1, u_2 \rightarrow v_2, u_3 \rightarrow v_3, u_4 \rightarrow v_4, u_5 \rightarrow v_1\}$ , are homomorphisms, and only  $f_i$  is subgraph isomorphism, as both  $u_1$  and  $u_5$  in  $f_j$  map to the same node  $v_1$  in  $G$ . With  $f_i$  and  $f_j$ , there is an iso-match,  $G_{f_i}$ , and a homo-match,  $G_{f_j}$ , in  $G$ .

The automorphism group of  $p$ ,  $\text{Aut}(p)$ , is in size of  $|\text{Aut}(p)| = 12$ . As an example, one automorphism  $\gamma$  is  $(u_1, u_5)(u_2, u_3, u_4)$  based on the representation of the product of disjoint cycles regarding permutation. For example, consider the disjoint cycle of  $(u_2, u_3, u_4)$ , it refers to  $u_2 \rightarrow u_3, u_3 \rightarrow u_4$ , and  $u_4 \rightarrow u_2$ . Some other examples are  $(u_1)(u_2)(u_3)(u_4)(u_5)$  and  $(u_1, u_5)(u_2)(u_3)(u_4)$ , where the former is an identity automorphism (or trivial automorphism) by which  $u_i \rightarrow u_i$  for every  $u_i$ . By  $\text{Aut}(p)$ , we have two automorphism orbits,  $\vartheta_1 = \{u_1, u_5\}$  and  $\vartheta_2 = \{u_2, u_3, u_4\}$ .

There are 6 iso-matches,  $G_{f_k}$  of  $p$  in  $G$  that map  $u_1$  to  $v_1$  based on 6 subgraph isomorphisms,  $f_k$ , for  $1 \leq k \leq 6$ , in which  $u_1$  maps to  $v_1$ ,  $u_5$  maps to  $v_5$ ,  $u_2$  maps to any of the 3 nodes in  $\{v_2, v_3, v_4\}$ ,  $u_3$  maps to any of the 2 nodes in  $\{v_2, v_3, v_4\}$  that  $u_2$  does not map to, and  $u_4$  maps to one left that both  $u_2$  and  $u_3$  do not map to. Therefore,  $|\text{ISO}_{p,o}(v_1)| = 6$ , where the 6 iso-matches are induced from the same node set  $\{v_1, v_2, v_3, v_4, v_5\}$ . As observed in Fig. 2, we have  $|\text{Sub}_{G,p,o}(v_1)| = 1$  by  $|\text{Sub}_{G,p,o}(v_1)| = |\text{ISO}_{p,o}(v_1)| \cdot |\vartheta|/|\text{Aut}(p)|$ .

Here,  $|\text{ISO}_{p,o}(v_1)| = 6$ ,  $|\vartheta| = 2$  for  $\vartheta = \{u_1, u_5\}$  where  $o = u_1$  is an orbit that maps to  $v_1$ , and  $|\text{Aut}(p)| = 12$ .  $\square$

### 3 AN APPROACH BY HOMOMORPHISM

HOM queries are considered faster to process than ISO queries due to the fact that homomorphisms are not required to be injective. In [101], it shows that, for a pattern graph  $p$  and a node orbit  $o$  of  $p$ , an ISO query  $Q$  can be systematically divided into a set of HOM queries ( $p' = (V_{p'}, E_{p'})$ ,  $o'$ ) with  $|V_{p'}| \leq |V_p|$  to process.

We present DISC [101], which is a general approach by homomorphism to compute SubG queries in brief. First, the count of  $\text{ISO}_{p,o}(v)$  is less than or equal to the count of  $\text{HOM}_{p,o}(v)$ , as HOM is less restrictive than ISO. Hence, the count of  $\text{HOM}_{p,o}(v)$  is equal to the count of  $\text{ISO}_{p,o}(v)$  plus the count of  $\nabla_{p,o}(v)$ , where  $\nabla_{p,o}(v)$  is the set of all non-injective homomorphisms  $f$  of  $p$  with  $f(o) = v$ . Note that by non-injective homomorphism, two nodes in  $p$  map to the same node in data graph  $G$ . Second, it shows that  $|\nabla_{p,o}(v)|$  is the sum of  $|\text{ISO}_{p',o(p')}(v)|$  for any subpattern  $p'$  of  $p$ , where  $o(p')$  denotes the node in  $p'$  that contains  $o$ , as there is a 1:1 connection between a subpattern  $p'$  of  $p$  and a non-injective homomorphism. Third, the count of  $\text{HOM}_{p,o}(v)$  is the sum of  $|\text{ISO}_{p',o(p')}(v)|$  for any subpattern  $p'$  of  $p$  and  $p$  itself.

$$|\text{HOM}_{p,o}(v)| = \sum_{p' \in \{p\} \cup \text{Sub}(p)} |\text{ISO}_{p',o(p')}(v)| \quad (1)$$

With the assistance of the Möbius inversion formula,  $|\text{ISO}|$  can be computed by  $|\text{HOM}|$  as follows.

$$|\text{ISO}_{p,o}(v)| = \sum_{p' \in \{p\} \cup \text{Sub}(p)} \mu(p, p') \cdot |\text{HOM}_{p',o(p')}(v)| \quad (2)$$

where  $\mu(p, p')$  is the corresponding Möbius function and  $\text{Sub}(p)$  is the set of all subpatterns of  $p$ .

The key idea behind DISC is to process  $|\text{HOM}_{p,o}(v)|$  for a pattern graph  $p = (V_p, E_p)$  by relational algebra efficiently using database techniques as follows, supposing that an undirected data graph  $G$  is stored in an edge relation  $R_G(\text{from}, \text{to})$ .

$$oY_{\text{count}(\ast)}(R_1 \bowtie R_2 \bowtie \dots \bowtie R_{|E_p|}) \quad (3)$$

In Eq. (3),  $R_i$  is a renamed relation of  $R_G$  for an edge  $e_i \in E_p$  of  $p$ , for  $1 \leq i \leq |E_p|$ , and  $oY_F(\cdot)$  is to aggregate using the function  $F$  over each group grouped by the group-by attribute  $o$ .

To process Eq. (3) for a specific pattern graph  $p$  by HOM, DISC decomposes a complex cyclic join  $(R_1 \bowtie R_2 \bowtie \dots \bowtie R_{|E_p|})$  to a join tree based on tree decomposition [32]. With the tree decomposition  $T$ , DISC processes the complex joins  $(R_1 \bowtie R_2 \bowtie \dots \bowtie R_{|E_p|})$  in two steps. In the first step, it processes the joins for each node  $\tau$  in the join tree, which results in an intermediate relation  $\text{rel}(\tau)$ , and in the second step it processes the joins over all the intermediate relations. Following the AGM bound [8], the size of the intermediate result  $(\text{rel}(\tau))$  for a node  $\tau$  is bounded by  $|R_G|^{\text{fhw}}$ , where fhw is called the fractional hypertree width and is the minimum real number such that every node  $\tau$  in the join tree  $T$  has a fractional edge cover of weight fhw. By pushing down group-by and aggregations, the time complexity of processing  $T$  is  $O(|V_p| \cdot |R_G|^{\text{fhw}})$  in [101].

The workflow of DISC [101] is depicted in the upper part in Fig. 3.

① The ISO count of  $p$  equals to the HOM count of  $p$  minus the

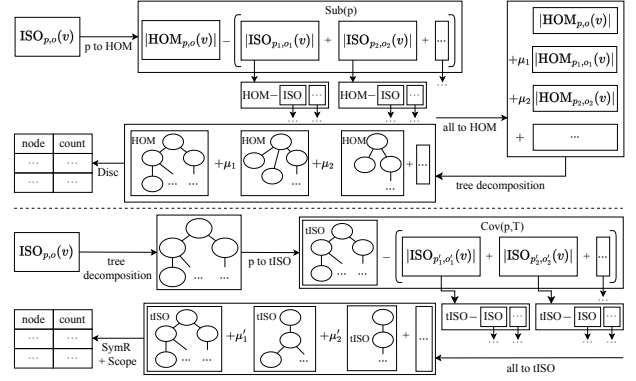


Figure 3: An overview of DISC and SCOPE

ISO count of a set of subpatterns ( $\text{Sub}(p)$ ). ② The ISO count for a subpattern,  $p'$ , in  $\text{Sub}(p)$  is counted in a similar manner recursively. Hence,  $\text{ISO}_p(o, v)$  becomes a linear combination of the HOM counts of patterns. ③ DISC uses tree decomposition [32] to compute HOM count for each of such patterns with joins and aggregations. ④ DISC proposes a multi-join algorithm Disc to compute joins and aggregations for each tree decomposition.

### 4 A NEW APPROACH BY PARAMETERIZED SUBGRAPH ISOMORPHISM

Subgraph isomorphism is studied in parameterized complexity [21, 61]. Here, a fixed-parameter algorithm is to deal with an NP-hard problem in running time  $f(\kappa) \cdot n^c$  where  $\kappa$  is a parameter or a set of parameters,  $f(\cdot)$  is a computable function which can be exponential on  $\kappa$ , and  $c$  is a constant which is independent of  $\kappa$  and  $n$ . Mark and Phlipczuk in [61] investigate the different parameters on the complexity of parameterized subgraph isomorphisms. In this paper, we focus on a parameterized algorithm by treewidth based on tree decomposition on the pattern graph  $p$  for local subgraph counting.

**Definition 4.1: (Tree Decomposition)** Given an undirected pattern graph  $p = (V_p, E_p)$ , a tree decomposition of  $p$  is a tree  $T = (V_T, E_T)$ . We use  $\tau_i$  (simply  $\tau$ ) to denote a node in  $V_T$ , where  $\tau_i$  maintains a nonempty subset of  $V_p$ , denoted as  $V(\tau_i)$  ( $\subseteq V_p$ ), with which an induced subgraph of  $p$  can be constructed, denoted as  $p(\tau_i)$ . We say a node  $v$  in  $V_p$  appears in  $\tau_i$  if  $v \in V(\tau_i)$ . The three conditions on  $T$  are as follows. (1) Every node in  $p$  is covered by  $T$  such that  $V_p = \bigcup_{\tau_i \in V_T} V(\tau_i)$ . (2) Every edge in  $p$  is covered by  $T$  such that for every edge  $(u, v) \in E_p$ , both  $u$  and  $v$  appear in at least one  $\tau_i$ . (3) Nodes in  $T$  are connected if they all contain a pattern node. That is, if a node  $v \in V_p$  appears in both  $\tau_i$  and  $\tau_j$ , then  $v$  appears in every  $\tau_k$  on the path that connects  $\tau_i$  and  $\tau_j$  in  $T$ .

Given a tree decomposition  $T$  for a query  $Q = (p, o)$ , the root node of  $T$  is selected from a node that contains  $o$ . We denote a subtree rooted at node  $\tau_i$  as  $T_i$ . The subtree  $T_i$  represents a subgraph of  $p$ , denoted as  $p_i$  (or  $p(T_i)$ ), which is an induced subgraph of  $p$  over  $\bigcup_{\tau_j \in T_i} V(\tau_j)$ . We also use  $\text{parent}(\tau_i)$  to denote the parent node of  $\tau_i$  in  $T$ , and  $V_C(\tau_i)$  to denote the common nodes that appear in both  $V(\tau_i)$  and  $V(\text{parent}(\tau_i))$ , for a node  $\tau_i$  in  $T$ . An induced subgraph  $G(V_C(\tau_i))$ , denoted as  $G_C(\tau_i)$ , can be constructed on  $V_C(\tau_i)$ .

**Treewidth:** The width of a tree decomposition  $T$  is the largest size

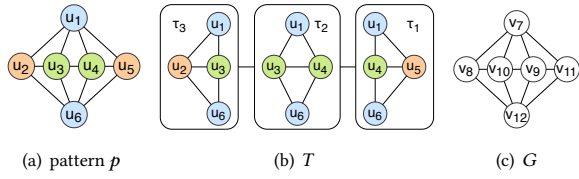


Figure 4: A pattern  $p$ , a tree decomposition  $T$ , and a graph  $G$

of  $|V(\tau_i)|$  minus 1. There are several other metrics defined, including generalized hypertree width and fractional hypertree width [32]. Our approach can use any width. The treewidth of  $p$ , denoted as  $\text{tw}(p)$ , is the smallest width over all possible tree decompositions for  $G$ . Below, we use  $\text{tw}(p)$  to denote a tree decomposition  $T$  whose width is  $\text{tw}(p)$ . It is known that subgraph isomorphism for a pattern graph  $p = (V_p, E_p)$  and a data graph  $G = (V, E)$  can be solved in time  $2^{O(|V_p|)} \cdot |V|^{O(\text{tw}(p))}$  in [31, 61].

**Example 4.1:** Consider  $\text{ISO}_{p,o}(v)$  where  $p$  is a 6-node pattern shown in Fig. 4(a) and  $o = u_1$ . The tree decomposition  $T$  for  $p$  is shown in Fig. 4(b). In  $T$ , there are 3 nodes,  $\tau_i$ , for  $1 \leq i \leq 3$ , that represent three subgraphs,  $p(\tau_i)$ . The root of  $T$  is  $\tau_3$  as it contains the orbit  $u_1$ . The subtree  $T_1$  is  $\tau_1$ , the subtree  $T_2$  is the subtree of  $T$  rooted at  $\tau_2$ , and the subtree  $T_3 = T$ . Here,  $\text{parent}(\tau_1) = \tau_2$ ,  $\text{parent}(\tau_2) = \tau_3$ . The common nodes of  $\tau_3$  and  $\tau_2$  with its parent are  $V_C(\tau_1) = \{u_1, u_4, u_6\}$ , and  $V_C(\tau_2) = \{u_1, u_3, u_6\}$ , respectively; and  $G_C(\tau_1)$  and  $G_C(\tau_2)$  are simple paths,  $u_1-u_4-u_6$  and  $u_1-u_3-u_6$ , respectively. The data graph  $G$  is also shown in Fig. 4(c).

**TISO-based counting:** We propose *tree-decomposition-based counting*. Let  $T_i$  be a subtree in a tree decomposition  $T$  for a pattern graph  $p$ . A tiso-match over  $p(T_i)$  is a homo-match of  $p(T_i)$  in  $G$  on the condition that  $p(\tau_j)$  are iso-matches for every  $\tau_j \in T_i$ . Below, we use  $|\text{tISO}_{p,o,T}(v)|$  to denote the number of tiso-matches that match a node  $v$  in  $G$  to  $o$ , given a tree decomposition  $T$ .

**Lemma 4.1:** For a query  $Q = (p, o)$ ,  $|\text{ISO}_{p,o}(v)| \leq |\text{tISO}_{p,o,T}(v)| \leq |\text{HOM}_{p,o}(v)|$  for every node  $v$  in  $G$ .

**Proof Sketch:** First, some nodes,  $u$  and  $u'$ , in  $p$  that map to the same node in  $G$  by HOM cannot map to the same node by tISO if both appear in a node  $\tau$  in  $T$  due to iso-matches in  $\tau$ . We have  $|\text{tISO}_{p,o,T}(v)| \leq |\text{HOM}_{p,o}(v)|$ . Second, some nodes,  $u$  and  $u'$ , that do not appear in any node  $\tau$  in  $T$  together may map to the same node in  $G$  by tISO, which cannot happen by iso-matches. We have  $|\text{ISO}_{p,o}(v)| \leq |\text{tISO}_{p,o,T}(v)|$ .  $\square$

Next, we discuss how to compute  $|\text{tISO}|$  by enumerating iso-matches for every tree node  $\tau$  in  $T$ . Here, for every  $\tau$  in  $T$ , we find iso-matches of  $\tau$  in  $G$ , and maintain it in a relation  $\mathbf{R}_\tau = \text{ISO}(p(\tau))$ . For a leaf node  $\tau$  in  $T$ , we keep the count of tiso-matches in a relation  $\mathbf{X}(V_C(\tau), C)$  as follows.

$$\mathbf{X}_\tau(V_C(\tau), C) = V_{C(\tau)} \mathcal{Y}_{\text{count}(\ast) \rightarrow C}(\mathbf{R}_\tau) \quad (4)$$

Here  $\mathcal{Y}_{h \rightarrow A}(\cdot)$  is to apply an aggregate function  $h$  for a group  $v$  and rename the output of  $h$  as  $A$ . For a non-leaf node  $\tau$  in  $T$ , let  $\mathbf{J}_\tau$  be as follows.

$$\begin{aligned} \mathbf{J}_\tau &= \mathbf{R}_\tau \bowtie_{V_C(\tau_1)} \rho_{C \rightarrow C_1}(\mathbf{X}_{\tau_1}) \bowtie_{V_C(\tau_2)} \rho_{C \rightarrow C_1}(\mathbf{X}_{\tau_2}) \\ &\quad \cdots \bowtie_{V_C(\tau_k)} \rho_{C \rightarrow C_k}(\mathbf{X}_{\tau_k}) \end{aligned} \quad (5)$$

	$\mathbf{J}_{\tau_1} = \mathbf{R}_{\tau_1}$					$\mathbf{X}_{\tau_1}$					$\mathbf{R}_{\tau_2}$					$\mathbf{J}_{\tau_2} = \mathbf{R}_{\tau_2} \bowtie \mathbf{X}_{\tau_1}$					$\mathbf{X}_{\tau_2}$					$\mathbf{X}_{\tau_1}(u_1)$				
	$u_1$	$u_4$	$u_5$	$u_6$	$u_6$	$u_1$	$u_4$	$u_6$	$C$	$u_1$	$u_3$	$u_4$	$u_6$	$u_1$	$u_3$	$u_4$	$u_6$	$C_1$	$u_1$	$u_3$	$u_6$	$C$	$u_1$	$u_3$	$u_6$	$C$	$u_1$	$C$		
$f_1$	$v_8$	$v_7$	$v_4$	$v_9$	$v_9$	$v_8$	$v_7$	$v_9$	1	$v_8$	$v_7$	$v_{10}$	$v_9$	$v_8$	$v_7$	$v_{10}$	$v_9$	2	$v_8$	$v_7$	$v_9$	2	$v_8$	$v_7$	$v_9$	2	$v_7$	16		
$f_2$	$v_8$	$v_{10}$	$v_7$	$v_9$	$v_9$	$v_8$	$v_{10}$	$v_9$	2	$v_8$	$v_{10}$	$v_7$	$v_9$	$v_8$	$v_{10}$	$v_7$	$v_9$	1	$v_8$	$v_{10}$	$v_9$	2	$v_8$	$v_{10}$	$v_9$	2	$v_8$	8		
$f_3$	$v_8$	$v_{10}$	$v_{12}$	$v_9$	$v_9$	$v_8$	$v_{12}$	$v_9$	1	$v_8$	$v_{10}$	$v_{12}$	$v_9$	$v_8$	$v_{10}$	$v_{12}$	$v_9$	1	$v_8$	$v_{12}$	$v_9$	2	$v_9$	$v_7$	$v_8$	2	$v_{10}$	8		
$f_4$	$v_8$	$v_{12}$	$v_{10}$	$v_9$	$v_9$	$v_8$	$v_{12}$	$v_{10}$	$v_9$	$v_8$	$v_{12}$	$v_{10}$	$v_9$	$v_8$	$v_{12}$	$v_{10}$	$v_9$	2	$v_9$	$v_7$	$v_8$	2	$v_{10}$	8						
$f_5$	$v_9$	$v_7$	$v_{10}$	$v_8$	$v_8$	$v_9$	$v_7$	$v_{10}$	$v_8$	$v_9$	$v_7$	$v_{10}$	$v_8$	$v_9$	$v_7$	$v_{10}$	$v_8$	2	$v_9$	$v_{10}$	$v_8$	2	$v_{11}$	8						
$f_6$	$v_9$	$v_{10}$	$v_7$	$v_8$	$v_8$	$v_9$	$v_{10}$	$v_7$	$v_8$	$v_9$	$v_{10}$	$v_7$	$v_8$	$v_9$	$v_{10}$	$v_7$	$v_8$	1	$v_9$	$v_{12}$	$v_8$	2	$v_{10}$	8						
$f_7$	$v_9$	$v_{10}$	$v_{12}$	$v_8$	$v_8$	$v_9$	$v_{10}$	$v_{12}$	$v_8$	$v_9$	$v_{10}$	$v_{12}$	$v_8$	$v_9$	$v_{10}$	$v_{12}$	$v_8$	1	...	...	...	...	...	...	...	...	...	...		
$f_8$	$v_9$	$v_{12}$	$v_{10}$	$v_8$	$v_8$	$v_9$	$v_{12}$	$v_{10}$	$v_8$	$v_9$	$v_{12}$	$v_{10}$	$v_8$	$v_9$	$v_{12}$	$v_{10}$	$v_8$	2	...	...	...	...	...	...	...	...	$v_{12}$	16		

Figure 5:  $|\text{tISO}_{p,o,T}(v)|$  Computing

where  $\tau_i$  for  $1 \leq i \leq k$  is a child of  $\tau$  such that  $\tau = \text{parent}(\tau_i)$ , a join is a join on the common nodes of  $V_C(\tau_i)$ , and  $\rho_{C \rightarrow C'}(X)$  is a rename operator to rename  $C$  in  $X$  to be  $C'$ . Then, for a non-leaf node  $\tau$  that is not the root of  $T$ , we have

$$\mathbf{X}_\tau(V_C(\tau), C) = V_{C(\tau)} \mathcal{Y}_{\text{sum}(C_1 \times C_2 \times \cdots \times C_k) \rightarrow C}(\mathbf{J}_\tau) \quad (6)$$

For the root  $\tau$  of  $T$ , we have

$$\mathbf{X}_\tau(o, C) = \mathcal{O}_{\text{sum}(C_1 \times C_2 \times \cdots \times C_k) \rightarrow C}(\mathbf{J}_\tau) \quad (7)$$

where  $o$  is the orbit  $o$  of the given query  $Q = (p, o)$ . It is important to note that  $|\text{tISO}_{p,o,T}(v)|$  for every node  $v$  in  $G$  is the tuple in the table  $\mathbf{X}_\tau(o, C)$  for  $v = f(o)$  when  $\tau$  is the root of  $T$ . As a special case, when there is only one node  $\tau$  in  $T$  which is both the root and the leaf, we compute it as  $\mathbf{X}_\tau(o, C) = \mathcal{O}_{\text{count}(\ast) \rightarrow C}(\mathbf{R}_\tau)$ .

**Example 4.2:** Reconsider Example 4.1. We show  $|\text{tISO}_{p,o,T}(v)|$  computing for every node  $v$  in  $G$  in Fig. 5, focusing on when  $f(u_1) = v_8$  or  $f(u_1) = v_9$ . ① For  $\tau_1$ ,  $\mathbf{R}_{\tau_1}$  shows the iso-matches of  $V_{\tau_1}$ ,  $\mathbf{X}_{\tau_1}$  shows the tISO-counts of  $p(\tau_1)$  that have the same matches to  $V_C(\tau_1) = \{u_1, u_4, u_6\}$  based on Eq. (4) as  $\tau_1$  is the leaf node in  $T$ . For example, the 2nd tuple of  $(v_8, v_{10}, v_9, 2)$  in  $\mathbf{X}_{\tau_1}$  shows that there are 2 iso-matches of  $p(\tau_1)$  that contain the same three nodes in  $G$  by  $\{u_1 \rightarrow v_8, u_4 \rightarrow v_{10}, u_6 \rightarrow v_9, u_5 \rightarrow v_7\}$  and  $\{u_1 \rightarrow v_8, u_4 \rightarrow v_{10}, u_6 \rightarrow v_9, u_5 \rightarrow v_{12}\}$ . ② For  $\tau_2$ ,  $\mathbf{R}_{\tau_2}$  shows the iso-matches of  $p(\tau_2)$  in  $G$ . In  $\mathbf{R}_{\tau_2}$ , there are 8 iso-matches of  $p(\tau_2)$  in  $G$  when  $u_1$  matches either  $v_8$  or  $v_9$ . Following Eq. (5),  $\mathbf{J}_{\tau_2} = \mathbf{R}_{\tau_2} \bowtie \mathbf{X}_{\tau_1}$ . In a similar manner,  $\mathbf{X}_{\tau_2}$  shows the tiso-counts of  $p(\tau_2)$  grouped by  $V_C(\tau_2) = \{u_1, u_3, u_6\}$  based on Eq. (6) as  $\tau_2$  is the non-root/non-leaf node in  $T$ . ③ In a similar manner we can compute  $\tau_3$ .

We discuss how to compute  $|\text{ISO}|$  by  $|\text{tISO}|$  given a tree decomposition  $T$  for a pattern graph  $p$ . As shown in Eq. (1),  $|\text{HOM}_{p,o}(v)|$  is the sum of  $|\text{ISO}_{p,o}(v)|$  and  $|\text{ISO}_{p',o}(v)|$  for any subpattern  $p'$  of  $p$ , where the set of subpatterns of  $p$ ,  $\text{Sub}(p)$ , is defined in [101] (Definition 4.2). In brief, a subpattern of  $p = (V_p, E_p)$ ,  $p' = (V_{p'}, E_{p'})$ , is a pattern with less number of nodes ( $|V_{p'}| < |V_p|$ ), due to the reason that some nodes in  $p$  map to the same node in  $p'$  by homomorphism. We follow [101] to describe such homomorphism below. In [101], it specifies  $V_{p'}$  by a partition of nodes in  $V_p$ , denoted as  $\mathcal{I} = \{I_1, I_2, \dots, I_{|V_{p'}|}\}$ , where  $I_k$ , for  $1 \leq k \leq |V_{p'}|$ , is an independent subset of  $V_p$  so that there are no edges in  $p$  between any two nodes in  $I_k$ . In other words, all the nodes in an  $I_k$  may map to the same node by homomorphism. A node  $v_k$  in  $V_{p'}$  corresponds to one distinct  $I_k$ . An edge  $(v_i, v_j)$  exists in  $p'$  if there is an edge in  $p$  between some node in  $I_i$  and some node in  $I_j$ . As there is a one-to-one connection between a subpattern  $p'$  and a partition  $\mathcal{I}$ , it becomes possible to discuss subpatterns by such partitions. Below, we use  $\mathbb{I} = \{\mathcal{I}_1, \mathcal{I}_2, \dots\}$  to denote all subpatterns,  $p_i$ , in  $\text{Sub}(p)$ .

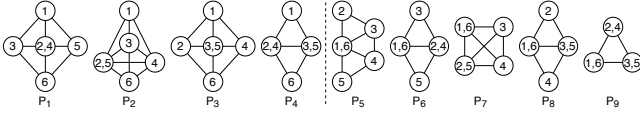


Figure 6: Cov and non-Cov of  $p$  with  $T$  in Fig. 4

**Definition 4.2:** Given a tree decomposition  $T$  on  $p$ , let  $p_i$  be a subpattern of  $p$ , where its corresponding partition is  $\mathcal{I}_i = \{I_1, I_2, \dots\}$  in  $\mathbb{I}$ . We say  $T$  covers  $p_i$ , if any two nodes,  $u$  and  $v$ , in the same  $I_k$  do not appear together in any node  $\tau$  in  $T$ .

By Definition 4.2, the set of subpatterns of  $p$  covered by  $T$  is a subset of  $\text{Sub}(p)$ , which we denote as  $\text{Cov}(p, T)$  ( $\subseteq \text{Sub}(p)$ ).

**Example 4.3:** Given a pattern  $p$  and its tree decomposition  $T$  in Fig. 4, the 9 subpatterns of  $p$ ,  $\text{Sub}(p)$ , are shown in Fig. 6, and the first 4 subpatterns are covered by  $T$  such that  $\text{Cov} = \{p_1, p_2, p_3, p_4\}$ . Such 4 subpatterns covered by  $T$  are induced by the partitions  $\mathcal{I}_1 = \{\{u_2, u_4\}, \{u_1\}, \{u_3\}, \{u_5\}, \{u_6\}\}$ ,  $\mathcal{I}_2 = \{\{u_2, u_5\}, \{u_1\}, \{u_3\}, \{u_4\}, \{u_6\}\}$ ,  $\mathcal{I}_3 = \{\{u_2\}, \{u_1\}, \{u_3, u_5\}, \{u_4\}, \{u_6\}\}$ ,  $\mathcal{I}_4 = \{\{u_2, u_4\}, \{u_1\}, \{u_3, u_5\}, \{u_6\}\}$ , respectively.

**Proposition 4.1:** For a given pattern graph  $p$ , an orbit  $o$ , and a tree decomposition  $T$ , we have

$$|\text{HOM}_{p,o}(v)| - |\text{tISO}_{p,o,T}(v)| = \sum_{p' \in \text{Sub}(p) \setminus \text{Cov}(p,T)} |\text{ISO}_{p',o}(p')(v)| \quad (8)$$

**Proof Sketch:** It can be proved in a similar way as to prove Proposition 4.2 in [101]. We give the proof sketch below. To prove the  $\leq$  part, we consider a homomorphism  $f$  that is not by tISO. First, with  $f$ , we can find a matching  $G_f$  in  $G$ . By  $G_f$ , we can find a subpattern  $p' \in \text{Sub}(p) \setminus \text{Cov}(p, T)$ . Second, given  $p'$ , we can find an injective homomorphism (or subgraph isomorphism)  $f'$  with which we find a matching  $G_{f'}$  that is isomorphic to  $G_f$ . Hence, its count is included on the right by subgraph isomorphism regarding  $p'$ . To prove the  $\geq$  part, we consider an subgraph isomorphism  $f'$  of a subpattern  $p' \in \text{Sub}(p) \setminus \text{Cov}(p, T)$ , and we can find a homomorphism  $f$  that is not by tISO for the pattern  $p$ . Hence its count is included on the left. The proposition is proved by the two inequalities.  $\square$

By combing Eq. (1) and Eq. (8), we have

$$|\text{tISO}_{p,o,T}(v)| = \sum_{p' \in \{p\} \cup \text{Cov}(p,T)} |\text{ISO}_{p',o}(p')(v)| \quad (9)$$

Based on Eq. (9), we can compute  $|\text{ISO}_{p,o}(v)|$  by tISO using Algorithm 1 with which we have the following formula.

$$|\text{ISO}_{p,o}(v)| = \sum_{p_i \in \mathcal{P}} \mu_i \cdot |\text{tISO}_{p_i,o(p_i),T_i}(v)| \quad (10)$$

Here,  $\mathcal{P}$  is the set of all distinct patterns in  $\Lambda$ ,  $\mu_i = \sum_{x=p_i} (-1)^{d_x+1}$ ,  $d_x$  is the depth of  $x$  in  $\Lambda$ , and the depth of the root is 1.

The workflow of our approach SCOPE is presented in the lower part in Fig. 3. ① We first get the tree decomposition for a given pattern  $p$ . ② We show that ISO count of  $p$  can be obtained by subtracting the tISO count with the ISO count of  $\text{Cov}(p, T)$ , a subset of  $\text{Sub}(p)$  (Eq. (9)). ③ With Algorithm. 1, we recursively apply Eq. (9) to compute ISO counts by tISO counts (Eq. (10)). ④ We process tree decompositions by symmetry-breaking rules (Section 5) and our

---

### Algorithm 1: tISOToISO ( $p$ )

---

**Input:** pattern graph  $p$   
**Output:** A formula to compute  $|\text{ISO}_{p,o}(v)|$  by tISO

- 1  $\Lambda \leftarrow p$ , push  $p$  to  $Q$ ;
- 2 **while**  $Q \neq \emptyset$  **do**
- 3     pop  $Q$  to  $p'$ ; compute  $T_{p'}$  and  $\text{Cov}(p', T_{p'})$ ;
- 4     **foreach**  $p'' \in \text{Cov}(p', T_{p'})$  **do**
- 5         add a node  $p''$  and an edge  $(p', p'')$  into  $\Lambda$ ; push  $p''$  to  $Q$ ;
- 6     **return**  $\text{GenEq}(\Lambda.root)$ ;
- 7 **Procedure**  $\text{GenEq}(x)$
- 8     **return**  $\text{tISO}_{x,o(x),T_x} - \sum_{y \in x.child} \text{GenEq}(y)$ ;

---

new multi-join algorithm Scope (Section 6). Notably, symmetry-breaking rules are exclusive to ISO with which we define tISO. In particular, we use ISO for nodes in  $T$ .

## 5 TISO-BASED COUNTING

As shown in Eq. (4)-Eq. (7), a main cost in tISO-based counting is to compute ISO-matches ( $\mathbf{R}_\tau$ ) for every tree node  $\tau$  in a tree decomposition  $T$ . And the key issue is how to compute ISO-matches  $\mathbf{R}_\tau$  using automorphism orbits. A common technique to compute ISO for a given pattern graph  $p$  is by symmetry-breaking, which is used in subgraph enumeration to reduce the number of iso-matches of  $p$  in  $G$  if there exist automorphism orbits in  $p$  [33, 57, 77, 96]. This is due to the fact that a subgraph in  $G$  can be iso-matched multiple times given automorphism orbits. We introduce symmetry-breaking rules below.

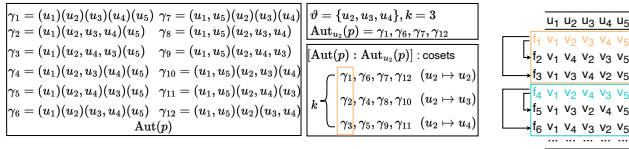
**Symmetry-breaking rules:** For a pattern graph  $p$  with automorphisms  $\text{Aut}(p)$ , a symmetry-breaking is an automorphism orbit of  $\text{Aut}(p)$ , namely,  $\vartheta = \{u_1, u_2, \dots, u_k\}$ . For each pair  $(u_i, u_j)$  in  $\vartheta$ , for  $2 \leq i < j \leq k$ , a partial order  $<$  is imposed such that  $u_1 < u_i$ . A symmetry-breaking rule ( $\text{SymR}$ ) for a given  $\vartheta$  is presented in the form of  $\theta = \{u_1 < u_2, \dots, u_1 < u_k\}$ . Assume that there is a total order ( $<$ ) on nodes in the data graph  $G$ . By a symmetry-breaking rule, it enforces  $f(u_1) < f(u_i)$  in the data graph  $G$ , if  $u_1 < u_i$  in  $\theta$ . In other words, An iso-match of  $p$  in  $G$  is valid if  $f[u_1] < f[u_i]$  for all  $i \in [2, k]$  in a  $\text{SymR } \theta$ .

Existing works enumerate subgraphs using a set of  $\text{SymRs}$  [6, 26, 33, 35, 45, 57–59, 77, 79, 84, 96, 100] for a pattern graph without tree decomposition. In this work, we study how to compute ISO for every tree node  $\tau$  using automorphism orbits and its  $\text{SymRs}$ , given a tree decomposition. It is important to note that this technique cannot be used for homomorphisms.

### 5.1 More about Automorphism Orbits

Given a pattern graph  $p$ , consider an automorphism orbit  $\vartheta = \{u_1, u_2, \dots, u_k\}$  in  $\text{Aut}(p)$  which can be represented by  $\text{SymR } \theta = \{u_1 < u_2, \dots, u_1 < u_k\}$ . We define a stabilizer subgroup of  $\text{Aut}_\vartheta(p) = \{\gamma \mid \gamma(u_1) = u_1 \text{ for } \gamma \in \text{Aut}(p)\}$  [24]. Note that  $\text{Aut}_\vartheta(p)$  is defined by fixing  $u_1$  in  $\vartheta$ .

With the stabilizer, the (left) cosets of  $\text{Aut}_\vartheta(p)$  in  $\text{Aut}(p)$ , denoted as  $[\text{Aut}(p) : \text{Aut}_\vartheta(p)]$ , are disjoint and are in the same size obtained by composing each automorphism of  $\text{Aut}_\vartheta(p)$  by a  $\gamma$  in  $\text{Aut}(p)$  such as  $[\text{Aut}(p) : \text{Aut}_\vartheta(p)] = \cup_{\gamma \in \text{Aut}(p)} \{\gamma \circ \gamma' \mid \gamma' \in \text{Aut}_\vartheta(p)\}$ . By the orbit-stabilizer theorem in group theory, the number of cosets is  $|\vartheta|$



(a) automorphisms, stabilizer subgroup and cosets (b) two msets

Figure 7: automorphisms and the  $\theta$  mapping set

and  $|\text{Aut}_g(p)| = |\text{Aut}(p)|/k$  for  $k = |\vartheta|$ . That is, there are  $k$  disjoint cosets in the same size, and in the  $i$ -th coset the automorphisms send  $u_1$  to the same  $u_i$  in  $\vartheta$ .

**Example 5.1:** Consider the pattern graph  $p$  in Fig. 2(a). Its  $\text{Aut}(p)$  is in Fig. 7(a) where  $|\text{Aut}(p)| = 12$ . There is an automorphism orbit  $\vartheta = \{u_2, u_3, u_4\}$  in  $p$ , which can serve the role of stabilizer as  $\text{Aut}_g(p) = \{\gamma_1, \gamma_6, \gamma_7, \gamma_{12}\}$ . Here,  $|\text{Aut}_g(p)| = |\text{Aut}(p)|/|\vartheta| = 12/3 = 4$ . We explain how a coset in  $[\text{Aut}(p) : \text{Aut}_g(p)]$  is constructed. Consider  $\gamma_6 \in \text{Aut}(p)$ . By composing  $\gamma_6$  with each of  $\text{Aut}_g(p)$ , we have the first coset,  $\{\gamma_1, \gamma_6, \gamma_7, \gamma_{12}\}$ . In detail,  $\gamma_6 \circ \gamma_1 = \gamma_6$ ,  $\gamma_6 \circ \gamma_6 = \gamma_1$ ,  $\gamma_6 \circ \gamma_7 = \gamma_{12}$ ,  $\gamma_6 \circ \gamma_{12} = \gamma_7$ . Hence  $\gamma_6$  is one that produces the first coset. Note that in the 1st coset all automorphisms send  $u_2$  to  $u_2$ . The 2nd coset sends  $u_2$  to  $u_3$ , and the 3rd coset sends  $u_2$  to  $u_4$  for  $\vartheta = \{u_2, u_3, u_4\}$ .

The stabilizer and the corresponding cosets are discussed regarding a pattern graph  $p$ . Given the  $k$  cosets by the stabilizer  $\text{Aut}_g(p)$  where  $|\vartheta| = k$ , we define a set of iso-mappings from  $p$  to  $G$  called an  $\vartheta$  mapping set ( $\vartheta$ -mapset).

**Definition 5.1:** ( $\vartheta$ -mapset) For a stabilizer  $\text{Aut}_g(p)$  over  $\vartheta = \{u_1, u_2, \dots, u_k\}$ , there are  $k$  cosets. Let  $\tilde{\gamma}_i$  be the automorphism selected from the  $i$ -th cosets for  $1 \leq i \leq k$ . An  $\vartheta$  mapping set ( $\vartheta$ -mapset) is a set of  $k$  iso-mappings  $\Phi = \{\phi_1, \phi_2, \dots, \phi_k\}$  where  $\phi_i$  is an iso-mapping by  $\tilde{\gamma}_i$  over the same set of nodes in  $G$ . For the  $\text{SymR } \theta = \{u_1 < u_2, \dots, u_1 < u_k\}$  over  $\vartheta$ ,  $\phi_1 \in \Phi$  is called the representative of  $\vartheta$ -mapset ( $\Phi$ ), if it satisfies  $\text{SymR } \theta$ .

**Reconstruction of  $\vartheta$  mapping set:** It is obvious that  $\phi_1$  is the only iso-mapping in  $\Phi$  that satisfies  $\text{SymR } \theta$ . Given  $\phi_1$  for  $\vartheta$ -mapset, we can reconstruct iso-mappings  $\phi_i$  as  $\phi_i(u_j) = \phi_1(\tilde{\gamma}_i(u_j))$  for  $u_j \in V_p$ . It is important to note that we do not need to find other iso-mappings if we find  $\phi_1$ , we can reconstruct the other mappings.

**Example 5.2:** Reconsider Example 5.1. In Fig. 7(a), there are 3 cosets for  $\vartheta = \{u_2, u_3, u_4\}$  in Fig. 7(a), and we have  $\tilde{\gamma}_1 = \gamma_1$ ,  $\tilde{\gamma}_2 = \gamma_2$ , and  $\tilde{\gamma}_3 = \gamma_3$ . Over  $\vartheta$ , the  $\text{SymR } \theta = \{u_2 < u_3, u_2 < u_4\}$ . Given the data graph  $G$  in Fig. 2(b), we show two  $\vartheta$ -mapsets in Fig. 7(b) where one is formed by the first 3 iso-mappings, and one is formed by the second 3 iso-mappings. Consider the 2nd  $\vartheta$ -mapset in which  $f_4$  ( $\phi_1$ ) is the representative that satisfies the  $\text{SymR}$ , and we can construct  $f_5$  ( $\phi_2$ ) and  $f_6$  ( $\phi_3$ ) by composing  $f_4$  with  $\tilde{\gamma}_2$  and  $\tilde{\gamma}_3$ , respectively.

**Generating multiple  $\text{SymR}$ s:** There are multiple sets of  $\text{SymR}$ s for a pattern graph  $p$ . An algorithm in [33] generates one set of  $\text{SymR}$ s for  $p$  randomly. Existing works show the efficiency of enumerating subgraphs using such a set of  $\text{SymR}$ s [3, 6, 26, 33, 35, 45, 57–59, 77, 79, 84, 96, 100]. Different from the existing work, we explore how to utilize  $\text{SymR}$ s for every tree node  $\tau$  in a tree decomposition  $T$ , and we need to select  $\text{SymR}$ s for  $\tau$ 's and for the entire  $T$ . There exists

## Algorithm 2: GenAllRules ( $p$ )

**Input:** a pattern graph  $p$   
**Output:** all sets of  $\text{SymR}$ s for  $p$

- 1  $A \leftarrow \text{Aut}(p); \mathcal{R} \leftarrow \emptyset;$
- 2  $\text{GenRules}(p, A, \emptyset, \emptyset);$
- 3 **return**  $\mathcal{R};$

**Procedure**  $\text{GenRules}(p, A, F, R)$

- 4 **if**  $|A| = 1$  **then**  $\mathcal{R} \leftarrow \mathcal{R} \cup \{R\};$
- 5 **else**
- 6  $\Theta \leftarrow$  the set of equivalent classes in  $A$  constrained by  $F;$
- 7 **foreach**  $\vartheta \in \Theta$  **do**
- 8 let  $u_i$  be the one with the smallest id in  $\vartheta;$
- 9  $F' \leftarrow F \cup \{u_i\};$
- 10 let  $\theta$  be the set of  $\text{SymR}$ s by  $\vartheta;$
- 11  $R' \leftarrow R \cup \{\theta\};$
- 12  $A' \leftarrow A$  constrained by  $F';$
- 13  $\text{GenRules}(p, A', F', R');$

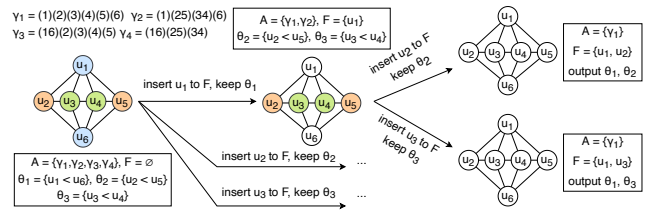


Figure 8: Generating  $\text{SymR}$ s

some  $\text{SymR}$  that cannot be efficiently used with  $T$ , which we will discuss later. Such an issue does not occur when applying  $\text{SymR}$ s to  $p$  without  $T$ .

We give an algorithm GenAllRules (Algorithm 2) to generate all sets of  $\text{SymR}$ s based on the algorithm given in [33]. We explain GenAllRules using the pattern graph  $p$  in Fig. 4(a) as an example, where we consider  $p$  as  $\tau$  in  $T$ . There are 4 automorphisms such that  $\text{Aut}(p) = \{\gamma_1, \gamma_2, \gamma_3, \gamma_4\}$ , where  $\gamma_1 = (u_1)(u_2)(u_3)(u_4)(u_5)(u_6)$ ,  $\gamma_2 = (u_1)(u_2u_5)(u_3u_4)(u_6)$ ,  $\gamma_3 = (u_1u_6)(u_2)(u_3)(u_4)(u_5)$ , and  $\gamma_4 = (u_1u_6)(u_2u_5)(u_3u_4)$ . With  $\text{Aut}(p)$ , there are 3 symmetry-breakings,  $\vartheta_1 = \{u_1, u_6\}$ ,  $\vartheta_2 = \{u_2, u_5\}$ , and  $\vartheta_3 = \{u_3, u_4\}$ . In GenAllRules,  $\mathcal{R}$  is the set of  $\text{SymR}$ s to be generated, which is initialized to be empty. It initially calls the procedure GenRules with the inputs of the pattern graph  $p$  and  $A = \text{Aut}(p)$ . In the procedure,  $F$  is the constraints when selecting  $\text{SymR}$ s based on the automorphisms  $A$ , and  $R$  is one set of  $\text{SymR}$ s to be generated. Initially, there are no constraints, so  $F = \emptyset$ ,  $\Theta = \{\vartheta_1, \vartheta_2, \vartheta_3\}$  (line 7). Suppose that  $\vartheta_1 = \{u_1, u_6\}$  is selected (line 8),  $u_1$  is added into  $F'$  as a constraint, with which it constrains that only an automorphism,  $\gamma_i$ , with  $(u_1)$  can be further explored next. In other words,  $u_1$  is fixed. Here, the  $\text{SymR } \theta = \{u_1 < u_6\}$  is added into  $R'$  (lines 11-12), the automorphism in  $A$  constrained by  $F'$  is  $A' = \{\gamma_1, \gamma_2\}$  in which  $(u_1)$  appears as constrained. It recursively calls GenRules. The procedure is illustrated in Fig. 8. The output of GenAllRules is  $\mathcal{R} = \{\mathcal{R}_1, \mathcal{R}_2\}$ , where  $\mathcal{R}_1 = \{\theta_1, \theta_2\}$  and  $\mathcal{R}_2 = \{\theta_1, \theta_3\}$ . Below, we use  $\mathcal{R}_\tau$  to denote the set of sets of  $\text{SymR}$ s for a tree node  $\tau$ .

**The  $\vartheta$  independency:** Let  $\vartheta$  and  $\vartheta'$  be any two automorphism orbits in the same set of  $\text{SymR}$ s (e.g.,  $\mathcal{R}_i$ ) generated by GenAllRules

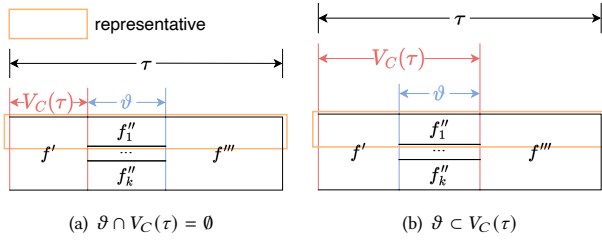


Figure 9: SymR with  $\tau$  and  $T$

$\mathbf{R}'_{\tau_1}$	$\mathbf{X}'_{\tau_1}$	$\mathbf{R}'_{\tau_2}$	$\mathbf{J}'_{\tau_2} = \mathbf{R}'_{\tau_2} \bowtie \mathbf{X}'_{\tau_1}$	$\mathbf{X}'_{\tau_2}$
u1 u4 u5 u6 v8 v7 v10 v9	u1 u4 u6 C v8 v7 v9 1	u1 u3 u4 u6 v8 v7 v10 v9	u1 u3 u4 u6 C v8 v7 v10 v9 2	u1 u3 u6 C v8 v7 v9 2
v8 v10 v7 v9 v8 v10 v12 v9	v8 v10 v9 2 v8 v12 v9 1	v8 v10 v7 v9 v8 v10 v12 v9	v8 v10 v7 v9 1 v8 v10 v12 v9 1	v8 v10 v9 2 v8 v12 v9 2
v8 v12 v10 v9 ...	... ...	v8 v12 v10 v9 ...	v8 v12 v10 v9 2 ...	... ...

(a)  $\{u_1 < u_6\}, \tau_1$

(b)  $\{u_1 < u_6\}, \tau_2$

Figure 10: Apply SymRs to  $\tau$ 's given  $p$  and  $T$  in Fig. 4(a)

(Algorithm 2). There are only two cases between  $\vartheta$  and  $\vartheta'$ : one is contained in another, and the other is the two are disjoint.  $\vartheta$  and  $\vartheta'$  are independent since  $\vartheta'$  is an automorphism orbit when we fix a node in  $\vartheta$  or vice-versa. The corresponding  $\tilde{\gamma}$  of the  $\vartheta$ -mapsets are independent. Therefore, any subset of a set of SymRs can be used. We call this  $\vartheta$  independency.

## 5.2 Automorphism Orbits in Tree Nodes

Following the discussion on automorphisms for a pattern graph, with a tree decomposition  $T$ , we compute iso-matches ( $\mathbf{R}_\tau$ ) for each tree node  $\tau \in T$  using SymRs by considering  $\tau$  as a pattern graph. Here, a key issue is the relationship between SymR  $\theta$  for  $\tau$  and  $V_C(\tau)$ . Note that  $V_C(\tau)$  is used as the group-by attributes in Eq. (4) and Eq. (6) to compute tISO counts. We need to ensure that the application of SymRs for a tree node  $\tau$  does not affect the aggregation by the group-by attributes  $V_C(\tau)$ .

There are 4 cases for an automorphism  $\vartheta$  (or its SymR  $\theta$ ) regarding a tree node  $\tau$ , namely, ①  $\vartheta \not\subset \tau$ , ②  $\vartheta \cap V_C(\tau) = \emptyset$ , ③  $\vartheta \subset V_C(\tau)$ , and ④  $\vartheta \cap V_C(\tau) \neq \emptyset$  and  $\vartheta \not\subset V_C(\tau)$ . The last 3 cases are the cases when  $\vartheta \subset V(\tau)$ . Assume there is an iso-match  $f$  of the entire  $\tau$  such that  $f = f' \| f'' \| f'''$ , where  $f'$  is for the part of mapping by  $V_C \setminus \vartheta$ ,  $f''$  is the part of mapping by  $\vartheta$ , and  $f'''$  is for the part of mapping by  $V_\tau \setminus (V_C \cup \vartheta)$ . By SymR  $\theta$  for  $\vartheta$ , it affects  $f''$ , or more precisely, its  $\vartheta$ -mapset of size  $k = |\vartheta|$ . For ①, as  $\vartheta \not\subset \tau$ ,  $\vartheta$  has no impacts on  $V_C$ . For ②, it only finds the representative  $f''_1$  using SymR  $\theta$ . As  $\vartheta$  and  $V_C(\tau)$  is disjoint, its count without SymR  $\theta$  is its count with SymR  $\theta$  multiplied by  $k$ . For ③, it only finds the representative  $f''_1$  using SymR  $\theta$  for  $\vartheta$ . As  $\vartheta$  is included in  $V_C(\tau)$ , it needs to reconstruct the other mappings,  $f''_i$ , in the  $\vartheta$ -mapset. For ④, it is prohibited for the reason that there may exist two different  $\vartheta$ -mapsets whose match to  $V_C$  overlap but are not the same. The count by such  $V_C$  is incorrect. We show ② and ③ in Fig. 9.

**Example 5.3:** Consider Example 5.2 where we take this pattern graph  $p$  as  $\tau$  in a tree decomposition  $T$ . As shown in Fig. 7(a), there are 3 cosets for  $\vartheta = \{u_2, u_3, u_4\}$  in Fig. 7(a), where  $k = |\vartheta| = 3$ . Over  $\vartheta$ , SymR  $\theta = \{u_2 < u_3, u_2 < u_4\}$ . For ②, suppose  $V_C(\tau) = \{u_1\}$ , and we only get  $f_4$  for the 2nd  $\theta$  mapping set. There should be additional

$k - 1$  iso-mappings. The count for  $V_C = \{u_1\}$  needs to be multiplied by  $k$  for a  $\vartheta$ -mapset. For ③, suppose  $V_C(\tau) = \{u_1, u_2, u_3, u_4\}$ , and we only get  $f_4$  for the 2nd  $\vartheta$ -mapset. There should be additional  $k - 1$  iso-mappings to be reconstructed.

We discuss how to do count-correction/reconstruction. Here, we discuss it by assuming that  $\mathbf{R}_\tau$  is for a single  $\vartheta$ -mapset, which is a set of  $k$  iso-mappings  $\Phi = \{\phi_1, \phi_2, \dots, \phi_k\}$  where  $\phi_i$  is an iso-mapping by  $\tilde{\gamma}_i$  over the same set of nodes in  $G$  (Definition 5.1). Among all in  $\Phi$ ,  $\phi_1$  is the one that satisfies the SymR  $\theta$ .

First, for ②, we only need to do count-correction, as illustrated in Fig. 9(a). We give the details for a non-leaf node  $\tau$  that is not the root (Eq. (5) and Eq. (6)) regarding a given  $\vartheta$  where  $\vartheta \cap V_C(\tau) = \emptyset$ . The others can be dealt with in a similar manner. We show how we correct it for a single  $\vartheta$ -mapset. If we can do it for a single  $\vartheta$ -mapset, the overall sum in Eq. (6) is correct as it is the sum of the counts for all  $\vartheta$ -mapsets by the group-by attributes  $V_C(\tau)$ . It is worth noting that the only place that changes is  $\mathbf{R}_\tau$  in Eq. (5) if we enforce SymR  $\theta$  w.r.t  $\vartheta$ , where every child  $\mathbf{X}_{\tau_i}$  remains unchanged. Without the SymR  $\theta$ , the size of  $\mathbf{R}_\tau$  is  $k = |\vartheta|$ , and the size of the corresponding  $\mathbf{J}_\tau$  is  $k$  because it can only join one tuple from  $\mathbf{X}_{\tau_i}$ . By enforcing SymR  $\theta$ , the size of  $\mathbf{R}_\tau$  becomes 1 for the representative  $\phi_1$ , and its count  $C$  in  $\mathbf{X}_\tau(V_C(\tau), C)$  in Eq. (6) becomes  $C/k$ . We correct its count by multiplying it by  $k$ . There is no need to do reconstruction.

Second, for ③, there is no need to do count-correction. We explain it using Fig. 9(b). The count for  $f$ , where  $f''_1$  is the representative of a  $\vartheta$ -mapset, is reduced by  $1/k$  for  $k = |\vartheta|$ . That is the correct count of  $f$  with  $f''_1$ . We only need to reconstruct the other mappings following the discussion of reconstruction given in Section 5.1. That is, we reconstruct the other  $f''_i$  in  $\Phi$  by  $f''_1$ .

**Example 5.4:** Consider Example 4.1 where its pattern graph  $p$ , tree decomposition  $T$ , and data graph  $G$  are in Fig. 4. There are 3 tree nodes,  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$ , and the root is  $\tau_3$ . For the leaf node  $\tau_1$  with  $V_C(\tau_1) = \{u_1, u_4, u_6\}$ , there are two SymRs  $\theta_{11} = \{u_1 < u_6\}$  over  $\vartheta_{11} = \{u_1, u_6\}$ , and  $\theta_{12} = \{u_4 < u_5\}$  over  $\vartheta_{12} = \{u_4, u_5\}$ . Each has 2 cosets. For  $\theta_{11}$ ,  $\tilde{\gamma}_1(\theta_{11}) = (u_1)(u_4)(u_5)(u_6)$ , and  $\tilde{\gamma}_2(\theta_{11}) = (u_1, u_6)(u_4)(u_5)$ . For  $\theta_{12}$ ,  $\tilde{\gamma}_1(\theta_{12}) = (u_1)(u_4)(u_5)(u_6)$ , and  $\tilde{\gamma}_2(\theta_{12}) = (u_1)(u_6)(u_4, u_5)$ . We explain  $\vartheta_{11}$  which is the case ③, because  $\vartheta_{11} \subset V_C(\tau_1)$ , for the leaf node  $\tau_1$ . For a leaf-node  $\tau_1$ ,  $\mathbf{J}_{\tau_1} = \mathbf{R}_{\tau_1}$ . First,  $\mathbf{R}_{\tau_1}$  with  $\mathbf{X}_{\tau_1}$  are shown in Fig. 5 without SymRs. Second, consider  $\mathbf{R}_{\tau_1}$  with SymR  $\theta_{11} = \{u_1 < u_6\}$ . One of its  $\vartheta$ -mapset is  $\{f_1, f_5\}$ , where for example  $f_1 = \{u_1 \rightarrow v_8, u_4 \rightarrow v_7, u_5 \rightarrow v_{10}, u_6 \rightarrow v_9\}$ , and the other  $\vartheta$ -mapsets are  $\{f_2, f_6\}$ ,  $\{f_3, f_7\}$ , and  $\{f_4, f_8\}$ . Third, we have  $\mathbf{R}'_{\tau_1}$  by enforcing SymR  $\theta_{11}$  as shown in Fig. 10(a), where it only keeps the representatives of  $\vartheta$ -mapsets,  $\{f_1, f_2, f_3, f_4\}$ . The corresponding  $\mathbf{X}'_{\tau_1}$  is also shown in Fig. 10(a). Comparing  $\mathbf{R}_{\tau_1}$  in Fig. 5, the number of iso-mappings in  $\mathbf{R}'_{\tau_1}$  in Fig. 10(a) is reduced by a half as  $|\vartheta_{11}| = 2$ . And  $\mathbf{X}'_{\tau_1}$  is also reduced by a half comparing  $\mathbf{X}_{\tau_1}$ . We can reconstruct  $\mathbf{X}_{\tau_1}$  from  $\mathbf{X}'_{\tau_1}$  using the representatives  $\{f_1, f_2, f_3, f_4\}$ . For example, with  $f_1$ , we can get  $f_5$  (e.g.,  $(v_9, v_7, v_8)$  from  $(v_8, v_7, v_9)$ ) by swapping the node mapped by  $u_1$  and  $u_6$ , and keep the count as the one in  $f_1$ .  $\square$

**The  $\vartheta$  independency in a tree node  $\tau$ :** We have discussed the  $\vartheta$  independency for a pattern graph in Section 5.1. We also need to impose such dependency for a tree node  $\tau$  in  $T$ . We need certain conditions regarding  $V_C(\tau)$  for  $\tau$  and  $V_C(\tau_i)$  for each child  $\tau_i$  of

$\tau$  in  $T$ . First, for  $V_C(\tau)$ , the condition is that for any single  $\vartheta$  in  $\tau$ , the automorphism  $\tilde{\gamma}_l$  selected from the  $l$ -th coset must satisfy  $\tilde{\gamma}_l(u_j) = u_j$  for any  $u_j \in V_C(\tau) \setminus \vartheta$ . The main idea is to ensure that group-by attributes not in  $\vartheta$  are fixed when computing the aggregation so that the aggregation result can be correctly used in its parent. For example, in Fig. 7(a),  $\vartheta = \{u_2, u_3, u_4\}$  for a  $\tau$  over the set of nodes  $\{u_1, u_2, u_3, u_4, u_5\}$ . Suppose it is a leaf node  $\tau$  in  $T$ , and  $V_C(\tau) = \{u_2, u_3, u_4, u_5\}$ , we have  $V_C(\tau) \setminus \vartheta = \{u_5\}$ . Consider the 2nd coset,  $\gamma_2$  can be  $\tilde{\gamma}_2$  as  $\gamma_2(u_5) = u_5$ , and  $\gamma_8$  cannot be selected as  $\tilde{\gamma}_2$ . Second, for  $V_C(\tau_i)$ , we ensure that  $\tilde{\gamma}_l(u_j) = u_j$  for any  $u_j \in V_C(\tau_i) \setminus \vartheta$ , and that  $\vartheta \cap V_C(\tau_i) = \emptyset$  or  $\theta \in \mathcal{R}_{\tau_i}$ . Note that  $\theta \in \mathcal{R}_{\tau_i}$  implies  $\vartheta \subset V_C(\tau_i)$ . These conditions ensure that in Eq. (5), mappings in the same  $\vartheta$ -mapset of  $\tau$  have the same  $C_1, \dots, C_w$ , therefore establishing the correctness of computing aggregation in  $\tau$ . If such a  $\tilde{\gamma}_i$  does not exist, we do not use this  $\vartheta$ . If two orbits  $\vartheta$  and  $\vartheta'$  both satisfy these conditions, they are independent and both can be used.

**The automorphism orbit at the root of  $T$ :** We discuss how we use *SymRs* at the root node  $\tau$  in  $T$ , where its  $V_C(\tau) = \emptyset$  as it does not have any parent. The group-by attribute in the root is  $o$ , a single node, which is the orbit in a local subgraph counting query  $Q = (p, o)$ . There are two cases. One is  $o \notin \vartheta$ , and one is  $o \in \vartheta$ . For the former, it is a similar case that  $\vartheta$  and  $V_C(\tau)$  are disjoint to be dealt with. For the latter, for any *SymR*  $\theta = (o < u_i)$ , we only need to reconstruct by swapping  $o$  with  $u_i$  w.r.t the mapping by *SymR*  $\theta$ .

### 5.3 The Optimizations

We have discussed how to correct the count and reconstruct a single  $\vartheta$ -mapset for a tree node  $\tau_i$  in  $T$ . We can ensure all the counts are correct if we do so for every  $\tau_i$  in  $T$ . We propose an optimization technique in a way that we do not need to do so for every  $\tau_i$  in  $T$ , and we can delay it from  $\tau_i$  to its parent  $\tau$  under certain conditions.

First, consider the case  $\textcircled{2}$  for  $\tau_i$  where  $\vartheta \cap V_C(\tau_i) = \emptyset$ , its parent  $\tau = \text{parent}(\tau_i)$  must not have  $\vartheta$  by tree decomposition. The count-correction can be delayed from  $\tau_i$  to  $\tau$ . We explain it below. Suppose we have to do count-correction by multiplying  $k$  for a  $\vartheta$ -mapset in  $\tau_i$  when a *SymR*  $\theta$  is used. That is we have to correct  $X_{\tau_i}(V_C(\tau_i), C)$  with the  $\vartheta$ -mapset to be  $X_{\tau_i}(V_C(\tau_i), C \cdot k)$  by Eq. (6). Furthermore, consider  $\tau$  which is the parent of  $\tau_i$ . As shown in Eq. (6), its  $C_i$  before correction under  $\tau_i$  becomes  $C_i \cdot k$  after correction, for the parent  $\tau$ . We have  $V_C(\tau) \Upsilon_{\text{sum}}(C_1 \times \dots \times C_i \cdot k \times \dots) \rightarrow C = V_C(\tau) \Upsilon_{\text{sum}}(C_1 \times \dots \times C_i \times \dots) \cdot k \rightarrow C$  w.r.t Eq. (6), where  $C_i$  is the count before count-correction in  $\tau_i$ . In other words, it is possible that we do not do count-correction by multiplying  $k$  in  $\tau_i$  but do it in its parent  $\tau$ .

Second, for the case  $\textcircled{3}$  where  $\vartheta \subset V_C(\tau_i)$ , its parent  $\tau$  must have the same  $\vartheta$ . There are 3 sub-cases with  $\tau$ ,  $\vartheta$  is an automorphism orbit in  $\tau$  ( $\textcircled{2}$ ,  $\textcircled{3}$ ), and  $\vartheta$  is not an automorphism orbit in  $\tau$ . When  $\vartheta$  is an automorphism in  $\tau$  for the sub-cases  $\textcircled{2}/\textcircled{3}$ , the counts by  $V_C(\tau)$  will not be affected by  $\vartheta$ -mapsets in  $\tau_i$ , because, for any  $\vartheta$ -mapset in  $\tau_i$ , there will be one and only one  $\vartheta$ -mapset in  $\tau$  that map to the same nodes in  $G$  by the same  $\vartheta$ . This is similar to the discussion in Section 5.2. We can delay reconstruction. Third, for the case that  $\vartheta$  is not an automorphism orbit in  $\tau$ , we have to do reconstruction for  $\tau_i$ , and we cannot delay it.

**Example 5.5:** Continue Example 5.4. As given in Example 5.4, for the leaf node  $\tau_1$  with  $V_C(\tau_1) = \{u_1, u_4, u_6\}$ , there are two *SymRs*

**Table 1: Some statistics about *SymRs***

$k$ -node	# of $p$	Symmetry( $p$ )	All-in- $T$	$\geq 1$ -in- $T$	%
5	58	58	24	36	51.1%
6	407	359	144	238	53.1%
7	4,306	3,298	1,405	2,171	54.3%

$\theta_{11} = \{u_1 < u_6\}$  over  $\vartheta_{11} = \{u_1, u_6\}$ , and  $\theta_{12} = \{u_4 < u_5\}$  over  $\vartheta_{12} = \{u_4, u_5\}$ . For  $\tau_2$ , the *SymRs* are  $\theta_{21} = \{u_1 < u_6\}$  and  $\theta_{22} = \{u_3 < u_4\}$ . Here,  $V_C(\tau_2) = \{u_1, u_3, u_6\}$ . For  $\tau_3$  (the root), the *SymR* is  $\theta_{31} = \{u_1 < u_6\}$ , and  $o = u_1 \in \tau_3$  where  $o$  is the orbit of the given query  $Q = (p, o)$ . Note that  $\theta_{11}$  in  $\tau_1$ ,  $\theta_{21}$  in  $\tau_2$ , and  $\theta_{31}$  are identical.

The results by count-correction/reconstruction for every tree node are shown in Fig. 5. In Fig. 10(a), we show the result in  $X'_{\tau_1}$  by enforcing the *SymR*  $\theta_{11}$  in  $\tau_1$ , together with its  $R'_{\tau_1}$ . Here, in  $\tau_1$ , it is the case  $\textcircled{3}$ , for  $\vartheta_{11} \subset V_C(\tau_1)$ , and in  $\tau$ , it is the case of  $\textcircled{3}$  as well, such that  $\vartheta_{21} \subset V_C(\tau_2)$ . Note  $\vartheta_{21} = \vartheta_{11}$ . We can delay reconstruction from  $\tau_1$  to  $\tau_2$ . The result is presented in Fig. 10(b).

We can also delay reconstruction to  $\tau_3$ . Note that  $\tau_3$  is the root of  $T$ , which does not have  $V_C(\tau_3)$ , and we can treat it in a specific way as discussed. That is, for the *SymR*  $\theta_{31}$ , we reconstruct  $\phi_2(u_1)$  from  $\phi_1(u_6)$  for each  $\vartheta$ -mapset.  $\square$

**Selecting *SymRs*:** We select *SymRs* for  $T$  as follows. First, we obtain  $\mathcal{R}_\tau$  for each  $\tau$  in  $T$  by GenAllRules. Second, we check each combination of all rule sets in all tree nodes. Let  $\mathcal{R}'_\tau \in \mathcal{R}_\tau$  be the selected rule set of  $\tau$ . We check the constraints for  $\vartheta$ -mapsets independency and remove invalid rules in  $\mathcal{R}'_\tau$ . Third, we use a simple cost function  $\sum_{\tau \in T, \theta \in \mathcal{R}'_\tau} |\theta| * |V_\tau|$  to select one set of rules that maximizes this function. Here, we prefer rules with larger  $|\theta|$  since they can reduce more matches for a given tree node. Additionally, we prefer rules applied to larger tree nodes because the induced subgraphs are more challenging to enumerate.

We show some statistics about *SymRs* in Table 1. The 1st column is  $k$  for  $k$ -node pattern graph  $p$ , the 2nd column is the total number of pattern graphs with such  $k$ , the 3rd column is the total number of patterns that have *SymRs*, the 4th column is the number of tree decompositions that can use all *SymRs* that appear in  $p$  in their nodes, the 5th column is the number of tree decompositions that can use at least one *SymR* that appears in  $p$  in their nodes, and the last column is the percentage of *SymRs* that appear in  $p$  to be used in  $T$  on average. A majority of pattern graphs are with *SymRs* with tree decomposition.

## 6 MULTI-JOIN ALGORITHMS

In this section, we discuss how to process a pattern graph  $p = (V_p, E_p)$  given its tree decomposition  $T$  based on the multi-join algorithm Leapfrog [91], which is a state-of-the-art worst-case optimal algorithm that is also used in [101]. In brief, Leapfrog is to process a join query over  $m$  relations based on the join attribute order using iterators. To process it for  $p$  by Leapfrog, we can represent an edge  $e_i \in E_p$  as a relation for  $|E_p| = m$ .

There are several ways to process the aggregations given  $T$ . First, following Eq. (4)-Eq. (7), we can process every tree node  $\tau \in T$  using Leapfrog, maintain its result in  $X_\tau$ , and join all such  $X_\tau$  relations. As pointed out in [101], this approach cannot be taken when the



sizes of such relations are too large to keep in the main memory. For example, in Fig. 4(b), there are 3 tree nodes, and for  $i = 1$  or  $i = 2$ , the relation  $X_{\tau_i}(V_C(\tau_i), C)$  is a 4 attribute relation where  $V_C(\tau_i)$  is the group-by attributes and  $C$  maintains its aggregation. In this example, the size of  $X_{\tau_i}(V_C(\tau_i), C)$  is  $O(n^k)$  for  $k = |V_C(\tau_i)|$  where  $n$  is the number of nodes in the data graph  $G$ . We call it an  $O(n^k)$ -approach, which is related to the space complexity.

In [101], DISC proposed an  $O(1)$ -approach regarding the memory, which we show in Algorithm 3. We call it Disc and explain it using an example decomposition  $T$ .

**Example 6.1:** Consider a tree decomposition  $T$  for a pattern graph  $p$  with 3 tree nodes,  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$ , where  $\tau_3$  is the root as shown at the top in Fig. 11. We assume that  $\tau_1$  is over two subgraphs  $B$  and  $D$  with  $V_C(\tau_1) = B$ ,  $\tau_2$  is over three subgraphs  $A$ ,  $E$ , and  $B$  with  $V_C(\tau_2) = A$ , and  $\tau_3$  is over two subgraphs  $F$  and  $A$ . Furthermore, we assume that  $V_C(\tau_1) = B$  can be divided into two disjoint sets,  $V_C^1(\tau_1) = B_1$  and  $V_C^2(\tau_1) = B_2$ , and  $V_C(\tau_2) = A$  can be divided into two disjoint sets,  $V_C^1(\tau_2) = A_1$  and  $V_C^2(\tau_2) = A_2$ .

For  $T$  in Example 6.1, Disc (Algorithm 3) processes it as follows starting from the root  $\tau_3$  in  $T$ . ❶ It finds a partial match  $f_{AB} = f_A \parallel f_B$  that matches  $V_C(\tau_2) = A$  in  $\tau_2$  and then  $V_C(\tau_1) = B$  in  $\tau_1$ . ❷ Given  $f_{AB}$ , it enumerates all  $D$  in  $\tau_1$ , and stores its count in  $X_{\tau_1}(B, C)$ . ❸ Given  $f_{AB}$ , it enumerates all  $E$  in  $\tau_2$ , and updates its count in  $X_{\tau_1}(A, C)$  regarding  $f_{AB}$  using the count done in ❷. ❹ For  $f_A$  in  $\tau_2$ , it repeats ❷ and ❸ to compute its final count in  $X_{\tau_2}(A, C)$  regarding  $f_A$ . ❺ It indicates that the final count is stored in  $X_{\tau_2}(A, C)$  for  $f_A$ . ❻ In  $\tau_3$ , with the count of  $X_{\tau_2}(A, C)$  for  $f_A$  done in ❺, it enumerates  $F$  and updates its count regarding the orbit of  $p$ . ❼ By repeating ❶-❻ for all possible matches of  $f_A$ , it gets the final count for the orbit  $o$ . Disc repeats it for every possible match of  $f_A$  by maintaining it with 2 counts. Here, the number of counts need to be maintained is  $|V_T| - 1$ , which is considered as a constant.

Disc is space-efficient. However, the cost of computing  $T$  is high, which is related to the number of iterations. The number of iterations depends on  $\alpha = |\cup_i V_C(\tau_i)|$  over a path in  $T$ . In Example 6.1, it is  $|V_C(\tau_2) \cup V_C(\tau_1)| = |A \cup B|$ , and the number of exploring  $\tau_1$  in the data graph  $G$  is  $O(n^\alpha)$ . We propose a new  $O(m)$ -approach regarding the memory, which we call Scope. The main idea is to reduce  $\alpha$  by only using a part of  $V_C(\tau_i)$ , while keeping the other part of  $V_C(\tau_i)$  in memory. The memory used is bounded by  $O(m)$ . We give the algorithm in Algorithm 4, and explain it using Example 6.1.

As given in Example 6.1, we have  $A = A_1, A_2$ , and  $B = B_1, B_2$ . Suppose  $|A_1 \cup B_1| < |A \cup B|$ . Here, we take  $A_2$  and  $B_2$  as an edge to ensure  $O(m)$  memory. Scope (Algorithm 4) processes it as follows starting from the root  $\tau_3$  in  $T$ . ❶ It finds a partial match  $f_{A_1 B_1} = f_{A_1} \parallel f_{B_1}$  that matches  $A_1 \subset V_C(\tau_2)$  in  $\tau_2$  and then  $B_1 \subset V_C(\tau_1)$  in  $\tau_1$ . ❷ Given  $f_{A_1 B_1}$ , it enumerates all  $B_2$  and  $D$  in  $\tau_1$ , and stores a count for every  $f_{B_2}$  that can expand from  $f_{A_1 B_1}$  in  $X_{\tau_1}(B_1 B_2, C)$ . ❸ Given  $f_{A_1 B_1}$ , it enumerates all  $A_2, B_2$ , and  $E$  in  $\tau_2$ , and updates the counts in  $X_{\tau_2}(A_1 A_2, C)$  regarding  $f_{A_1}$ . To update, it needs to find the count done in ❷ by hash-join. As all the matches  $f_{B_2}$  given  $f_{A_1 B_1}$  are distinct, the join cost is constant. ❹ For  $f_{A_1}$  in  $\tau_2$ , it repeats ❷ and ❸ to compute its final count in  $X_{\tau_2}(A_1 A_2, C)$  regarding  $f_{A_1}$ . ❺ It indicates that the final count is stored in  $X_{\tau_2}(A_1 A_2, C)$  regarding  $f_{A_1}$ . ❻ In  $\tau_3$ , with the count of  $X_{\tau_2}(A_1 A_2, C)$  regarding  $f_{A_1}$  done in ❺, it enumerates  $A_2$  and  $F$  and updates its count regarding the orbit

### Algorithm 3: Disc( $f_i, \tau, T$ )

**Input:** an  $i$ -mapping  $f_i$ , pattern graph  $\tau = (V_\tau, E_\tau)$ , a tree decomposition  $T$

**Output:**  $X_\tau$

- 1 Disc( $f_i, \tau_j, T$ ) if  $i$  is the smallest number for  $f_i$  to contain  $V_C(\tau_j)$  for every child  $\tau_j$  of  $\tau$  in  $T$ ;
- 2 **if**  $i = |V_\tau|$  **then**
- 3 | update  $X_\tau$  based on  $V_C(\tau)$ ;
- 4 **else**
- 5 | let  $u_{i+1}$  be the  $(i+1)$ -th node in  $V_\tau$  in order;
- 6 | find all  $f(u_{i+1})$  matches that can expand from  $f_i$  constrained by  $\tau$ , denoted as  $\text{val}(f_i \rightarrow u_{i+1})$ ;
- 7 | **for each**  $v$  in  $\text{val}(f_i \rightarrow u_{i+1})$  **do**
- 8 | | Disc( $f_i \parallel v, \tau, T$ );

### Algorithm 4: Scope( $f_i, \tau, T$ )

**Input:** an  $i$ -mapping  $f_i$ , pattern graph  $\tau = (V_\tau, E_\tau)$ , a tree decomposition  $T$

**Output:**  $X_\tau$

- 1 let  $V_C(\tau_j) = V_C^1(\tau_j) \cup V_C^2(\tau_j)$  where  $V_C^1(\tau_j) \cap V_C^2(\tau_j) = \emptyset$  for every child  $\tau_j$  of  $\tau$ ;
- 2 Scope( $f_i, \tau_j, T$ ) if  $i$  is the smallest number for  $f_i$  to contain  $V_C^1(\tau_j)$  for every child  $\tau_j$  of  $\tau$  in  $T$ ;
- 3 **if**  $i = |V_\tau|$  **then**
- 4 | update  $X_\tau$  based on  $V_C^2(\tau)$  using  $X_{\tau_j}$  for every child  $\tau_j$  of  $\tau$ ;
- 5 **else**
- 6 | let  $u_{i+1}$  be the  $(i+1)$ -th node in  $V_\tau$  in order;
- 7 | find all  $f(u_{i+1})$  matches that can expand from  $f_i$  constrained by  $\tau$ , denoted as  $\text{val}(f_i \rightarrow u_{i+1})$ ;
- 8 | **for each**  $v$  in  $\text{val}(f_i \rightarrow u_{i+1})$  **do**
- 9 | | Scope( $f_i \parallel v, \tau, T$ );

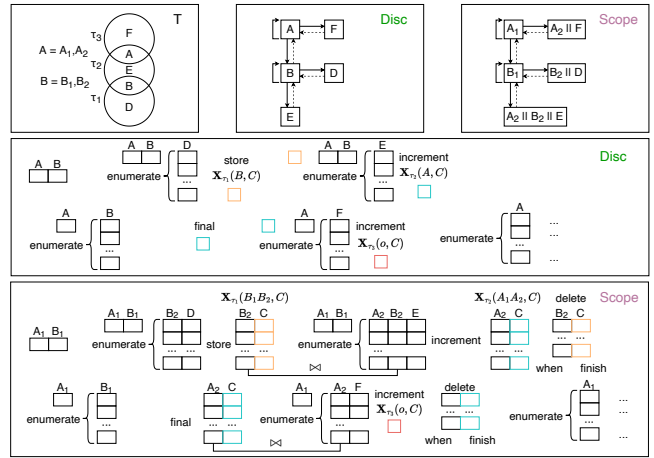


Figure 11: Disc vs Scope

of  $p$ . ❼ By repeating ❶-❻ for all possible matches of  $f_{A_1}$ , it gets the final count for the orbit  $o$ . Scope repeats it for every possible match by maintaining  $O(m)$  matches and counts supposing  $A_2$  and  $B_2$  are for an edge. As indicated in Fig. 11, when we start processing  $\tau_2$  in ❷, we can release the memory used for its child in ❸.

Scope is more efficient than Disc. For this example, we have  $|A_1 \cup B_1| < |A \cup B|$ . As an optimization technique, we share the  $V_C^1$  of a child with the  $V_C^1$  of its parent. In this example, if we make it as  $A_1 = B_1$ , we have  $|A_1 \cup B_1| = |A_1|$ , and we only have  $O(n^{|A_1|})$  iterations instead of  $O(n^{|A_1 \cup B_1|})$ . Recall that a similar optimization is to share *SymRs* between a child and its parent to reduce the number of reconstructions. Consider the pattern graph  $p$  and the tree decomposition  $T$  in Fig. 4. Here,  $A_1 = B_1 = \{u_1\}$ ,  $A_2 = \{u_3, u_6\}$ ,  $B_2 = \{u_4, u_6\}$ ,  $D = \{u_5\}$ ,  $E = \emptyset$ , and  $F = \{u_2\}$ . The number of iterations is  $O(n^{|A_1|}) = O(n)$ , and  $\tau_1$  can be processed efficiently. We put the complexity analysis of the three multi-join algorithms in the full version in our GitHub repository.

## 7 RELATED WORK

**Subgraph Counting.** The recent survey on subgraph counting [80] outlines three primary approaches to exact subgraph counting: enumeration-based, matrix-based, and decomposition-based. Enumeration-based approaches [33, 43, 44, 50, 52, 67, 71, 81, 97, 98] count by enumerating all matches of the pattern graph in the data graph. Matrix-based approaches [23, 41, 42, 63–65] rely on resolving linear algebra equations, which are grounded in the enumeration of other pattern graphs. JESSE [63–65] is a representative matrix-based approach that can automatically generate and select equations, but it is limited to computing the local counts of all  $k$ -node patterns collectively for a specified  $k$ . Decomposition-based approaches [4, 60, 62, 73, 74, 101] count  $p$  based on enumerating smaller graphs that are obtained by decomposing  $p$ . DISC [101], EVOKE [73] and SCOPE are in this category. For approximate subgraph counting, there are sampling-based approaches [11, 14, 15, 17, 18, 29, 30, 37, 46, 51, 56, 75, 78, 94, 95, 99] and learning-based approaches [93, 102]. Our approach is a new decomposition-based approach for general local subgraph counting.

**Subgraph Matching.** Subgraph matching enumerates iso-matches from the pattern graph to the data graph, with many works founded on Ullmann’s backtracking [90] or Leapfrog [91]. Research efforts have been devoted to filtering candidates [10, 12, 19, 33, 38, 39, 53, 54, 87, 88, 103], optimizing matching order [12, 38, 39, 53, 54, 66, 83, 85, 87, 88], and using previous matching results for pruning [7, 38, 47, 53, 54]. Different backtracking algorithms have been proposed [48, 85] to reduce set intersections. There are also distributed approaches that decompose the query into sub-structures and assembly the matches of sub-structures to obtain the pattern’s results [58, 77, 79, 82, 96, 100]. Experimental studies for subgraph matching can be found in [59, 86]. We extend symmetry-breaking [33] in subgraph matching to decomposition-based counting.

**Worst-case optimal join.** The *AGM* bound [8] gives the worst-case output size of a multi-join. In worst-case scenarios, executing a series of binary joins is inefficient since their complexity exceeds the *AGM* bound. Conversely, worst-case optimal join algorithms such as NPRR [69], GenericJoin [70], and Leapfrog [91] have a time complexity that matches this *AGM* bound. Owing to its remarkable efficiency, Leapfrog finds extensive application in both subgraph matching and subgraph counting. Various studies [3, 49, 66, 89] have amalgamated the worst-case optimal join with a binary join, steered by the principles of tree decomposition for better performance. Based on Leapfrog, Disc [101] handles aggregations in tree decompositions, and our Scope improves Disc.

Table 2: The 12 datasets

Graph	Notation	V	E	avg. degree
web-spam	WS	$4.8 \times 10^3$	$3.7 \times 10^4$	15.7
rec-movielens-user-movies-10m	RM	$7.6 \times 10^3$	$5.5 \times 10^4$	14.6
bio-grid-yeast	BY	$6.0 \times 10^3$	$1.6 \times 10^5$	52.2
ca-AstroPh	CA	$1.9 \times 10^4$	$2.0 \times 10^5$	21.1
rec-github	RG	$1.2 \times 10^5$	$4.4 \times 10^5$	7.2
soc-gowalla	SG	$2.0 \times 10^5$	$9.5 \times 10^5$	9.7
soc-youtube	SY	$1.1 \times 10^6$	$3.0 \times 10^6$	5.3
web-wiki-ch-internal	WW	$1.9 \times 10^6$	$9.0 \times 10^6$	9.3
web-hudong	WH	$2.0 \times 10^6$	$1.4 \times 10^7$	14.6
ca-coauthors-dblp	CC	$5.4 \times 10^5$	$1.5 \times 10^7$	56.4
soc-livejournal1	SL	$4.8 \times 10^6$	$4.3 \times 10^7$	17.7
soc-orkut-dir	SO	$3.1 \times 10^6$	$1.2 \times 10^8$	76.3

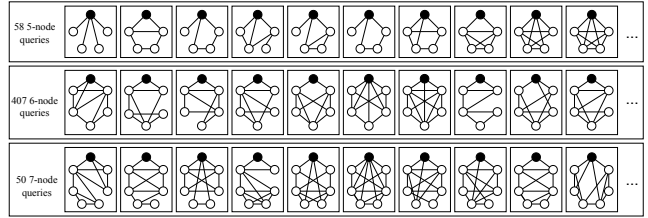


Figure 12: The 515 queries

## 8 EXPERIMENTS

**Algorithms:** We implemented SCOPE to compute a local subgraph counting query by tISO and *SymRs*, and evaluate tree decompositions by the Scope (Algorithm 4). To fully understand SCOPE, we have implemented its variants: SCOPE-Td, SCOPE-Tsd, and SCOPE-Ts. Here, SCOPE-Td is by tISO without *SymRs* and is evaluated by the Disc algorithm (Algorithm 3), SCOPE-Tsd is by both tISO and *SymRs* and is evaluated by the Disc algorithm, and SCOPE-Ts is by tISO without *SymRs* and is evaluated by the Scope algorithm. It is important to note that the Disc algorithm (Algorithm 3) is the algorithm used in DISC [101] to compute the aggregations in a tree decomposition, whereas DISC is the overall algorithm to compute subgraph counts. We also implemented 2 baselines, isoS and DISC1. Here, isoS is to count by directly enumerating iso-matches with *SymRs*. DISC1 is our implementation of DISC [101]. We also compare with EVOKE [73], DISC [101] and JESSE [63–65]. Like SCOPE, EVOKE and DISC are decomposition-based approaches, where EVOKE can only handle  $k$ -node pattern graphs for  $k \leq 5$ , and DISC is a general approach. JESSE is a matrix-based approach on a single machine that can handle any  $k$ -node pattern graphs. JESSE can only count all  $k$ -node patterns collectively for a specified  $k$ . It can not count selected queries. We omit other algorithms since they either only support global counting or are outperformed by DISC and EVOKE, as reported in [73, 101].

**12 Datasets:** We use 12 data graphs, including web graphs, recommendation networks, biological networks, collaboration networks, and social networks (Table 2). All data graphs are taken from [1, 2]. We deal with all graphs as simple undirected graphs.

**515 Queries:** We conduct testing using 515 local subgraph counting queries in total, as shown in Fig. 12: all 58 5-node queries, all 407 6-node queries, and 50 7-node queries. We randomly select 50 7-node queries from 2,423 7-node queries whose  $tw = 3$ . On average, each selected pattern has 11.7 edges, and uses 10.9 HOM counts

**Table 3: Results of the 6-node pattern in Fig. 4**

	Enumerated matches ( $\times 10^8$ )					Elapsed Time(seconds)				
	WS	BY	CA	SG	WW	WS	BY	CA	SG	WW
isoS	17.5	372	886	356	3,724	36.3	710	1,286	851	15,200
DISC1	4.6	65.4	113	89.8	1,127	15.8	237	153	572	16,400
SCOPE-Td	4.1	59.8	107	78.7	934	16.0	257	172	542	17,596
SCOPE-TSd	1.9	28.8	51.0	37.7	458	7.4	115	81.3	265	8,871
SCOPE-Ts	1.4	15.7	21.6	29.2	483	6.8	80.3	69.4	159	5,011
SCOPE	0.6	6.7	8.1	12.9	232	3.3	38.2	32.9	80.6	2,536

or 9.6 tISO counts to compute its ISO count. It is challenging. As an indication, the general approach DISC can only handle up to some simple 6-node queries [101], and there is no report published to test all 6-node queries in real-world data graphs.

**Settings:** We conduct all experiments on a single machine running CentOS 8 with Intel Xeon Silver 4215 32-core 2.5GHz CPU and 128GB memory. SCOPE and the variants, isoS, DISC1 and EVOKE are in C++. All these C++ implementations are compiled by g++ 8.5.0 with -O3 enabled and run with one thread. JESSE is in Java (Java 1.8). We use the default configuration and run it with one thread. DISC is a distributed system built on *Spark* (Spark 2.4.3). To remedy the difference in the programming language, we use the single machine configuration in [101] and run DISC with 32 threads. We also show the results of our C++ implementation DISC1. Like the previous works [63, 73, 101], we report the total time of running all queries in a batch. The time limit is 1 day. We also present the number of matches enumerated. We compare two algorithms  $\mathcal{A}_1$  and  $\mathcal{A}_2$  by the speedup of  $\mathcal{A}_2$  relative to  $\mathcal{A}_1$ , defined as the ratio of  $\mathcal{A}_1$ 's execution time to that of  $\mathcal{A}_2$ , and the reduction in enumerated matches, defined as the ratio of the number of matches enumerated by  $\mathcal{A}_1$  to the number by  $\mathcal{A}_2$ . We put the results for memory usage, preprocessing time and scalability tests in the full version.

### 8.1 A Simple Case Study: HOM, ISO, or tISO

As a case study to start, we consider a local subgraph counting query  $Q = (p, o)$ , where  $p$  is the 6-node pattern graph in Fig. 4 and  $o = u_1$ . We test 3 different approaches, namely, ISO, HOM, and tISO. For ISO, we use isoS which takes an ISO-based approach on  $p$  by directly enumerating iso-matches with *SymRs*. For HOM, we use DISC1, which generates 8 distinct trees for  $p$  by tree decomposition and processes each of the 8 trees by Disc [101]. For tISO, we use SCOPE and its variants, which generate 6 distinct trees by the tISOTOISO algorithm (Algorithm 1). This results in 11 tree nodes, and each of them is a 3/4-graph. The results are presented in Table 3. Decomposition approaches (DISC1, SCOPE, and the variants) enumerate much less matches compared to the enumeration approach isoS. Our SCOPE that combines tree decomposition, *SymRs*, and the Scope algorithm outperforms the others significantly.

### 8.2 The Three Batches of Queries

We have conducted testing using three batches of queries, namely, all 58 5-node queries, all 407 6-node queries, and 50 selected 7-node queries. The results are shown in Fig. 13, Fig. 14(a), and Fig. 14(c), respectively. Cases where the algorithm exceeded the time or memory limits are excluded from the figures.

**All 5-node Queries:** As shown in Fig. 13, JESSE can only complete all 5-node queries on WS and CA, and SCOPE outperforms JESSE

by more than 2 orders. SCOPE also outperforms DISC. Here, on the one machine setting SCOPE runs using one thread, whereas DISC runs using 32 threads on *Spark*. SCOPE is 133 $\times$  faster than DISC on average across 6 datasets, while DISC cannot compute the other 6 datasets. The maximum speedup observed is 227 $\times$  on the WS graph. EVOKE outperforms SCOPE in many cases. The main reason is that EVOKE does its best to deal with each of the 5-node queries in implementation, even though EVOKE takes a simple way to construct a 2-level tree for each  $p$ . However, the differences between EVOKE and SCOPE are not big. On average, EVOKE is about 1.9 $\times$  faster than SCOPE. But note that SCOPE is better than EVOKE in the three largest data graphs. SCOPE can compute all 5-node queries on SO, but EVOKE cannot. Also, EVOKE cannot support  $k$ -node queries when  $k > 5$ .

**All 6-node Queries:** For all 407 6-node queries in Fig. 14(a), only SCOPE can compute the batch of all 407 6-node queries. EVOKE does not support  $k = 6$ . JESSE and DISC run out of memory. To compare with DISC, we referred to our implementation DISC1 in Fig. 14(a). SCOPE outperforms all variants and DISC1 significantly. We also select 50 of the 6-node queries whose  $tw$  are 3. The results are in Fig. 14(b). SCOPE consistently outperforms DISC by more than 1 order of magnitude, in terms of both the elapsed time and the number of enumerated matches.

**The 50 7-node Queries:** The results are shown in Fig. 14(c). Like in the batch of 6-node queries, DISC can not run these 7-node queries due to memory exhaustion, so we use DISC1 in Fig. 14(c). SCOPE significantly outperforms the other variants and DISC1.

### 8.3 Effect of Proposed Techniques

**tISO-based counting:** We compare DISC1 and SCOPE-Td. Here, the former counts by HOM and the latter counts by tISO. Both use tree decomposition and use Disc to process each tree. SCOPE-Td does not use *SymR*. SCOPE-Td consistently performs better, showing average speedups of 1.1 $\times$ . Also, the average reductions in enumerated matches are 1.3 $\times$ , 1.4 $\times$ , and 1.3 $\times$ , respectively for  $k = 5, 6, 7$ . On one hand, ISO has an overhead for checking the injectivity of mappings. On the other hand, it benefits from having fewer iso-matches, and there are less trees to compute in tISO. Also, note that tISO can be used with *SymRs*, whereas HOM cannot. Overall, our experiments show that tISO performs better than HOM.

**The symmetry rules:** We investigate the benefit of symmetry rules by comparing SCOPE-Ts (without *SymRs*) and SCOPE (with *SymRs*). SCOPE significantly outperforms SCOPE-Ts. Take the CA graph as an example. SCOPE is 18.3 $\times$  faster than SCOPE-Ts for the total time of the 5-node batch query, 10.9 $\times$  faster for the 6-node batch query, and 5.7 $\times$  faster for the 7-node batch query. To further investigate the effectiveness of *SymRs* in different queries, we run all 6-node queries separately and compare SCOPE-Ts and SCOPE in Fig. 15. Here, we use a scatter plot where the x-axis is the density of each pattern, the colored points denote the  $tw$ , and the y-axis is the speedup/reduction in enumerated matches. The average speedup of the 407 individual queries in BY, CA, and RG are 7.2 $\times$ , 5.8 $\times$ , and 7.0 $\times$ , respectively, and the average reduction in enumerated matches are 5.8 $\times$ , 6.4 $\times$ , and 5.3 $\times$ , respectively. We observe that the speedup and reduction in enumerated matches are more significant for queries with larger  $tw$ . A tree node is likely to

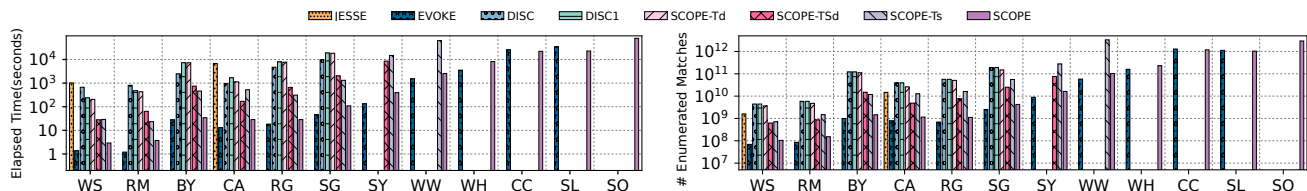


Figure 13: The batch with 58 5-node queries

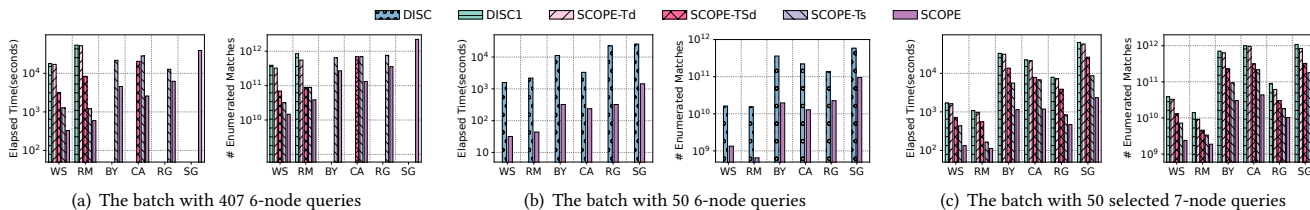


Figure 14: 6-node and 7-node queries

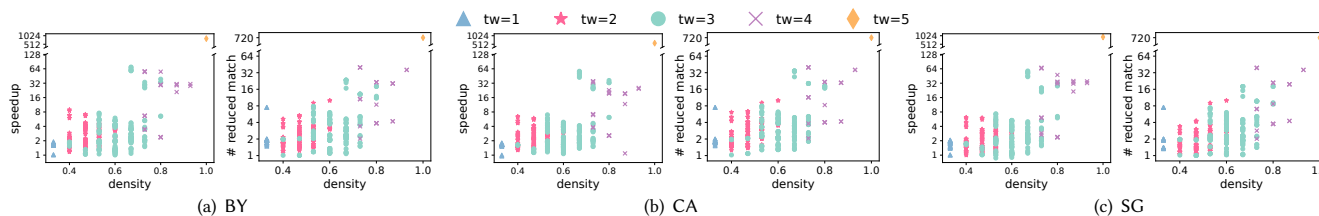


Figure 15: SCOPE vs SCOPE-Ts over the 407 individual 6-node queries

Table 4: Mean absolute error (MAE) for the ZINC dataset

model	GCN	GraphSage	GAT	MoNet	GatedGCN
Baseline	.356±.011	.455±.023	.464±.005	.260±.008	.340±.006
$\mathcal{F}$ -MPNN	.198±.003	.235±.005	.209±.006	.190±.002	.135±.010
$\mathcal{F}^+$ -MPNN	<b>.190±.021</b>	<b>.226±.014</b>	<b>.200±.003</b>	<b>.168±.011</b>	<b>.126±.009</b>

induce a dense subgraph since it can not be further partitioned. If the  $tw$  is large, then there exist some large and dense tree nodes in the tree decomposition, and there are many automorphism orbits and *SymRs* to be used. Note that a large  $tw$  means that the query is difficult to compute, so hard queries benefit more from *SymRs*.

**The multi-join algorithms:** We discussed two algorithms, namely, Disc (Algorithm 3) and Scope (Algorithm 4). We study the efficiency of the two by comparing SCOPE-TSd and SCOPE. The primary distinction lies in the algorithmic choice; SCOPE-TSd uses Disc, while SCOPE uses the Scope algorithm. In the three batch queries, SCOPE is more effective than SCOPE-TSd. For example, for the 6-node queries, SCOPE-TSd can only finish in 3 graphs, where SCOPE is 10.6x faster on average. This is due to the fact that Scope enumerates much less matches. SCOPE-TSd enumerates more matches since it can repeatedly enumerate the same tree nodes.

#### 8.4 Applying Subgraph Counts to GNN

In [9], the authors show that augmenting node feature with local subgraph counts can increase the expressive power of GNNs theoretically and empirically. Here, we further augment node features with 5-node and 6-node local subgraph counts that they do not use. We conducted extensive experimental studies to study the two

tasks conducted in [9], with five GNN architectures: GCN [55], GraphSage [36], GAT [92], MoNet [68], and GatedGCN [16], following the settings [9]. We study predicting the solubility of molecules in the ZINC dataset [27]. It has 12,000 graphs and each graph is a particular molecule. Table 4 shows the results. Here, Baseline uses atom types as node features.  $\mathcal{F}$ -MPNN [9] adds 3-10 cycle counts to node features, and we further add 95 non-zero counts taken from all 5/6-node patterns that  $\mathcal{F}$ -MPNN does not use, denoted as  $\mathcal{F}^+$ -MPNN.  $\mathcal{F}^+$ -MPNN has the smallest mean absolute error in all cases, and the improvement over the Baseline is significant. We also studied the node classification task in the full version.

## 9 CONCLUSION

We propose a novel decomposition-based approach for local subgraph counting,  $Q = (p, o)$ , by tree-decomposition-based counting (tISO-based counting), which can handle any  $k$ -node pattern graph  $p$  with the node orbit  $o$ . We confirm the efficiency of tISO-based counting by comparing our SCOPE with two state-of-the-art approaches, EVOKE and DISC, using 12 large datasets. EVOKE only supports pattern graphs up to 5 nodes, and SCOPE outperforms EVOKE in 3 large data graphs. For the batch of 5-node queries, SCOPE is 133 times faster than DISC on average over 6 datasets, while DISC cannot compute the other 6 datasets in the given time limit. For the batch of all 407 6-node queries, SCOPE is the only one that can compute on real large graphs in the given time limit.

## ACKNOWLEDGEMENT

This work was supported by the Research Grants Council of Hong Kong, China, No.14205520.

## REFERENCES

- [1] [n.d.]. Network Repository. Last Accessed: 18/04/2024. <https://networkrepository.com>.
- [2] [n.d.]. SNAP: Stanford Large Network Dataset Collection. Last Accessed: 18/04/2024. <https://snap.stanford.edu/data/index.html>.
- [3] Christopher R. Aberger, Susan Tu, Kunle Olukotun, and Christopher Ré. 2016. EmptyHeaded: A Relational Engine for Graph Processing. In *SIGMOD*. ACM, 431–446.
- [4] Nesreen K. Ahmed, Jennifer Neville, Ryan A. Rossi, and Nick G. Duffield. 2015. Efficient Graphlet Counting for Large Networks. In *2015 IEEE International Conference on Data Mining, ICDM 2015, Atlantic City, NJ, USA, November 14-17, 2015*. IEEE Computer Society, 1–10.
- [5] Omid Amini, Fedor V Fomin, and Saket Saurabh. 2009. Counting subgraphs via homomorphisms. In *Proc. of ICALP'09*. 71–82.
- [6] Khaled Ammar, Frank McSherry, Semih Salihoglu, and Manas Joglekar. 2018. Distributed Evaluation of Subgraph Queries Using Worst-case Optimal Low-memory Dataflows. *Proc. VLDB Endow.* 11, 6 (2018), 691–704.
- [7] Junya Arai, Yasuhiro Fujiwara, and Makoto Onizuka. 2023. GuP: Fast Subgraph Matching by Guard-based Pruning. *Proceedings of the ACM on Management of Data* 1, 2 (2023), 1–26.
- [8] Albert Aserias, Martin Grohe, and Dániel Marx. 2008. Size Bounds and Query Plans for Relational Joins. In *49th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2008, October 25-28, 2008, Philadelphia, PA, USA*. IEEE Computer Society, 739–748.
- [9] Pablo Barceló, Floris Geerts, Juan L. Reutter, and Maksimilian Ryschkov. 2021. Graph Neural Networks with Local Graph Parameters. In *NeurIPS*. 25280–25293.
- [10] Bibek Bhattacharai, Hang Liu, and H. Howie Huang. 2019. CECI: Compact Embedding Cluster Index for Scalable Subgraph Matching. In *SIGMOD*. ACM, 1447–1462.
- [11] Mansurul Bhuiyan, Mahmudur Rahman, Mahmuda Rahman, and Mohammad Al Hasan. 2012. GUISE: Uniform Sampling of Graphlets for Large Graph Analysis. In *12th IEEE International Conference on Data Mining, ICDM 2012, Brussels, Belgium, December 10-13, 2012*. IEEE Computer Society, 91–100.
- [12] Fei Bi, Lijun Chang, Xuemin Lin, Lu Qin, and Wenjie Zhang. 2016. Efficient Subgraph Matching by Postponing Cartesian Products. In *SIGMOD*. ACM, 1199–1214.
- [13] Christian Borgs, Jennifer Chayes, László Lovász, Vera T Sós, and Katalin Vesztegombi. 2006. Counting graph homomorphisms. In *Topics in Discrete Mathematics*. 315–371.
- [14] Marco Bressan, Flavio Chierichetti, Ravi Kumar, Stefano Leucci, and Alessandro Panconesi. 2018. Motif Counting Beyond Five Nodes. *ACM Trans. Knowl. Discov. Data* 12, 4 (2018), 48:1–48:25.
- [15] Marco Bressan, Stefano Leucci, and Alessandro Panconesi. 2019. Motivo: Fast Motif Counting via Succinct Color Coding and Adaptive Sampling. *Proc. VLDB Endow.* 12, 11 (2019), 1651–1663.
- [16] Xavier Bresson and Thomas Laurent. 2017. Residual Gated Graph ConvNets. *CoRR* abs/1711.07553 (2017).
- [17] Xiaowei Chen, Yongkun Li, Pinghui Wang, and John C. S. Lui. 2016. A General Framework for Estimating Graphlet Statistics via Random Walk. *Proc. VLDB Endow.* 10, 3 (2016), 253–264.
- [18] Xiaowei Chen and John C. S. Lui. 2018. Mining Graphlet Counts in Online Social Networks. *ACM Trans. Knowl. Discov. Data* 12, 4 (2018), 41:1–41:38.
- [19] Luigi P. Cordella, Pasquale Foggia, Carlo Sansone, and Mario Vento. 2004. A (Sub)Graph Isomorphism Algorithm for Matching Large Graphs. *IEEE Trans. Pattern Anal. Mach. Intell.* 26, 10 (2004), 1367–1372.
- [20] Radu Curticapean, Holger Dell, and Dániel Marx. 2017. Homomorphisms are a good basis for counting small subgraphs. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017, Montreal, QC, Canada, June 19-23, 2017*. ACM, 210–223.
- [21] Marek Cygan, Fedor V. Fomin, Lukasz Kowalik, Daniel Lokshtanov, Dániel Marx, Marcin Pilipczuk, Michal Pilipczuk, and Saket Saurabh. 2015. *Parameterized Algorithms*. Springer.
- [22] Victor Dalmau and Peter Jonsson. 2004. The complexity of counting homomorphisms seen from the other side. *Theoretical Computer Science* 329, 1-3 (2004), 315–323.
- [23] Vachik S. Dave, Nesreen K. Ahmed, and Mohammad Al Hasan. 2017. E-CLOg: Counting edge-centric local graphlets. In *2017 IEEE International Conference on Big Data (IEEE BigData 2017), Boston, MA, USA, December 11-14, 2017*. IEEE Computer Society, 586–595.
- [24] Richard M. Foote David S. Dummit. 2003. *Abstract Algebra* (3 ed.). Wiley.
- [25] Josep Diaz, Maria Serna, and Dimitrios M Thilikos. 2002. Counting H-colorings of partial k-trees. *Theoretical Computer Science* 281, 1-2 (2002), 291–309.
- [26] Vinicius Vitor dos Santos Dias, Carlos H. C. Teixeira, Dorgival O. Guedes, Wagner Meira Jr., and Srinivasan Parthasarathy. 2019. Fractal: A General-Purpose Graph Pattern Mining System. In *SIGMOD*. ACM, 1357–1374.
- [27] Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. 2023. Benchmarking Graph Neural Networks. *J. Mach. Learn. Res.* 24 (2023), 43:1–43:48.
- [28] Martin Dyer and Catherine Greenhill. 2000. The complexity of counting graph homomorphisms. In *Proc. of SODA'00*. 246–255.
- [29] Ethan R. Elenberg, Karthikeyan Shanmugam, Michael Borokhovich, and Alexandros G. Dimakis. 2015. Beyond Triangles: A Distributed Framework for Estimating 3-profiles of Large Graphs. In *SIGKDD*. 229–238.
- [30] Ethan R. Elenberg, Karthikeyan Shanmugam, Michael Borokhovich, and Alexandros G. Dimakis. 2016. Distributed Estimation of Graph 4-Profiles. In *Proc. of WWW'16*. 483–493.
- [31] Jörg Flum and Martin Grohe. 2006. *Parameterized Complexity Theory*. Springer.
- [32] Georg Gottlob, Gianluigi Greco, Nicola Leone, and Francesco Scarcello. 2016. Hypertree Decompositions: Questions and Answers. In *Proc. of PODS'16*. 57–74.
- [33] Joshua A. Grochow and Manolis Kellis. 2007. Network Motif Discovery Using Subgraph Enumeration and Symmetry-Breaking. In *Research in Computational Molecular Biology, 11th Annual International Conference, RECOMB 2007, Oakland, CA, USA, April 21-25, 2007, Proceedings (Lecture Notes in Computer Science)*, Vol. 4453. Springer, 92–106.
- [34] Martin Grohe. 2007. The complexity of homomorphism and constraint satisfaction problems seen from the other side. *J. ACM* 54, 1 (2007), 1–24.
- [35] Wentian Guo, Yuchen Li, Mo Sha, Bingsheng He, Xiaokui Xiao, and Kian-Lee Tan. 2020. GPU-Accelerated Subgraph Enumeration on Partitioned Graphs. In *SIGMOD*. ACM, 1067–1082.
- [36] William L. Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *NeurIPS*. 1024–1034.
- [37] Guyue Han and Harish Sethu. 2016. Waddling Random Walk: Fast and Accurate Mining of Motif Statistics in Large Graphs. In *IEEE 16th International Conference on Data Mining, ICDM 2016, December 12-15, 2016, Barcelona, Spain*. IEEE Computer Society, 181–190.
- [38] Myoungji Han, Hyunjoon Kim, Geonmo Gu, Kunsoo Park, and Wook-Shin Han. 2019. Efficient Subgraph Matching: Harmonizing Dynamic Programming, Adaptive Matching Order, and Failing Set Together. In *SIGMOD*. 1429–1446.
- [39] Wook-Shin Han, Jinsoo Lee, and Jeong-Hoon Lee. 2013. Turboiso: Towards Ultrafast and Robust Subgraph Isomorphism Search in Large Graph Databases. In *SIGMOD*. 337–348.
- [40] Pavol Hell and Jaroslav Nešetřil. 1990. On the complexity of H-coloring. *Journal of Combinatorial Theory, Series B* 48, 1 (1990), 92–110.
- [41] Tomaz Hočevar and Janez Demšar. 2014. A combinatorial approach to graphlet counting. *Bioinformatics* 30, 4 (2014), 559–565.
- [42] Tomaz Hočevar and Janez Demšar. 2017. Combinatorial algorithm for counting small induced graphs and orbits. *PLoS one* 12, 2 (2017), e0171428.
- [43] Maarten Houbraeken, Sofie Demeyer, Tom Michoel, Pieter Audenaert, Didier Colle, and Mario Pickavet. 2014. The Index-based Subgraph Matching Algorithm with General Symmetries (ISMAGS): exploiting symmetry for faster subgraph enumeration. *PLoS one* 9, 5 (2014), e97896.
- [44] Royi Itzhack, Yelena Mogilevski, and Yoram Louzoun. 2007. An optimal algorithm for counting network motifs. *Physica A: Statistical Mechanics and its Applications* 381 (2007), 482–490.
- [45] Kasra Jamshidi, Rakesh Mahadasa, and Keval Vora. 2020. Peregrine: a pattern-aware graph mining system. In *EuroSys '20: Fifteenth EuroSys Conference 2020, Heraklion, Greece, April 27-30, 2020*. ACM, 13:1–13:16.
- [46] Madhav Jha, C. Seshadhri, and Ali Pinar. 2015. Path Sampling: A Fast and Provable Method for Estimating 4-Vertex Subgraph Counts. In *Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015*. ACM, 495–505.
- [47] Xun Jian, Zhiyuan Li, and Lei Chen. 2023. SUFF: Accelerating Subgraph Matching with Historical Data. *Proc. VLDB Endow.* 16, 7 (2023), 1699–1711.
- [48] Tatiana Jin, Boyang Li, Yichao Li, Qihui Zhou, Qianli Ma, Yunjian Zhao, Hongzhi Chen, and James Cheng. 2023. Circinus: Fast redundancy-reduced subgraph matching. *Proceedings of the ACM on Management of Data* 1, 1 (2023), 1–26.
- [49] Oren Kalinsky, Yoav Etsion, and Benny Kimelfeld. 2016. Flexible caching in trie joins. *arXiv preprint arXiv:1602.08721* (2016).
- [50] Zahra Razaghi Moghadam Kashani, Hayedeh Ahrabian, Elahe Elahi, Abbas Nowzari-Dalini, Elnaz Saberi Ansari, Sahar Asadi, Shahin Mohammadi, Falk Schreiber, and Ali Masoudi-Nejad. 2009. Kavosh: a new algorithm for finding network motifs. *BMC Bioinform.* 10 (2009), 318.
- [51] Nadav Kashtan, Shalev Itzkovitz, Ron Milo, and Uri Alon. 2004. Efficient sampling algorithm for estimating subgraph concentrations and detecting network motifs. *Bioinform.* 20, 11 (2004), 1746–1758.
- [52] Sahand Khakabimamaghani, Iman Sharafuddin, Norbert Dichter, Ina Koch, and Ali Masoudi-Nejad. 2013. Quatexelero: an accelerated exact network motif detection algorithm. *PLoS one* 8, 7 (2013), e68073.
- [53] Hyunjoon Kim, Yunyoung Choi, Kunsoo Park, Xuemin Lin, Seok-Hee Hong, and Wook-Shin Han. 2021. Versatile Equivalences: Speeding up Subgraph Query Processing and Subgraph Matching. In *SIGMOD*. ACM, 925–937.
- [54] Hyunjoon Kim, Yunyoung Choi, Kunsoo Park, Xuemin Lin, Seok-Hee Hong, and Wook-Shin Han. 2023. Fast subgraph query processing and subgraph matching via static and dynamic equivalences. *The VLDB journal* 32, 2 (2023), 343–368.
- [55] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with

- Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- [56] Tamara G. Kolda, Ali Pinar, and C. Seshadhri. 2013. Triadic Measures on Graphs: The Power of Wedge Sampling. In *Proceedings of the 13th SIAM International Conference on Data Mining, May 2-4, 2013, Austin, Texas, USA*. SIAM, 10–18.
- [57] Longbin Lai, Lu Qin, Xuemin Lin, and Lijun Chang. 2015. Scalable Subgraph Enumeration in MapReduce. *Proc. VLDB Endow.* 8, 10 (2015), 974–985.
- [58] Longbin Lai, Lu Qin, Xuemin Lin, and Lijun Chang. 2017. Scalable subgraph enumeration in MapReduce: a cost-oriented approach. *VLDB J.* 26, 3 (2017), 421–446.
- [59] Longbin Lai, Zhu Qing, Zhengyi Yang, Xin Jin, Zhengmin Lai, Ran Wang, Kongzhang Hao, Xuemin Lin, Lu Qin, Wenjie Zhang, Ying Zhang, Zhengping Qian, and Jingren Zhou. 2019. Distributed Subgraph Matching on Timely Dataflow. *Proc. VLDB Endow.* 12, 10 (2019), 1099–1112.
- [60] Dror Marcus and Yuval Shavitt. 2010. Efficient Counting of Network Motifs. In *30th IEEE International Conference on Distributed Computing Systems Workshops (ICDCS 2010 Workshops), 21-25 June 2010, Genova, Italy*. IEEE Computer Society, 92–98.
- [61] Dániel Marx and Michal Pilipczuk. 2014. Everything you always wanted to know about the parameterized complexity of Subgraph Isomorphism (but were afraid to ask). In *31st International Symposium on Theoretical Aspects of Computer Science (STACS 2014), STACS 2014, March 5-8, 2014, Lyon, France (LIPIcs)*, Vol. 25. 542–553.
- [62] Luis A. A. Meira, Vinicius R. Máximo, Álvaro Luiz Fazenda, and Arlindo Flávio da Conceição. 2014. acc-Motif: Accelerated Network Motif Detection. *IEEE ACM Trans. Comput. Biol. Bioinform.* 11, 5 (2014), 853–862.
- [63] Ine Melckenbeeck, Pieter Audenaert, Didier Colle, and Mario Pickavet. 2018. Efficiently counting all orbits of graphlets of any order in a graph using autogenerated equations. *Bioinform.* 34, 8 (2018), 1372–1380.
- [64] Ine Melckenbeeck, Pieter Audenaert, Tom Michoel, Didier Colle, and Mario Pickavet. 2016. An algorithm to automatically generate the combinatorial orbit counting equations. *PLoS one* 11, 1 (2016), e0147078.
- [65] Ine Melckenbeeck, Pieter Audenaert, Thomas Van Parys, Yves Van de Peer, Didier Colle, and Mario Pickavet. 2019. Optimising orbit counting of arbitrary order by equation selection. *BMC Bioinform.* 20, 1 (2019), 27:1–27:13.
- [66] Amine Mhedhbi and Semih Salihoglu. 2019. Optimizing subgraph queries by combining binary and worst-case optimal joins. *Proc. VLDB Endow.* 12, 11 (2019), 1692–1704.
- [67] Ron Milo, Shai Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri Chklovskii, and Uri Alon. 2002. Network motifs: simple building blocks of complex networks. *Science* 298, 5594 (2002), 824–827.
- [68] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodolà, Jan Svoboda, and Michael M. Bronstein. 2017. Geometric Deep Learning on Graphs and Manifolds Using Mixture Model CNNs. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 5425–5434.
- [69] Hung Q. Ngo, Ely Porat, Christopher Ré, and Atri Rudra. 2012. Worst-case Optimal Join Algorithms: [Extended Abstract]. In *Proc. of PODS'12*. 37–48.
- [70] Hung Q Ngo, Christopher Ré, and Atri Rudra. 2014. Skew strikes back: new developments in the theory of join algorithms. *ACM SIGMOD Record* 42, 4 (2014), 5–16.
- [71] Pedro Paredes and Pedro Manuel Pinto Ribeiro. 2013. Towards a faster network-centric subgraph census. In *Advances in Social Networks Analysis and Mining 2013, ASONAM '13, Niagara, ON, Canada - August 25 - 29, 2013*. ACM, 264–271.
- [72] Ha-Myung Park, Sung-Hyon Myaeng, and U Kang. 2016. Pte: Enumerating trillion triangles on distributed systems. In *SIGKDD*. 1115–1124.
- [73] Noujan Pashanasangi and C. Seshadhri. 2020. Efficiently Counting Vertex Orbits of All 5-vertex Subgraphs, by EVOKE. In *WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020*. ACM, 447–455.
- [74] Ali Pinar, C. Seshadhri, and Vaidyanathan Vishal. 2017. ESCAPE: Efficiently Counting All 5-Vertex Subgraphs. In *Proceedings of the 26th International Conference on World Wide Web, WWW 2017, Perth, Australia, April 3-7, 2017*. ACM, 1431–1440.
- [75] Natasa Przulj, Derek G. Corneil, and Igor Jurisica. 2006. Efficient estimation of graphlet frequency distributions in protein-protein interaction networks. *Bioinform.* 22, 8 (2006), 974–980.
- [76] Chendi Qian, Gaurav Rattan, Floris Geerts, Mathias Niepert, and Christopher Morris. 2022. Ordered Subgraph Aggregation Networks. In *NeurIPS*.
- [77] Miao Qiao, Hao Zhang, and Hong Cheng. 2017. Subgraph Matching: On Compression and Computation. *Proc. VLDB Endow.* 11, 2 (2017), 176–188.
- [78] Mahmudur Rahman, Mansurul Alam Bhuiyan, and Mohammad Al Hasan. 2014. Graft: An Efficient Graphlet Counting Method for Large Graph Analysis. *IEEE Trans. Knowl. Data Eng.* 26, 10 (2014), 2466–2478.
- [79] Xuguang Ren, Junhu Wang, Wook-Shin Han, and Jeffrey Xu Yu. 2019. Fast and Robust Distributed Subgraph Enumeration. *Proc. VLDB Endow.* 12, 11 (2019), 1344–1356.
- [80] Pedro Ribeiro, Pedro Paredes, Miguel E. P. Silva, David Aparício, and Fernando M. A. Silva. 2022. A Survey on Subgraph Counting: Concepts, Algorithms, and Applications to Network Motifs and Graphlets. *ACM Comput. Surv.* 54, 2 (2022), 28:1–28:36.
- [81] Pedro Manuel Pinto Ribeiro and Fernando M. A. Silva. 2010. g-tries: an efficient data structure for discovering network motifs. In *Proceedings of the 2010 ACM Symposium on Applied Computing (SAC), Sierre, Switzerland, March 22-26, 2010*. ACM, 1559–1566.
- [82] P. Griffiths Selinger, M. M. Astrahan, D. D. Chamberlin, R. A. Lorie, and T. G. Price. 1979. Access Path Selection in a Relational Database Management System. In *SIGMOD*. 23–34.
- [83] Haichuan Shang, Ying Zhang, Xuemin Lin, and Jeffrey Xu Yu. 2008. Taming Verification Hardness: An Efficient Algorithm for Testing Subgraph Isomorphism. *Proc. VLDB Endow.* 1, 1 (2008), 364–375.
- [84] Tianhui Shi, Mingshu Zhai, Yi Xu, and Jidong Zhai. 2020. GraphPi: high performance graph pattern matching through effective redundancy elimination. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2020, Virtual Event / Atlanta, Georgia, USA, November 9-19, 2020*. IEEE/ACM, 100.
- [85] Shixuan Sun, Yulin Che, Lipeng Wang, and Qiong Luo. 2019. Efficient Parallel Subgraph Enumeration on a Single Machine. In *ICDE*. IEEE, 232–243.
- [86] Shixuan Sun and Qiong Luo. 2020. In-Memory Subgraph Matching: An In-depth Study. In *SIGMOD*. ACM, 1083–1098.
- [87] Shixuan Sun and Qiong Luo. 2022. Subgraph Matching With Effective Matching Order and Indexing. *IEEE Trans. Knowl. Data Eng.* 34, 1 (2022), 491–505.
- [88] Shixuan Sun, Xibo Sun, Yulin Che, Qiong Luo, and Bingsheng He. 2020. Rapid-Match: A Holistic Approach to Subgraph Query Processing. *Proc. VLDB Endow.* 14, 2 (2020), 176–188.
- [89] Susan Tu and Christopher Ré. 2015. Duncemap: Query plans using generalized hypertree decompositions. In *SIGMOD*. 2077–2078.
- [90] Julian R. Ullmann. 1976. An Algorithm for Subgraph Isomorphism. *J. ACM* 23, 1 (1976), 31–42.
- [91] Todd L. Veldhuizen. 2012. Leapfrog Triejoin: a worst-case optimal join algorithm. *CoRR* abs/1210.0481 (2012).
- [92] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- [93] Hanchen Wang, Rong Hu, Ying Zhang, Lu Qin, Wei Wang, and Wenjie Zhang. 2022. Neural Subgraph Counting with Wasserstein Estimator. In *SIGMOD*. ACM, 160–175.
- [94] Pinghui Wang, Yiyang Qi, John C. S. Lui, Don Towsley, Junzhou Zhao, and Jing Tao. 2019. Inferring Higher-Order Structure Statistics of Large Networks from Sampled Edges. *IEEE Trans. Knowl. Data Eng.* 31, 1 (2019), 61–74.
- [95] Pinghui Wang, Junzhou Zhao, Xiangliang Zhang, Zhenguo Li, Jiefeng Cheng, John C. S. Lui, Don Towsley, Jing Tao, and Xiaohong Guan. 2018. MOSS-5: A Fast Method of Approximating Counts of 5-Node Graphlets in Large Graphs. *IEEE Trans. Knowl. Data Eng.* 30, 1 (2018), 73–86.
- [96] Zhaokang Wang, Rong Gu, Weiwei Hu, Chunfeng Yuan, and Yihua Huang. 2019. BENU: Distributed Subgraph Enumeration with Backtracking-Based Framework. In *ICDE*. IEEE, 136–147.
- [97] Sebastian Wernicke. 2005. A Faster Algorithm for Detecting Network Motifs. In *Algorithms in Bioinformatics, 5th International Workshop, WABI 2005, Mallorca, Spain, October 3-6, 2005, Proceedings (Lecture Notes in Computer Science)*, Vol. 3692. Springer, 165–177.
- [98] Sebastian Wernicke and Florian Rasche. 2006. FANMOD: a tool for fast network motif detection. *Bioinform.* 22, 9 (2006), 1152–1153.
- [99] Chen Yang, Min Lyu, Yongkun Li, Qianqian Zhao, and Yinlong Xu. 2018. SSRW: A Scalable Algorithm for Estimating Graphlet Statistics Based on Random Walk. In *Database Systems for Advanced Applications - 23rd International Conference, DASFAA 2018, Gold Coast, QLD, Australia, May 21-24, 2018, Proceedings, Part I (Lecture Notes in Computer Science)*, Vol. 10827. Springer, 272–288.
- [100] Zhengyi Yang, Longbin Lai, Xuemin Lin, Kongzhang Hao, and Wenjie Zhang. 2021. HUGE: An Efficient and Scalable Subgraph Enumeration System. In *SIGMOD*. ACM, 2049–2062.
- [101] Hao Zhang, Jeffrey Xu Yu, Yikai Zhang, Kangfei Zhao, and Hong Cheng. 2020. Distributed Subgraph Counting: A General Approach. *Proc. VLDB Endow.* 13, 11 (2020), 2493–2507.
- [102] Kangfei Zhao, Jeffrey Xu Yu, Hao Zhang, Qiyang Li, and Yu Rong. 2021. A Learned Sketch for Subgraph Counting. In *SIGMOD*. ACM, 2142–2155.
- [103] Peixiang Zhao and Jiawei Han. 2010. On Graph Query Optimization in Large Networks. *Proc. VLDB Endow.* 3, 1-2 (2010), 340–351.