

# An Approach for Addressing Internally-Disconnected Communities in Louvain Algorithm

Subhajit Sahu

subhajit.sahu@research.iiit.ac.in

IIIT Hyderabad

Hyderabad, Telangana, India

## ABSTRACT

Community detection is the problem of identifying densely connected clusters of nodes within a network. The Louvain algorithm is a widely used method for this task, but it can produce communities that are internally disconnected. To address this, the Leiden algorithm was introduced. In this technical report, we propose another approach to mitigate this issue. On a system with two 16-core Intel Xeon Gold 6226R processors, our new parallel algorithm GSP-Louvain, based on the Louvain algorithm, addresses this issue, and outperforms the original Leiden, igraph Leiden, and NetworKit Leiden by 341×, 83×, and 6.1× respectively - achieving a processing rate of 328M edges/s on a 3.8B edge graph. Furthermore, GSP-Louvain improves performance at a rate of 1.5× for every doubling of threads.

## KEYWORDS

Community detection, Internally-disconnected communities, Parallel Louvain implementation

## 1 INTRODUCTION

Community detection is the process of identifying groups of vertices characterized by dense internal connections and sparse connections between the groups [12]. These groups, referred to as communities or clusters [1], offer valuable insights into the organization and functionality of networks [12]. Community detection finds applications in various fields, including topic discovery in text mining, protein annotation in bioinformatics, recommendation systems in e-commerce, and targeted advertising in digital marketing [15].

Modularity maximization is a frequently employed method for community detection. Modularity measures the difference between the fraction of edges within communities and the expected fraction of edges under random distribution, and ranges from  $-0.5$  to  $1$  [12, 25]. However, optimizing for modularity is an NP-hard problem [4]. An additional challenge is the lack of prior knowledge about the number and size distribution of communities [3]. Heuristic-based approaches are thus used for community detection [3, 7, 10, 19, 27–29, 38, 39, 41, 44, 46]. The identified communities are considered intrinsic when solely based on network topology and disjoint when each vertex belongs to only one community [8, 15].

The Louvain method [3] is a popular heuristic-based approach for intrinsic and disjoint community detection, and has been identified as one of the fastest and top-performing algorithms [21, 45]. It utilizes a two-phase approach involving local-moving and aggregation phases to iteratively optimize the modularity metric [25].

Although widely used, the Louvain method has been noted for generating internally-disconnected and poorly connected communities [39]. In response, Traag et al. [39] introduce the Leiden algorithm, which incorporates an extra refinement phase between the local-moving and aggregation phases. This refinement step enables vertices to explore and potentially create sub-communities within the identified communities from the local-moving phase [39]. In this report, we propose another BFS-based approach for addressing the same issue. This is different from a number of earlier works, which tackled internally-disconnected communities as a post-processing step [15, 16, 24, 27, 43].

Furthermore, the proliferation of data and their graph representations has reached unprecedented levels in recent years. However, applying the original Leiden algorithm to massive graphs has posed computational challenges, primarily due to its inherently sequential nature, akin to the Louvain method [17]. In contexts prioritizing scalability, the development of optimized algorithms which address the issue of internally disconnected communities becomes essential, especially in the multicore/shared memory setting.

### 1.1 Our Contributions

In this report, we propose GSP-Louvain<sup>1</sup>, another approach to mitigate the issue of disconnected communities with the Louvain algorithm. On a system equipped with two 16-core Intel Xeon Gold 6226R processors, we demonstrate that GSP-Louvain resolves this issue and achieves a processing rate of 328M edges/s on a 3.8B edge graph. It outperforms the original Leiden, igraph Leiden, and NetworKit Leiden by 341×, 83×, and 6.1×, respectively. The identified communities are of similar quality as the first two implementations and 25% higher quality than NetworKit. Furthermore, GSP-Louvain improve performance at a rate of 1.5× for every doubling of threads.

## 2 RELATED WORK

Communities in networks represent functional units or meta-nodes [5]. Identifying these natural divisions in an unsupervised manner is a crucial graph analytics problem in many domains. Various schemes have been developed for finding communities [3, 7, 10, 19, 27–29, 38, 39, 41, 44, 46]. These schemes can be categorized into three basic approaches: bottom-up, top-down, and data-structure based, with further classification possible within the bottom-up approach [35]. They can also be classified into divisive, agglomerative, and multi-level methods [47]. To evaluate the quality of these methods, fitness scores like modularity are commonly used. Modularity [25] is a metric that ranges from  $-0.5$  to  $1$ , and measures the relative density of links within communities compared to those

<sup>1</sup><https://github.com/puzzlef/louvain-communities-openmp>

outside. Optimizing modularity theoretically results in the best possible clustering of the nodes [21]. However, as going through each possible clustering of the nodes is impractical [12], heuristic algorithms such as the Louvain method [3] are used.

The Louvain method, proposed by Blondel et al. [3] from the University of Louvain, is a greedy, multi-level, modularity-optimization based algorithm for intrinsic and disjoint community detection [21, 25]. It utilizes a two-phase approach involving local-moving and aggregation phases to iteratively optimize the modularity metric [3], and is recognized as one of the fastest and top-performing algorithms [21, 45]. Various algorithmic improvements [17, 23, 30, 31, 34, 40] and parallelization strategies [6, 11, 17, 23, 42, 48] have been proposed for the Louvain method. Several open-source implementations and software packages have been developed for community detection using Parallel Louvain Algorithm, including Vite [14], Grappolo [17], and NetworKit [36]. Note however that community detection methods which rely on modularity maximization are known to face the resolution limit problem, which hinders the identification of smaller communities [13].

Although favored for its ability to identify communities with high modularity, the Louvain method often produces internally disconnected communities due to vertices acting as bridges moving to other communities during iterations (see Section 3.4). This issue worsens with further iterations without decreasing the quality function [39]. To overcome these issues, Traag et al. [39] from the University of Leiden, propose the Leiden algorithm, which introduces a refinement phase after the local-moving phase of the Louvain method. In this refinement phase, vertices undergo additional local moves based on delta-modularity, allowing the discovery of sub-communities within the initial communities. Shi et al. [34] also introduce a similar refinement phase after the local-moving phase of the Louvain method to minimize poor clusters.

A few open-source implementations and software packages exist for community detection using the Leiden algorithm. The original implementation, `libleidenalg` [39], is written in C++ and has a Python interface called `leidenalg`. `NetworKit` [36], a software package designed for analyzing graph data sets with billions of connections, features a parallel implementation of the Leiden algorithm by Nguyen [26]. This implementation utilizes global queues for vertex pruning and vertex and community locking for updating communities. Another package, `igraph` [9], is written in C and has Python, R, and Mathematica frontends – also includes an implementation of the Leiden algorithm.

Internally disconnected communities are not a new problem. Internally-disconnected communities can be produced by label propagation algorithms [15, 27], multi-level algorithms [3], genetic algorithms [18], expectation minimization/maximization algorithms [2, 16], among others. Raghavan et al. [27] note that their Label Propagation Algorithm (LPA) for community detection may identify disconnected communities. In such a case, they suggest applying a Breadth-First Search (BFS) on the subnetworks of each individual group to separate the disconnected communities, with a time complexity of  $O(M + N)$ . Gregory [15] introduced the Community Overlap Propagation Algorithm (COPRA) as an extension of LPA. In their algorithm, they eliminate communities completely contained within others and use a similar method as Raghavan et al. [27] to split any returned disconnected communities into

connected ones. Hafez et al. [16] present a community detection algorithm leveraging Bayesian Network and Expectation Minimization (BNEM). They include a final step to examine the result for potentially containing disconnected communities within a single community. If this situation arises, new community labels are assigned to the disconnected components. This scenario occurs when the network may have more communities than the specified number  $k$  in the algorithm. Hesamipour et al. [18] propose a genetic algorithm for community detection, integrating similarity-based and modularity-based approaches. Their method uses an MST-based representation for addressing issues like disconnected communities and ineffective mutations.

### 3 PRELIMINARIES

Consider an undirected graph denoted as  $G(V, E, w)$ , where  $V$  stands for the vertex set,  $E$  represents the edge set, and  $w_{ij} = w_{ji}$  indicates the weight associated with each edge. In the scenario of an unweighted graph, we assume a unit weight for every edge ( $w_{ij} = 1$ ). Moreover, the neighbors of a vertex  $i$  are referred to as  $J_i = \{j \mid (i, j) \in E\}$ , the weighted degree of each vertex as  $K_i = \sum_{j \in J_i} w_{ij}$ , the total number of vertices as  $N = |V|$ , the total number of edges as  $M = |E|$ , and the sum of edge weights in the undirected graph as  $m = \sum_{i,j \in V} w_{ij}/2$ .

#### 3.1 Community detection

Disjoint community detection entails identifying a community membership mapping,  $C : V \rightarrow \Gamma$ , wherein each vertex  $i \in V$  is assigned a community ID  $c$  from the set of community IDs  $\Gamma$ . We denote the vertices of a community  $c \in \Gamma$  as  $V_c$ , and the community that a vertex  $i$  belongs to as  $C_i$ . Additionally, we denote the neighbors of vertex  $i$  belonging to a community  $c$  as  $J_{i \rightarrow c} = \{j \mid j \in J_i \text{ and } C_j = c\}$ , the sum of those edge weights as  $K_{i \rightarrow c} = \sum_{j \in J_{i \rightarrow c}} w_{ij}$ , the sum of weights of edges within a community  $c$  as  $\sigma_c = \sum_{(i,j) \in E \text{ and } C_i=C_j=c} w_{ij}$ , and the total edge weight of a community  $c$  as  $\Sigma_c = \sum_{(i,j) \in E \text{ and } C_i=c} w_{ij}$  [22].

#### 3.2 Modularity

Modularity is a metric used for assessing the quality of communities identified by heuristic-based community detection algorithms. It is calculated as the difference between the fraction of edges within communities and the expected fraction if edges were randomly distributed, and ranges from  $-0.5$  to  $1$ , where higher values indicate superior results [4]. The modularity  $Q$  of identified communities is determined using Equation 1, where  $\delta$  represents the Kronecker delta function ( $\delta(x, y) = 1$  if  $x = y$ ,  $0$  otherwise). The *delta modularity* of moving a vertex  $i$  from community  $d$  to community  $c$ , denoted as  $\Delta Q_{i:d \rightarrow c}$ , can be determined using Equation 2.

$$Q = \frac{1}{2m} \sum_{(i,j) \in E} \left[ w_{ij} - \frac{K_i K_j}{2m} \right] \delta(C_i, C_j) = \sum_{c \in \Gamma} \left[ \frac{\sigma_c}{2m} - \left( \frac{\Sigma_c}{2m} \right)^2 \right] \quad (1)$$

$$\Delta Q_{i:d \rightarrow c} = \frac{1}{m} (K_{i \rightarrow c} - K_{i \rightarrow d}) - \frac{K_i}{2m^2} (K_i + \Sigma_c - \Sigma_d) \quad (2)$$

### 3.3 Louvain algorithm

The Louvain method, introduced by Blondel et al. [3], is a greedy algorithm that optimizes modularity to identify high quality disjoint communities in large networks. It has a time complexity of  $O(L|E|)$ , where  $L$  is the total number of iterations performed, and a space complexity of  $O(|V| + |E|)$  [21]. This algorithm consists of two phases: the *local-moving phase*, wherein each vertex  $i$  greedily decides to join the community of one of its neighbors  $j \in J_i$  to maximize the gain in modularity  $\Delta Q_{i:C_i \rightarrow C_j}$  (using Equation 2), and the *aggregation phase*, during which all vertices within a community are combined into a single super-vertex. These phases constitute one pass, which is repeated until no further improvement in modularity is achieved [3, 22].

### 3.4 Possibility of Internally-disconnected communities with the Louvain algorithm

The Louvain method, though effective, has been noted to potentially identify internally disconnected communities [39]. This is illustrated with an example in Figure 1. Figure 1(a) shows the initial community structure after running a few iterations of the Louvain algorithm. It includes four communities labeled  $C_1$ ,  $C_2$ ,  $C_3$ ,  $C_4$ , and vertices 1 to 7 are grouped in community  $C_1$ . After a few additional iterations, in Figure 1(b), communities  $C_2$ ,  $C_3$ , and  $C_4$  merge together into  $C_3$ , due to strong connections among themselves. As vertex 4 is now more strongly connected to community  $C_3$ , it shifts from community  $C_1$  to join community  $C_3$  instead, in Figure 1(c). This results in the internal disconnection of community  $C_1$ , as vertices 1, 2, 3, 5, 6, and 7 retain their locally optimal assignments. Additionally, once all nodes are optimally assigned, the algorithm aggregates the graph. If an internally disconnected community becomes a node in the aggregated graph, it remains disconnected unless it combines with another community acting as a bridge. With subsequent passes, these disconnected communities are prone to steering the solution towards a lower local optima.

### 3.5 Leiden algorithm

The Leiden algorithm, proposed by Traag et al. [39], is a multi-level community detection technique that extends the Louvain method. It consists of three key phases. In the *local-moving phase*, each vertex  $i$  optimizes its community assignment by greedily selecting to join the community of one of its neighbors  $j \in J_i$ , aiming to maximize the modularity gain  $\Delta Q_{i:C_i \rightarrow C_j}$ , as defined by Equation 2, akin to the Louvain method. During the *refinement phase*, vertices within each community undergo further updates to their community memberships, starting from singleton partitions. Unlike the local-moving phase however, these updates are not strictly greedy. Instead, vertices may move to any community within their bounds where the modularity increases, with the probability of joining a neighboring community proportional to the delta-modularity of the move. The level of randomness in these moves is governed by a parameter  $\theta > 0$ . This randomized approach facilitates the identification of higher quality sub-communities within the communities established during the local-moving phase. Finally, in the *aggregation phase*, all vertices within each refined partition are combined into super-vertices, with an initial community assignment derived

from the local-moving phase [39]. The time complexity of the Leiden algorithm is  $O(L|E|)$ , where  $L$  represents the total number of iterations performed, and its space complexity is  $O(|V| + |E|)$ .

## 4 APPROACH

In preceding sections, we explored how the Louvain algorithm can produce internally-disconnected communities. However, this phenomenon is not exclusive to this algorithms and has been documented in various other community detection algorithms [2, 15, 16, 18, 27]. To mitigate this issue, a commonly employed approach involves splitting disconnected communities as a post-processing step [15, 16, 24, 27, 43], utilizing Breadth First Search (BFS) [15, 27]. We refer to this as Split Last (SL) using BFS, or *SL-BFS*. However, this strategy may exacerbate the problem of poorly connected communities for multi-level community detection techniques, such as the Louvain algorithm [39].

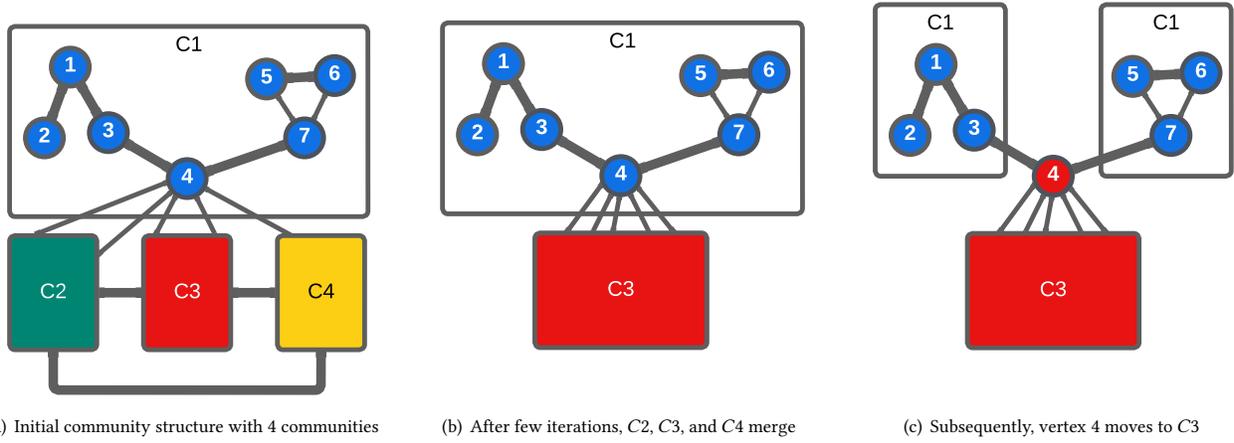
### 4.1 Our Split Pass (SP) approach

To tackle the aforementioned challenges encountered by the Louvain algorithm, we propose to split the disconnected communities in every pass. Specifically, this occurs after the local-moving phase in the Louvain algorithm. We refer to this as the *Split Pass (SP)* approach. Additionally, we explore the conventional approach of splitting disconnected communities as a post-processing step, i.e., after all iterations of the community detection algorithm have been completed and the vertex community memberships have converged. This traditional approach is referred to as *Split Last (SL)*.

In order to partition disconnected communities using either the SP or the SL approach, we explore three distinct techniques: minimum-label-based *Label Propagation (LP)*, minimum-label-based *Label Propagation with Pruning (LPP)*, and *Breadth First Search (BFS)*. The rationale behind investigating LP and LPP techniques for splitting disconnected communities as they are readily parallelizable.

With the LP technique, each vertex in the graph initially receives a unique label (its vertex ID). Subsequently, in each iteration, every vertex selects the minimum label among its neighbors within its assigned community, as determined by the community detection algorithm. This iterative process continues until labels for all vertices converge. Since each vertex obtains a unique label within its connected component and its community, communities comprising multiple connected components get partitioned. In contrast, the LPP technique incorporates a pruning optimization step where only unprocessed vertices are handled. Once a vertex is processed, it is marked as such, and gets reactivated (or marked as unprocessed) if one of its neighbors changes their label. The pseudocode for the LP and LPP techniques is presented in Algorithm 1, with detailed explanations given in Section 4.1.1.

On the other hand, the BFS technique for splitting internally-disconnected communities involves selecting a random vertex within each community and identifying all vertices reachable from it as part of one subcommunity. If any vertices remain unvisited in the original community, another random vertex is chosen from the remaining set, and the process iterates until all vertices within each community are visited. Consequently, the BFS technique facilitates the partitioning of connected components within each community. The pseudocode for the BFS technique for splitting disconnected



**Figure 1: An example demonstrating the possibility of internally disconnected communities with the Louvain algorithm. Here, C1, C2, C3, and C4 are four communities obtained after running a few iterations of the Louvain algorithm, with vertices 1 to 7 being members of community C1. Thick lines are used to denote higher edge weights.**

communities is outlined in Algorithm 2, with its in-depth explanation provided in Section 4.1.2.

**4.1.1 Explanation of LP/LPP algorithm.** We now discuss the pseudocode for the parallel minimum-label-based Label Propagation (LP) and Label Propagation with Pruning (LPP) techniques, given in Algorithm 1, that partition the internally-disconnected communities. These techniques can be employed either as a post-processing step (SL) at the end of the community detection algorithm, or after the refinement/local-moving phase (SP) in each pass. Here, the function `splitDisconnectedLp()` that is responsible for this task, takes as input the graph  $G(V, E)$  and the community memberships  $C$  of vertices, and returns the updated community memberships  $C'$  where all disconnected communities have been separated.

In lines 2-5, the algorithm starts by initializing the minimum labels  $C'$  of each vertex to their respective vertex IDs, and designates all vertices as unprocessed. Lines 6-23 represent the iteration loop of the algorithm. Initially, the number of changed labels  $\Delta N$  is initialized (line 7). This is followed by an iteration of label propagation (lines 8-21), and finally a convergence check (line 23). During each iteration (lines 8-21), unprocessed vertices are processed in parallel. For each unprocessed vertex  $i$ , it is marked as processed if the Label Propagation with Pruning (LPP) technique is utilized. The algorithm proceeds to identify the minimum label  $c'_{min}$  within the community of vertex  $i$  (lines 12-15). This is achieved by iterating over the outgoing neighbors of  $i$  in the graph  $G$  and considering only those neighbors belonging to the same community as  $i$ . The minimum label found among these neighbors, along with the label of vertex  $i$  itself, determines the minimum community label  $c'_{min}$  for  $i$ . If the obtained minimum label  $c'_{min}$  differs from the current minimum label  $C'[i]$  (line 16),  $C'[i]$  is updated,  $\Delta N$  is incremented to reflect the change, and neighboring vertices of  $i$  belonging to the same community as  $i$  are marked as unprocessed to facilitate their reassessment in subsequent iterations. The label propagation loop (lines 6-23) continues until there are no further changes in the

minimum labels. Finally, the updated labels  $C'$ , representing the updated community membership of each vertex with no disconnected communities, are returned in line 24.

**4.1.2 Explanation of BFS algorithm.** Next, we proceed to describe the pseudocode of the parallel Breadth First Search (BFS) technique, as presented in Algorithm 2, devised for the partitioning of disconnected communities. As with LP/LPP techniques, this technique can be applied either as a post-processing step (SL) at the end or after the refinement or local-moving phase (SP) in each pass. Here, the function `splitDisconnectedBfs()` accepts the input graph  $G(V, E)$  and the community membership  $C$  of each vertex, and returns the updated community membership  $C'$  of each vertex where all the disconnected communities have been split.

Initially, in lines 2-4, the flag vector `vis` representing visited vertices is initialized, and the labels  $C'$  for each vertex are set to their corresponding vertex IDs. Subsequently, each thread concurrently processes every vertex  $i$  in the graph  $G$  (lines 6-11). If the community  $c$  of vertex  $i$  is not present in the work-list  $work_t$  of the current thread  $t$ , or if vertex  $i$  has been visited, the thread proceeds to the next iteration (line 8). Conversely, if community  $c$  is in the work-list  $work_t$  of the current thread  $t$  and vertex  $i$  has not been visited, a BFS is performed from vertex  $i$  to explore vertices within the same community. This BFS utilizes lambda functions  $f_{if}$  to selectively execute BFS on vertex  $j$  if it belongs to the same community, and  $f_{do}$  to update the label of visited vertices after each vertex is explored during BFS (line 11). Upon completion of processing all vertices, threads synchronize, and the revised labels  $C'$  – representing the updated community membership of each vertex with no disconnected communities – are returned (line 12). It is pertinent to note that the work-list  $work_t$  for each thread identified by  $t$  is defined as a set encompassing communities  $[t\chi, t(\chi+1)) \cup [T\chi+t\chi, T\chi+t(\chi+1)) \cup \dots$ , where  $\chi$  denotes the chunk size, and  $T$  signifies the total number of threads. In our implementation, a chunk size of  $\chi = 1024$  is employed.

**Algorithm 1** Split disconnected communities using (min) LP.

---

```

1:  $G(V, E)$ : Input graph
2:  $C$ : Initial community membership/label of each vertex
3:  $C'$ : Updated community membership/label of each vertex
4:  $c'_{min}$ : Minimum connected label within the community
5:  $\Delta N$ : Number of changes in labels

1: function SPLITDISCONNECTEDLP( $G, C$ )
2:    $C' \leftarrow \{\}$ 
3:   for all  $i \in V$  in parallel do
4:     Mark  $i$  as unprocessed
5:      $C'[i] = i$ 
6:   loop
7:      $\Delta N \leftarrow 0$ 
8:     for all unprocessed  $i \in V$  in parallel do
9:       if is SL-LPP or SP-LPP then
10:        Mark  $i$  as processed
11:         $\triangleright$  Find minimum community label
12:         $c'_{min} \leftarrow C'[i]$ 
13:        for all  $j \in G.out(i)$  do
14:          if  $C[j] = C[i]$  then
15:             $c'_{min} \leftarrow \min(C'[j], c'_{min})$ 
16:          if  $c'_{min} = C'[i]$  then continue
17:           $\triangleright$  Update community label
18:           $C'[i] \leftarrow c'_{min}$ ;  $\Delta N \leftarrow \Delta N + 1$ 
19:          if is SL-LPP or SP-LPP then
20:            for all  $j \in G.out(i)$  do
21:              Mark  $j$  as unprocessed if  $C[j] = C[i]$ 
22:           $\triangleright$  Converged?
23:          if  $\Delta N = 0$  then break
24:   return  $C'$ 

```

---

**Algorithm 2** Split disconnected communities using BFS.

---

```

1:  $G(V, E)$ : Input graph
2:  $C$ : Initial community membership/label of each vertex
3:  $C'$ : Updated community membership/label of each vertex
4:  $f_{if}$ : Perform BFS to vertex  $j$  if condition satisfied
5:  $f_{do}$ : Perform operation after each vertex is visited
6:  $vis$ : Visited flag for each vertex
7:  $work_t$ : Work-list of current thread

1: function SPLITDISCONNECTEDBFS( $G, C$ )
2:    $C' \leftarrow \{\}$ ;  $vis \leftarrow \{\}$ 
3:   for all  $i \in V$  in parallel do
4:      $C'[i] = i$ 
5:   for all threads do
6:     for all  $i \in V$  do
7:        $c' \leftarrow C'[i]$ 
8:       if  $c \notin work_t$  or  $vis[i]$  then continue
9:        $f_{if} \leftarrow (j) \implies C[j] = C[i]$ 
10:       $f_{do} \leftarrow (j) \implies C'[j] \leftarrow c'$ 
11:       $bfsVisitForEach(vis, G, i, f_{if}, f_{do})$ 
12:   return  $C'$ 

```

---

Figure 2 illustrates an example of the BFS technique. Initially, Figure 2(a) displays two communities,  $C_1$  and  $C_2$ , derived after the local-moving phase. Here,  $C_1$  has become internally disconnected due to the inclusion of vertex 4 into community  $C_2$  – similar to the case depicted in Figure 1(c). Subsequently, employing the BFS technique, a thread selects a random vertex within community  $C_1$ , such as 2, and designates all vertices reachable within  $C_1$  from 2 with the label of 2, and marking them as visited (Figure 2(b)). Following this, as depicted in Figure 2(c), the same thread picks an unvisited vertex randomly within community  $C_1$ , for example, 7, and labels all vertices reachable within  $C_1$  from 7 with the label of 7, and marking them as visited. An analogous process is executed within community  $C_2$ . Consequently, all vertices are visited, and the labels assigned to them denote the updated community membership of each vertex with no disconnected communities. Note that each thread has a mutually exclusive work-list, ensuring that two threads do not simultaneously perform BFS within the same community.

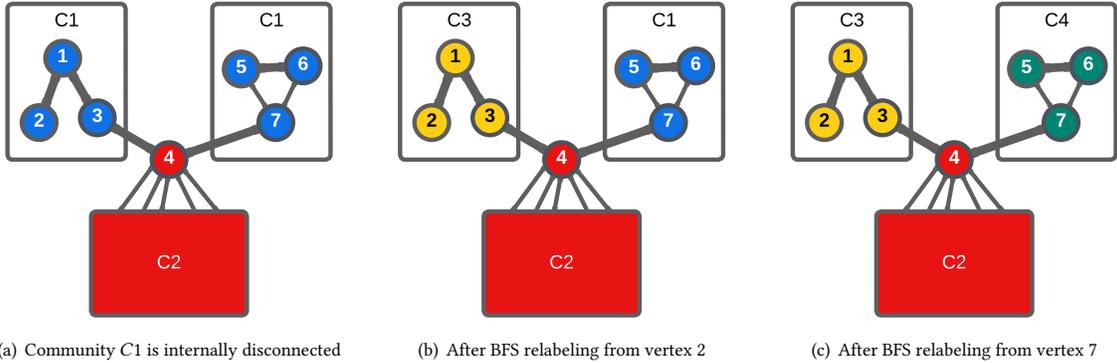
## 4.2 Our GSP-Louvain algorithm

To assess both our Split Pass (SP) approach and the conventional Split Last (SL) approach, utilizing minimum-label-based Label Propagation (LP), minimum-label-based Label Propagation with Pruning (LPP), and Breadth First Search (BFS) techniques for splitting disconnected communities with the Louvain algorithm, we use GVE-Louvain [33], our parallel implementation of Louvain algorithm.

*4.2.1 Determining suitable technique for splitting disconnected communities.* We now determine the optimal technique for partitioning internally-disconnected communities using GVE-Louvain. To achieve this, we investigate both the SL and SP approaches, employing LP, LPP, and BFS techniques. Figures 3(a), 3(b), and 3(c) illustrate the mean relative runtime, modularity, and fractions of disconnected communities for SL-LP, SL-LPP, SL-BFS, SP-LP, SP-LPP, SP-BFS, and the default (i.e., not splitting disconnected communities) approaches. As depicted in Figure 3(c), both SL and SP approaches result in non-disconnected communities. Additionally, Figure 3(b) reveals that the modularity of communities obtained through the SP approach surpasses that of the SL approach while maintaining proximity to the default approach. Finally, Figure 3(a) illustrates that SP-BFS, specifically the SP approach employing the BFS technique, demonstrates superior performance. Consequently, employing BFS to split disconnected communities in each pass (SP) of the Louvain algorithm emerges as the preferred choice.

*4.2.2 Explanation of the algorithm.* We refer to GVE-Louvain, which employs the Split Pass (SP) approach with Breadth-First Search (SP-BFS) technique to handle disconnected communities, as *GSP-Louvain*. The core procedure of GSP-Louvain, encapsulated in the `louvain()` function, is given in Algorithm 3. It consists of the following main steps: initialization, local-moving phase, splitting phase, and the aggregation phase. This function takes as input a graph  $G(V, E)$  and outputs the community membership  $C$  for each vertex in the graph, with none of the returned communities being internally disconnected.

First, in line 2, the community membership  $C$  is initialized for each vertex in  $G$ , and the algorithm conducts passes of the Louvain algorithm, limited to a maximum number of passes defined by



**Figure 2: An example illustrating the BFS technique for splitting internally-disconnected communities. Initially, two communities, C1 and C2, are shown, with C1 being internally disconnected due to vertex 4 joining C2. The BFS technique selects random vertices within each community and labels reachable vertices with the same label, indicated with a new community ID.**

*MAX\_PASSES*. During each pass, various metrics such as the total edge weight of each vertex  $K'$ , the total edge weight of each community  $\Sigma'$ , and the community membership  $C'$  of each vertex in the current graph  $G'$  are updated. Subsequently, in line 6, the local-moving phase is executed by invoking `louvainMove()` (Algorithm 4), which optimizes the community assignments. Following this, the algorithm proceeds to the splitting phase, where the internally disconnected communities in  $C'$  are separated. This is done using the parallel BFS technique, in line 7, with the `splitCommunitiesBfs()` function (Algorithm 2). Next, in line 8, global convergence is inferred if the local-moving phase converges in a single iteration. If so, we terminate the passes. Additionally, if there is minimal reduction in the number of communities ( $|\Gamma|$ ), indicating diminishing returns, the current pass is halted (line 10).

If convergence is not achieved, the algorithm proceeds with the following steps: renumbering communities (line 11), updating top-level community memberships  $C$  using dendrogram lookup (line 12), executing the aggregation phase via `louvainAggregate()` (Algorithm 5), and adjusting the convergence threshold for subsequent passes, known as threshold scaling (line 14). The subsequent pass initiates at line 3. Upon completion of all passes, a final update of the top-level community memberships  $C$  using dendrogram lookup occurs (line 15), followed by the return of the top-level community membership  $C$  of each vertex in graph  $G$ .

## 5 EVALUATION

### 5.1 Experimental Setup

**5.1.1 System used.** We utilize a server comprising two Intel Xeon Gold 6226R processors, with each processor housing 16 cores operating at 2.90 GHz. Each core is equipped with a 1 MB L1 cache, a 16 MB L2 cache, and a shared L3 cache of 22 MB. The system is configured with 376 GB of RAM and is running CentOS Stream 8.

**5.1.2 Configuration.** We employ 32-bit integers to represent vertex IDs and 32-bit floats for edge weights, while computations and hashtable values utilize 64-bit floats. We utilize 64 threads to match the number of available cores on the system, unless stated otherwise. Compilation is performed using GCC 8.5 and OpenMP 4.5.

**Algorithm 3** GSP-Louvain: Our Parallel Louvain algorithm which identifies communities that are not internally disconnected.

---

$\triangleright G(V, E)$ : Input graph  
 $\square G'(V', E')$ : Input/super-vertex graph  
 $\square C$ : Community membership of each vertex  
 $\square C'$ : Community membership of each vertex in  $G'$   
 $\square K'$ : Total edge weight of each vertex  
 $\square \Sigma'$ : Total edge weight of each community  
 $\square l_i$ : Number of iterations performed (per pass)  
 $\square l_p$ : Number of passes performed  
 $\square \tau, \tau_{agg}$ : Iteration, Aggregation tolerance

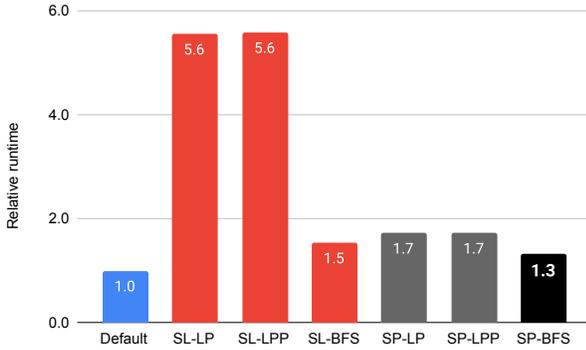
---

```

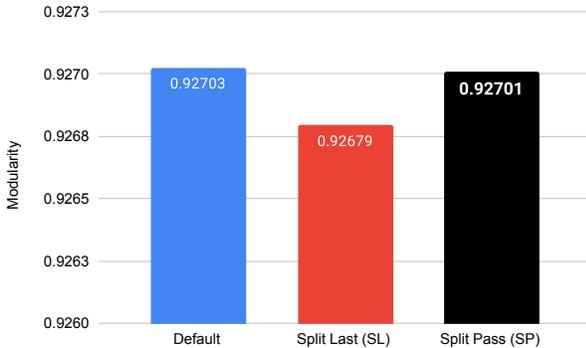
1: function LOUVAIN( $G$ )
2:   Vertex membership:  $C \leftarrow [0..|V|]$ ;  $G' \leftarrow G$ 
3:   for all  $l_p \in [0..MAX\_PASSES)$  do
4:      $\Sigma' \leftarrow K' \leftarrow vertexWeights(G')$ ;  $C' \leftarrow [0..|V'|]$ 
5:     Mark all vertices in  $G'$  as unprocessed
6:      $l_i \leftarrow louvainMove(G', C', K', \Sigma', \tau)$   $\triangleright$  Algorithm 4
7:      $C' \leftarrow splitDisconnectedBfs(G', C')$   $\triangleright$  Algorithm 2
8:     if  $l_i \leq 1$  then break  $\triangleright$  Globally converged?
9:      $|\Gamma|, |\Gamma_{old}| \leftarrow$  Number of communities in  $C, C'$ 
10:    if  $|\Gamma|/|\Gamma_{old}| > \tau_{agg}$  then break  $\triangleright$  Low shrink?
11:     $C' \leftarrow$  Renumber communities in  $C'$ 
12:     $C \leftarrow$  Lookup dendrogram using  $C$  to  $C'$ 
13:     $G' \leftarrow louvainAggregate(G', C')$   $\triangleright$  Algorithm 5
14:     $\tau \leftarrow \tau/TOLERANCE\_DROP$   $\triangleright$  Threshold scaling
15:     $C \leftarrow$  Lookup dendrogram using  $C$  to  $C'$ 
16:  return  $C$ 
  
```

---

**5.1.3 Dataset.** The graphs utilized in our experiments are listed in Table 1, sourced from the SuiteSparse Matrix Collection [20]. These graphs exhibit 3.07 to 214 million vertices, and 25.4 million to 3.80 billion edges. We ensure that the edges are undirected and weighted, with a default weight of 1.



(a) Relative runtime using different approaches for splitting disconnected communities with Parallel Louvain algorithm



(b) Modularity using different approaches for splitting disconnected communities with Parallel Louvain algorithm



(c) Fraction of disconnected communities (logarithmic scale) using different approaches for splitting disconnected communities with Parallel Louvain algorithm

**Figure 3: Mean relative runtime, modularity, and fraction of disconnected communities (log-scale) using *Split Last (SL)* and *Split Pass (SP)* approaches for splitting disconnected communities with Parallel Louvain algorithm [32] across all graphs in the dataset. Both *SL* and *SP* approaches employ *Label Propagation (LP)*, *Label Propagation with Pruning (LPP)*, or *Breadth First Search (BFS)* techniques for splitting.**

**Table 1: List of 13 graphs obtained SuiteSparse Matrix Collection [20] (directed graphs are marked with \*). Here,  $|V|$  is the number of vertices,  $|E|$  is the number of edges (after adding reverse edges),  $D_{avg}$  is the average degree, and  $|\Gamma|$  is the number of communities obtained with *GSP-Louvain*.**

Graph	$ V $	$ E $	$D_{avg}$	$ \Gamma $
<b>Web Graphs (LAW)</b>				
indochina-2004*	7.41M	341M	41.0	4.28K
uk-2002*	18.5M	567M	16.1	42.8K
arabic-2005*	22.7M	1.21B	28.2	3.58K
uk-2005*	39.5M	1.73B	23.7	20.4K
webbase-2001*	118M	1.89B	8.6	2.77M
it-2004*	41.3M	2.19B	27.9	5.10K
sk-2005*	50.6M	3.80B	38.5	3.86K
<b>Social Networks (SNAP)</b>				
com-LiveJournal	4.00M	69.4M	17.4	4.47K
com-Orkut	3.07M	234M	76.2	43
<b>Road Networks (DIMACS10)</b>				
asia_osm	12.0M	25.4M	2.1	2.50K
europa_osm	50.9M	108M	2.1	3.36K
<b>Protein k-mer Graphs (GenBank)</b>				
kmer_A2a	171M	361M	2.1	19.8K
kmer_V1r	214M	465M	2.2	7.37K

## 5.2 Performance Comparison

We now compare the performance of GSP-Louvain with the original Leiden [39], igraph Leiden [9], and NetworkKit Leiden [36]. For the original Leiden, we employ a C++ program to initialize a `ModularityVertexPartition` upon the loaded graph and invoke `optimise_partition()` to determine the community membership of each vertex. On graphs with high edge counts, such as *webbase-2001* and *sk-2005*, utilizing `ModularityVertexPartition` can result in disconnected communities due to numerical precision issues [37]. Despite a positive improvement in separating disconnected parts, the large total edge weight of the graph can render it effectively near zero. For such graphs, we employ `RBConfigurationVertexPartition`, as it uses unscaled modularity improvements, avoiding the occurrence of disconnected communities. In the case of igraph Leiden, we utilize `igraph_community_leiden()` with a resolution of  $1/2|E|$ , a beta value of 0.01, and specify that the algorithm to run until convergence. For NetworkKit Leiden, we create a Python script to call `ParallelLeiden()`, while constraining the number of passes to 10. For each graph, we measure the runtime of each implementation and the modularity of the resulting communities five times to obtain an average. Additionally, we store the community membership vector (for each vertex in the graph) in a file and subsequently determine the number of disconnected components using Algorithm 6. Throughout these evaluations, we optimize for modularity as the quality function.

Figure 4(a) presents the runtimes of the original Leiden, igraph Leiden, NetworkKit Leiden, and GSP-Louvain on each graph in the dataset. On the *sk-2005* graph, GSP-Louvain identifies communities in 11.6 seconds, achieving a processing rate of 328 million edges/s.

Figure 4(b) illustrates the speedup of GSP-Louvain relative to the original Leiden, igraph Leiden, and NetworKit Leiden. On average, GSP-Louvain exhibit speedups of 341 $\times$ , 83 $\times$ , and 6.1 $\times$ , respectively. Figure 4(c) displays the modularity of communities obtained using each implementation. On average, GSP-Louvain achieves modularity values that are 0.3% lower than those obtained by the original Leiden and igraph Leiden, but 25% higher than those obtained by NetworKit Leiden (particularly noticeable on road networks and protein k-mer graphs). Finally, Figure 4(d) illustrates the fraction of disconnected communities obtained by each implementation. The absence of bars indicates the absence of disconnected communities. Communities identified by the original Leiden, igraph Leiden, and GSP-Louvain exhibit no disconnected communities. However, on average, NetworKit Leiden exhibits fractions of disconnected communities amounting to  $1.5 \times 10^{-2}$ , particularly noticeable on web graphs and social networks. This is likely due to an error in its implementation. Thus, GSP-Louvain effectively tackles the issue of disconnected communities, while being significantly faster than existing alternatives, and attaining similar modularity scores. Figure 7 depicts the comparison of GVE-Louvain and GSP-Louvain. This comparison is explained in detail in Section A.3.

### 5.3 Performance Analysis

We proceed to analyze the performance of GSP-Louvain. The phase-wise and pass-wise split of GSP-Louvain is depicted in Figures 5(a) and 5(b). Figure 5(a) illustrates that GSP-Louvain devotes a considerable portion of its runtime to the local-moving phase on web graphs, road networks, and protein k-mer graphs, while it predominantly focuses on the aggregation phase on social networks, and on the splitting phase on road networks. The pass-wise breakdown for GSP-Louvain, shown in Figure 5(b), indicates that the initial pass is computationally intensive for high-degree graphs (such as web graphs and social networks), while subsequent passes take precedence in terms of execution time on low-degree graphs (such as road networks and protein k-mer graphs).

On average, GSP-Louvain dedicates 37% of its runtime to the local-moving phase, 21% to the splitting phase, 29% to the aggregation phase, and 13% to other steps (including initialization, renumbering communities, dendrogram lookup, and resetting communities). Additionally, the first pass of GSP-Louvain accounts for 71% of the total runtime. This initial pass in GSP-Louvain is computationally demanding due to the size of the original graph (subsequent passes operate on super-vertex graphs).

### 5.4 Strong Scaling

Finally, we evaluate the strong scaling performance of GSP-Louvain by varying the number of threads from 1 to 64 (in multiples of 2) for each input graph. We measure the total time taken for GSP-Louvain to identify communities, with its respective phase splits, including local-moving, splitting, aggregation, and other associated steps. The strong scaling of GSP-Louvain is illustrated in Figure 6.

With 32 threads, GSP-Louvain achieve an average speedup of 8.4 $\times$ , compared to single-threaded execution. This indicates a performance increase of 1.5 $\times$  for every doubling of threads. The scalability is limited due to the sequential nature of steps/phases in the

algorithm, as well as the lower scalability of splitting and aggregation phases. At 64 threads, GSP-Louvain is impacted by NUMA effects, resulting in speedups of only 9.4 $\times$ .

## 6 CONCLUSION

In this study, we proposed GSP-Louvain, another approach to mitigate the issue of disconnected communities with the Louvain algorithm. Utilizing a system featuring two 16-core Intel Xeon Gold 6226R processors, we demonstrated that GSP-Louvain not only rectifies this issue but also achieves an impressive processing rate of 328M edges/s on a 3.8B edge graph. Comparatively, it surpasses the original Leiden, igraph Leiden, and NetworKit Leiden by 341 $\times$ , 83 $\times$ , and 6.1 $\times$ , respectively. Moreover, the identified communities exhibit similar quality to the first two implementations and are 25% higher in quality than those produced by NetworKit. Additionally, GSP-Louvain exhibits a performance improvement rate of 1.5 $\times$  for every doubling of threads.

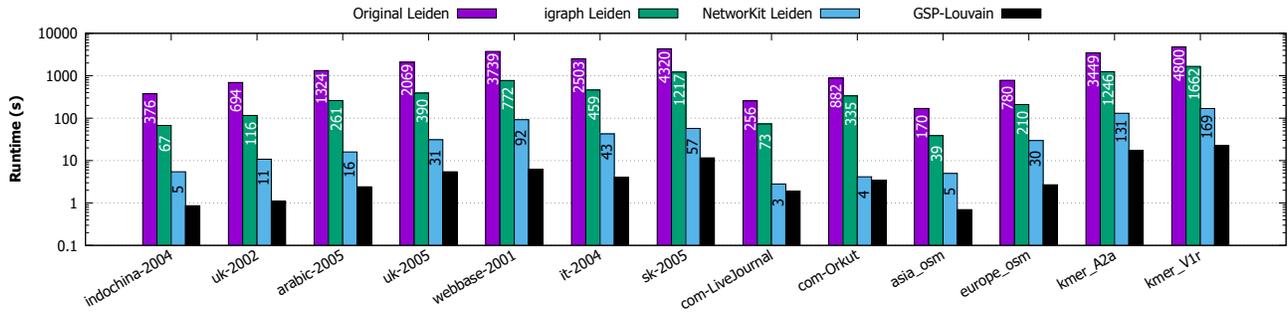
In this version of this report, we addressed issues in measuring disconnected communities for the original Leiden and igraph Leiden, which arose due to the number of vertices in a graph varying between the Matrix Market and the Edgelist formats (which does not have isolated vertices), and used the `RBConfigurationVertexPartition` with the original Leiden for large graphs (i.e., *webbase-2001* and *sk-2005*). Further, we removed any discussion on GSP-Leiden (it was similar to GSP-Louvain, but adapted to the Leiden algorithm).

## ACKNOWLEDGMENTS

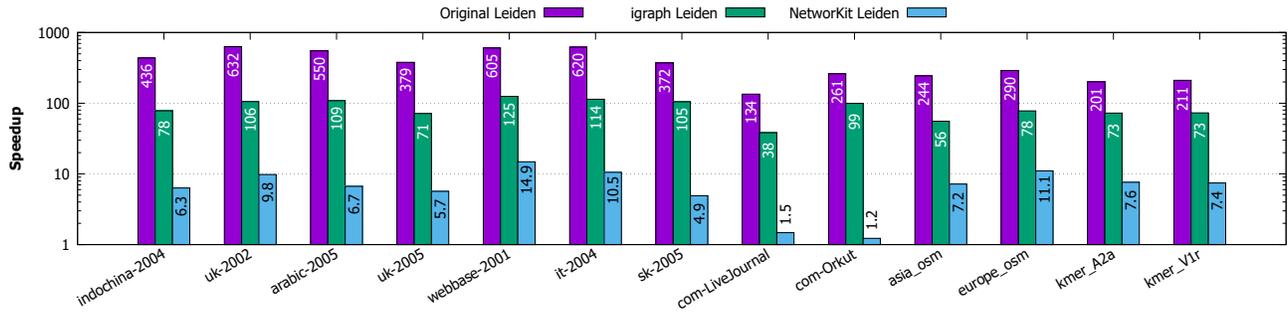
I would like to thank Prof. Kishore Kothapalli, Prof. Dip Sankar Banerjee, and Vincent Traag for their support.

## REFERENCES

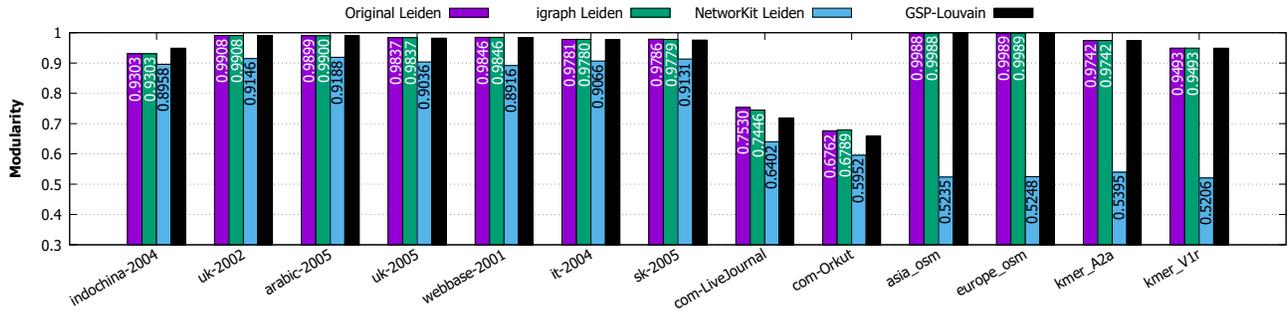
- [1] Emmanuel Abbe. 2018. Community detection and stochastic block models: recent developments. *Journal of Machine Learning Research* 18, 177 (2018), 1–86.
- [2] B. Ball, B. Karrer, and M. E. Newman. 2011. Efficient and principled method for detecting communities in networks. *Physical Review E* 84, 3 (2011), 036103.
- [3] V. Blondel, J. Guillaume, R. Lambiotte, and E. Lefebvre. 2008. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment* 2008, 10 (Oct 2008), P10008.
- [4] U. Brandes, D. Dellinger, M. Gaertler, R. Gorke, M. Hoefer, Z. Nikoloski, and D. Wagner. 2007. On modularity clustering. *IEEE transactions on knowledge and data engineering* 20, 2 (2007), 172–188.
- [5] B. Chatterjee and H. Saha. 2019. Detection of communities in large scale networks. In *IEEE 10th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)*. IEEE, 1051–1060.
- [6] C. Cheong, H. Huynh, D. Lo, and R. Goh. 2013. Hierarchical Parallel Algorithm for Modularity-Based Community Detection Using GPUs. In *Proceedings of the 19th International Conference on Parallel Processing (Aachen, Germany) (Euro-Par'13)*. Springer-Verlag, Berlin, Heidelberg, 775–787.
- [7] Aaron Clauset, Mark EJ Newman, and Christopher Moore. 2004. Finding community structure in very large networks. *Physical review E* 70, 6 (2004), 066111.
- [8] Michele Coscia, Fosca Giannotti, and Dino Pedreschi. 2011. A classification for community discovery methods in complex networks. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 4, 5 (2011), 512–546.
- [9] G. Csardi, T. Nepusz, et al. 2006. The igraph software package for complex network research. *InterJournal, complex systems* 1695, 5 (2006), 1–9.
- [10] Jordi Duch and Alex Arenas. 2005. Community detection in complex networks using extremal optimization. *Physical review E* 72, 2 (2005), 027104.
- [11] M. Fazlali, E. Moradi, and H. Malazi. 2017. Adaptive parallel Louvain community detection on a multicore platform. *Microprocessors and microsystems* 54 (Oct 2017), 26–34.
- [12] S. Fortunato. 2010. Community detection in graphs. *Physics reports* 486, 3-5 (2010), 75–174.
- [13] Santo Fortunato and Marc Barthélemy. 2007. Resolution limit in community detection. *Proceedings of the national academy of sciences* 104, 1 (2007), 36–41.



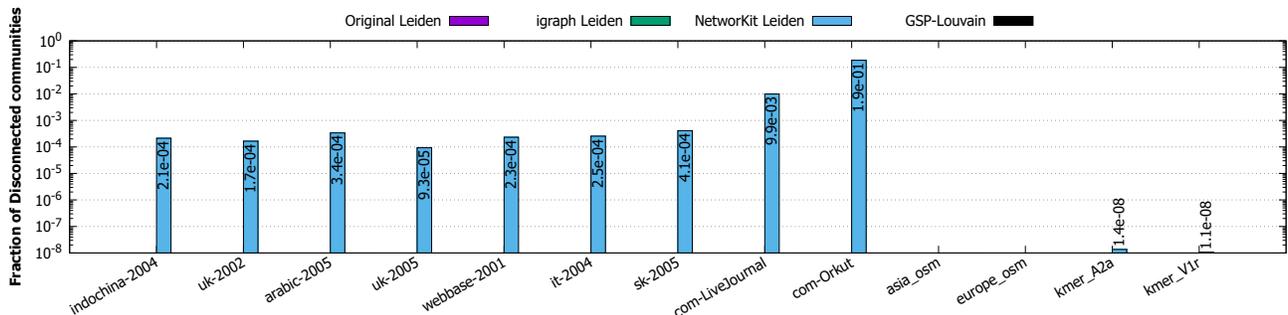
(a) Runtime in seconds (logarithmic scale) with *Original Leiden*, *igraph Leiden*, *NetworkKit Leiden*, and *GSP-Louvain*



(b) Speedup of *GSP-Louvain* (logarithmic scale) with respect to *Original Leiden*, *igraph Leiden*, and *NetworkKit Leiden*.



(c) Modularity of communities obtained with *Original Leiden*, *igraph Leiden*, *NetworkKit Leiden*, and *GSP-Louvain*.



(d) Fraction of disconnected communities (logarithmic scale) with *Original Leiden*, *igraph Leiden*, *NetworkKit Leiden*, and *GSP-Louvain*.

**Figure 4: Runtime in seconds (log-scale), speedup (log-scale), modularity, and fraction of disconnected communities (log-scale) with *Original Leiden*, *igraph Leiden*, *NetworkKit Leiden*, and *GSP-Louvain* for each graph in the dataset.**

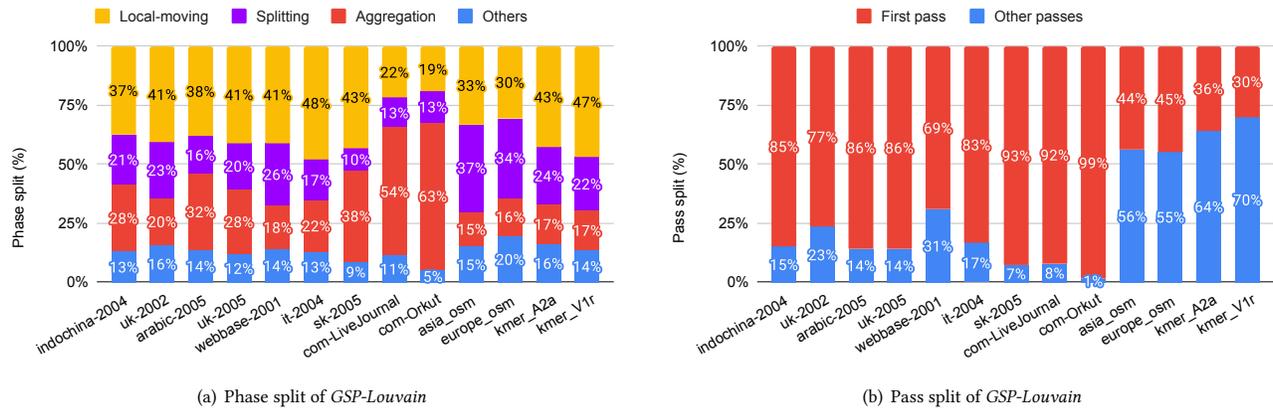


Figure 5: Phase split of GSP-Louvain shown on the left, and pass split shown on the right for each graph in the dataset.

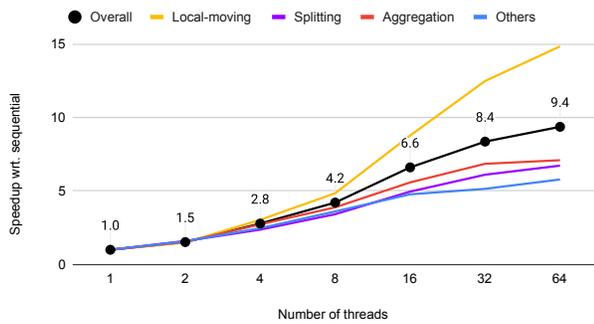


Figure 6: Overall speedup of GSP-Louvain, and its various phases (local-moving, splitting, aggregation, and others), with increasing number of threads (in multiples of 2).

[14] S. Ghosh, M. Halappanavar, A. Tumeo, A. Kalyanaraman, and A.H. Gebremedhin. 2018. Scalable distributed memory community detection using vite. In *2018 IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, 1–7.

[15] S. Gregory. 2010. Finding overlapping communities in networks by label propagation. *New Journal of Physics* 12 (10 2010), 103018. Issue 10.

[16] Ahmed Ibrahim Hafez, Aboul Ella Hassanien, and Aly A Fahmy. 2014. BNEM: a fast community detection algorithm using generative models. *Social Network Analysis and Mining* 4 (2014), 1–20.

[17] M. Halappanavar, H. Lu, A. Kalyanaraman, and A. Tumeo. 2017. Scalable static and dynamic community detection using Grappolo. In *IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, Waltham, MA USA, 1–6.

[18] Sajjad Hesamipour, Mohammad Ali Balafar, Saeed Mousazadeh, et al. 2022. Detecting communities in complex networks using an adaptive genetic algorithm and node similarity-based encoding. *Complexity* 2023 (2022).

[19] K. Kloster and D. Gleich. 2014. Heat kernel based community detection. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, New York, USA, 1386–1395.

[20] S. Kolodziej, M. Aznavah, M. Bullock, J. David, T. Davis, M. Henderson, Y. Hu, and R. Sandstrom. 2019. The SuiteSparse matrix collection website interface. *JOSS* 4, 35 (2019), 1244.

[21] A. Lancichinetti and S. Fortunato. 2009. Community detection algorithms: a comparative analysis. *Physical Review E, Statistical, Nonlinear, and Soft Matter Physics* 80, 5 Pt 2 (Nov 2009), 056117.

[22] J. Leskovec. 2021. CS224W: Machine Learning with Graphs | 2021 | Lecture 13.3 - Louvain Algorithm. <https://www.youtube.com/watch?v=0zuILBOlcsw>

[23] H. Lu, M. Halappanavar, and A. Kalyanaraman. 2015. Parallel heuristics for scalable community detection. *Parallel computing* 47 (Aug 2015), 19–37.

[24] Malte Luecken. 2016. *Application of multi-resolution partitioning of interaction networks to the study of complex disease*. Ph.D. Dissertation. University of Oxford.

[25] Mark EJ Newman. 2006. Modularity and community structure in networks. *Proceedings of the national academy of sciences* 103, 23 (2006), 8577–8582.

[26] Fabian Nguyen. [n. d.]. *Leiden-Based Parallel Community Detection*. Bachelor’s Thesis. Karlsruhe Institute of Technology, 2021 (zitiert auf S. 31).

[27] U. Raghavan, R. Albert, and S. Kumara. 2007. Near linear time algorithm to detect community structures in large-scale networks. *Physical Review E* 76, 3 (Sep 2007), 036106–1–036106–11.

[28] Jörg Reichardt and Stefan Bornholdt. 2006. Statistical mechanics of community detection. *Physical review E* 74, 1 (2006), 016110.

[29] M. Rosvall and C. Bergstrom. 2008. Maps of random walks on complex networks reveal community structure. *Proceedings of the national academy of sciences* 105, 4 (2008), 1118–1123.

[30] R. Rotta and A. Noack. 2011. Multilevel local search algorithms for modularity clustering. *Journal of Experimental Algorithmics (JEA)* 16 (2011), 2–1.

[31] S. Ryu and D. Kim. 2016. Quick community detection of big graph data using modified louvain algorithm. In *IEEE 18th International Conference on High Performance Computing and Communications (HPCC)*. IEEE, Sydney, NSW, 1442–1445.

[32] Subhajit Sahu. 2023. GVE-Leiden: Fast Leiden Algorithm for Community Detection in Shared Memory Setting. *arXiv preprint arXiv:2312.13936* (2023).

[33] Subhajit Sahu. 2023. GVE-Louvain: Fast Louvain Algorithm for Community Detection in Shared Memory Setting. *arXiv preprint arXiv:2312.04876* (2023).

[34] J. Shi, L. Dhulipala, D. Eisenstat, J. Łącki, and V. Mirrokni. 2021. Scalable community detection via parallel correlation clustering.

[35] S. Souravlas, A. Sifaleras, M. Tsintogianni, and S. Katsavounis. 2021. A classification of community detection methods in social networks: a survey. *International journal of general systems* 50, 1 (Jan 2021), 63–91.

[36] C.L. Staudt, A. Sazonovs, and H. Meyerhenke. 2016. NetworKit: A tool suite for large-scale complex network analysis. *Network Science* 4, 4 (2016), 508–530.

[37] V. Traag. 2024. Personal communication. (2024).

[38] V.A. Traag and L. Šubelj. 2023. Large network community detection by fast label propagation. *Scientific Reports* 13, 1 (2023), 2701.

[39] V. Traag, L. Waltman, and N. Eck. 2019. From Louvain to Leiden: guaranteeing well-connected communities. *Scientific Reports* 9, 1 (Mar 2019), 5233.

[40] L. Waltman and N. Eck. 2013. A smart local moving algorithm for large-scale modularity-based community detection. *The European physical journal B* 86, 11 (2013), 1–14.

[41] J. Whang, D. Gleich, and I. Dhillon. 2013. Overlapping community detection using seed set expansion. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*. 2099–2108.

[42] C. Wickramaarachchi, M. Frincu, P. Small, and V. Prasanna. 2014. Fast parallel algorithm for unfolding of communities in large graphs. In *IEEE High Performance Extreme Computing Conference (HPEC)*. IEEE, Waltham, MA USA, 1–6.

[43] F Alexander Wolf, Fiona K Hamey, Mireya Plass, Jordi Solana, Joakim S Dahlin, Berthold Göttgens, Nikolaus Rajewsky, Lukas Simon, and Fabian J Theis. 2019. PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome biology* 20 (2019), 1–9.

[44] J. Xie, B. Szymanski, and X. Liu. 2011. SLPA: Uncovering overlapping communities in social networks via a speaker-listener interaction dynamic process. In *IEEE 11th International Conference on Data Mining Workshops*. IEEE, IEEE, Vancouver, Canada, 344–349.

- [45] Zhao Yang, René Algesheimer, and Claudio J Tessone. 2016. A comparative analysis of community detection algorithms on artificial networks. *Scientific reports* 6, 1 (2016), 30750.
- [46] X. You, Y. Ma, and Z. Liu. 2020. A three-stage algorithm on community detection in social networks. *Knowledge-Based Systems* 187 (2020), 104822.
- [47] N. Zarayeneh and A. Kalyanaraman. 2021. Delta-Screening: A Fast and Efficient Technique to Update Communities in Dynamic Graphs. *IEEE transactions on network science and engineering* 8, 2 (Apr 2021), 1614–1629.
- [48] J. Zeng and H. Yu. 2015. Parallel Modularity-Based Community Detection on Large-Scale Graphs. In *IEEE International Conference on Cluster Computing*. 1–10.

## A APPENDIX

### A.1 Phases of GSP-Louvain

Here, we explain the local-moving, and aggregation phases of GSP-Louvain. For details on the splitting phase, please consult Section 4.1.2.

*A.1.1 Local-moving phase of GSP-Louvain.* The pseudocode for the local-moving phase of GSP-Louvain is outlined in Algorithm 4. Within this algorithm, vertices are iteratively moved between communities to maximize modularity. The `louvainMove()` function takes as input the current graph  $G'$ , community membership  $C'$ , total edge weight of each vertex  $K'$  and each community  $\Sigma'$ , and the iteration tolerance  $\tau$ . It returns the number of iterations performed  $l_i$ .

Lines 3-15 encapsulate the primary loop of the local-moving phase. Initially, all vertices are designated as unprocessed (line 2). Subsequently, in line 4, the total delta-modularity per iteration  $\Delta Q$  is initialized. Next, parallel iteration over unprocessed vertices is conducted (lines 5-14). For each unprocessed vertex  $i$ ,  $i$  is flagged as processed, i.e., vertex pruning (line 6), followed by scanning communities connected to  $i$ , excluding itself (line 7). Further, the best community  $c^*$  for moving  $i$  to is determined (line 9), and the delta-modularity  $\delta Q^*$  of moving  $i$  to  $c^*$  is computed (line 10). If a superior community is identified (with  $\delta Q^* > 0$ ), the community membership of  $i$  is updated (lines 12-13), and its neighbors are marked as unprocessed (line 14). If not,  $i$  stays in its original community. Line 15 verifies if the local-moving phase has converged, terminating the loop if so (or if `MAX_ITERATIONS` is reached). Finally, in line 16, the number of iterations performed  $l_i$  is returned.

*A.1.2 Aggregation phase of GSP-Louvain.* Finally, we provide the pseudocode for the aggregation phase in Algorithm 5, which aggregates communities into super-vertices in preparation for the subsequent pass of the Louvain algorithm, operating on the super-vertex graph. The `louvainAggregate()` function accepts the current graph  $G'$  and the community membership  $C'$  as input and returns the super-vertex graph  $G''$ .

In lines 3-4, the offsets array for the community vertices Compressed Sparse Row (CSR)  $G'C'.offsets$  is computed. Initially, this involves determining the number of vertices in each community using `countCommunityVertices()` and subsequently performing an exclusive scan on the array. Then, in lines 5-6, a parallel iteration over all vertices is conducted to atomically populate vertices belonging to each community into the community graph CSR  $G'C'$ . Following this, the offsets array for the super-vertex graph CSR is determined by estimating the degree of each super-vertex. This process includes calculating the total degree of each community with `communityTotalDegree()` and performing an exclusive scan on the array (lines 8-9). As a result, the super-vertex graph CSR exhibits sparsity, with gaps between the edges and weights arrays of each super-vertex in the CSR.

Following that, in lines 11-17, a parallel iteration over all communities  $c \in [0, |\Gamma|)$  is executed. For each vertex  $i$  belonging to community  $c$ , all communities  $d$  (along with associated edge weight  $w$ ) linked to  $i$ , as defined by `scanCommunities()` in Algorithm 4, are included in the per-thread hashtable  $H_t$ . Once  $H_t$  is populated

---

#### Algorithm 4 Local-moving phase of GSP-Louvain [33].

---

```

▷  $G'(V', E')$ : Input/super-vertex graph
▷  $C'$ : Community membership of each vertex
▷  $K'$ : Total edge weight of each vertex
▷  $\Sigma'$ : Total edge weight of each community
▷  $\tau$ : Iteration tolerance
□  $H_t$ : Collision-free per-thread hashtable
□  $l_i$ : Number of iterations performed

1: function LOUVAINMOVE( $G', C', K', \Sigma', \tau$ )
2:   Mark all vertices in  $G'$  as unprocessed
3:   for all  $l_i \in [0..MAX\_ITERATIONS)$  do
4:     Total delta-modularity per iteration:  $\Delta Q \leftarrow 0$ 
5:     for all unprocessed  $i \in V'$  in parallel do
6:       Mark  $i$  as processed (prune)
7:        $H_t \leftarrow scanCommunities(\{i\}, G', C', i, false)$ 
8:       ▷ Use  $H_t, K', \Sigma'$  to choose best community
9:        $c^* \leftarrow$  Best community linked to  $i$  in  $G'$ 
10:       $\delta Q^* \leftarrow$  Delta-modularity of moving  $i$  to  $c^*$ 
11:      if  $c^* = C'[i]$  then continue
12:       $\Sigma'[C'[i]] - = K'[i]; \Sigma'[c^*] + = K'[i]$  atomic
13:       $C'[i] \leftarrow c^*; \Delta Q \leftarrow \Delta Q + \delta Q^*$ 
14:      Mark neighbors of  $i$  as unprocessed
15:      if  $\Delta Q \leq \tau$  then break           ▷ Locally converged?
16:   return  $l_i$ 

17: function SCANCOMMUNITIES( $H_t, G', C', i, self$ )
18:   for all  $(j, w) \in G'.edges(i)$  do
19:     if not  $self$  and  $i = j$  then continue
20:      $H_t[C'[j]] \leftarrow H_t[C'[j]] + w$ 
21:   return  $H_t$ 

```

---

with all communities (and their associated weights) linked to community  $c$ , these are atomically added as edges to super-vertex  $c$  in the super-vertex graph  $G''$ . Finally, in line 18, the super-vertex graph  $G''$  is returned.

### A.2 Finding disconnected communities

We introduce our parallel algorithm designed to identify disconnected communities, given the original graph and the community membership of each vertex. The core principle involves assessing the size of each community, selecting a representative vertex from each community, navigating within the community from that vertex while avoiding neighboring communities, and designating a community as disconnected if all its vertices are unreachable. We investigate four distinct approaches, distinguished by their utilization of parallel Depth-First Search (DFS) or Breadth-First Search (BFS), and whether per-thread or shared *visited* flags are employed. When shared visited flags are utilized, each thread scans all vertices but exclusively processes its designated community based on the community ID. Our findings reveal that employing parallel BFS traversal with a shared flag vector yields the most efficient results. Given the deterministic nature of this algorithm, all approaches yield identical outcomes. Algorithm 6 outlines the pseudocode for

**Algorithm 5** Aggregation phase of GSP-Louvain [33].

---

```

1:  $G'(V', E')$ : Input/super-vertex graph
2:  $C'$ : Community membership of each vertex
3:  $G'_{C'}(V'_{C'}, E'_{C'})$ : Community vertices (CSR)
4:  $G''(V'', E'')$ : Super-vertex graph (weighted CSR)
5:  $*.offsets$ : Offsets array of a CSR graph
6:  $H_t$ : Collision-free per-thread hashtable

1: function LOUVAINAGGREGATE( $G', C'$ )
2:   ▷ Obtain vertices belonging to each community
3:    $G'_{C'}.offsets \leftarrow countCommunityVertices(G', C')$ 
4:    $G'_{C'}.offsets \leftarrow exclusiveScan(G'_{C'}.offsets)$ 
5:   for all  $i \in V'$  in parallel do
6:     Add edge ( $C'[i], i$ ) to CSR  $G'_{C'}$  atomically
7:   ▷ Obtain super-vertex graph
8:    $G''.offsets \leftarrow communityTotalDegree(G', C')$ 
9:    $G''.offsets \leftarrow exclusiveScan(G''.offsets)$ 
10:   $|\Gamma| \leftarrow$  Number of communities in  $C'$ 
11:  for all  $c \in [0, |\Gamma|)$  in parallel do
12:    if degree of  $c$  in  $G'_{C'} = 0$  then continue
13:     $H_t \leftarrow \{\}$ 
14:    for all  $i \in G'_{C'}.edges(c)$  do
15:       $H_t \leftarrow scanCommunities(H, G', C', i, true)$ 
16:    for all  $(d, w) \in H_t$  do
17:      Add edge  $(c, d, w)$  to CSR  $G''$  atomically
18:  return  $G''$ 

```

---

this approach. Here, the `disconnectedCommunities()` function takes the input graph  $G$  and the community membership  $C$  as input and returns the disconnected flag  $D$  for each community.

Let us now delve into Algorithm 6. Initially, in line 2, we initialize the disconnected community flag  $D$  and the visited vertices flags  $vis$ . Line 3 computes the size of each community  $S$  in parallel using the `communitySizes()` function. Subsequently, each thread processes each vertex  $i$  in the graph  $G$  in parallel (lines 5-14). In line 6, we determine the community membership of  $i$  ( $c$ ), and set the count of vertices reached from  $i$  to 0. If community  $c$  is either empty or not in the work-list of the current thread  $work_t$ , the thread proceeds to the next iteration (line 9). However, if community  $c$  is non-empty and in the work-list of the current thread  $work_t$ , we perform BFS from vertex  $i$  to explore vertices in the same community. This utilizes lambda functions  $f_{if}$  to conditionally execute BFS to vertex  $j$  if it belongs to the same community, and  $f_{do}$  to update the count of reached vertices after each vertex is visited during BFS (line 12). If the number of vertices  $reached$  during BFS is less than the community size  $S[c]$ , we mark community  $c$  as disconnected (line 13). Finally, we update the size of the community  $S[c]$  to 0, indicating that the community has been processed (line 14). It's important to note that the work-list  $work_t$  for each thread with ID  $t$  is defined as a set containing communities  $[t\chi, t(\chi+1)) \cup [T\chi+t\chi, T\chi+t(\chi+1)) \cup \dots$ , where  $\chi$  is the chunk size, and  $T$  is the number of threads. In our implementation, we utilize a chunk size of  $\chi = 1024$ .

**Algorithm 6** Finding disconnected communities in parallel [32].

---

```

1:  $G(V, E)$ : Input graph
2:  $C$ : Community membership of each vertex
3:  $D$ : Disconnected flag for each community
4:  $S$ : Size of each community
5:  $f_{if}$ : Perform BFS to vertex  $j$  if condition satisfied
6:  $f_{do}$ : Perform operation after each vertex is visited
7:  $reached$ : Number of vertices reachable from  $i$  to  $i$ 's community
8:  $vis$ : Visited flag for each vertex
9:  $work_t$ : Work-list of current thread

1: function DISCONNECTEDCOMMUNITIES( $G, C$ )
2:   $D \leftarrow \{\}$ ;  $vis \leftarrow \{\}$ 
3:   $S \leftarrow communitySizes(G, C)$ 
4:  for all threads in parallel do
5:    for all  $i \in V$  do
6:       $c \leftarrow C[i]$ ;  $reached \leftarrow 0$ 
7:      ▷ Skip if community  $c$  is empty, or
8:      ▷ does not belong to work-list of current thread.
9:      if  $S[c] = 0$  or  $c \notin work_t$  then continue
10:      $f_{if} \leftarrow (j) \implies C[j] = c$ 
11:      $f_{do} \leftarrow (j) \implies reached \leftarrow reached + 1$ 
12:      $bfsVisitForEach(vis, G, i, f_{if}, f_{do})$ 
13:     if  $reached < S[c]$  then  $D[c] \leftarrow 1$ 
14:      $S[c] \leftarrow 0$ 
15:  return  $D$ 

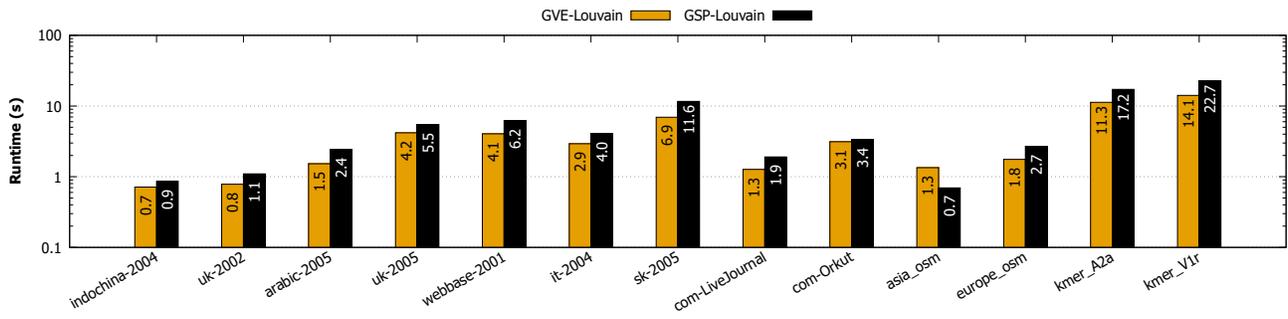
```

---

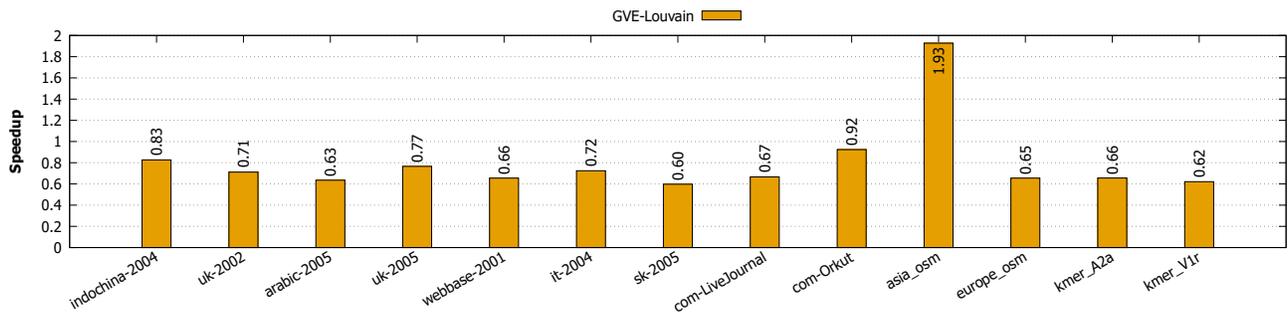
**A.3 Additional Performance comparison**

We proceed to compare the performance of GSP-Louvain with GVE-Louvain [33]. Similar to our previous approach, we execute each algorithm five times for every graph in the dataset to mitigate measurement noise and report the averages in Figures 7(a), 7(b), 7(c), and 7(d).

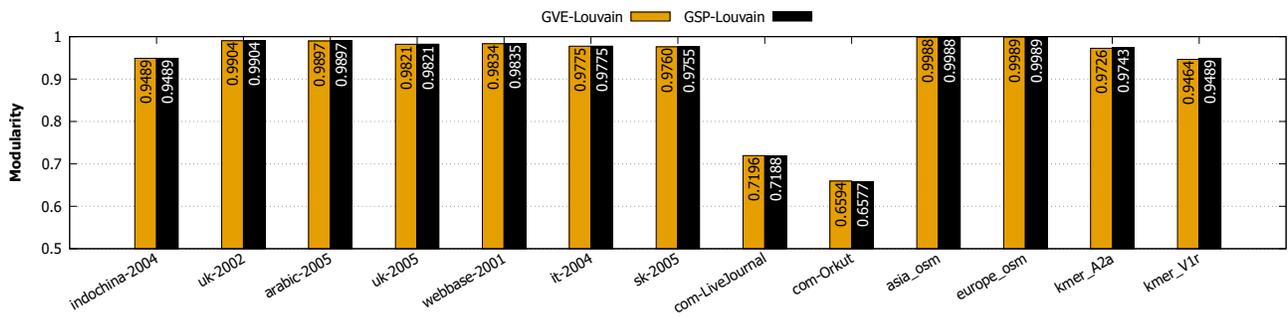
Figure 7(a) illustrates the runtimes of GSP-Louvain and GVE-Louvain on each graph in the dataset. On average, GSP-Louvain exhibits about a 33% increase in runtime compared to GVE-Louvain. This additional computational time is a compromise made to ensure the absence of internally disconnected communities. Figure 7(c) presents the modularity of communities obtained by each implementation. On average, the modularity of communities obtained using GSP-Louvain and GVE-Louvain remains roughly identical. Lastly, Figure 7(d) displays the fraction of internally disconnected communities identified by each implementation. Communities obtained with GSP-Louvain exhibit no disconnected communities, whereas communities identified with GVE-Louvain feature on average 3.9% disconnected communities.



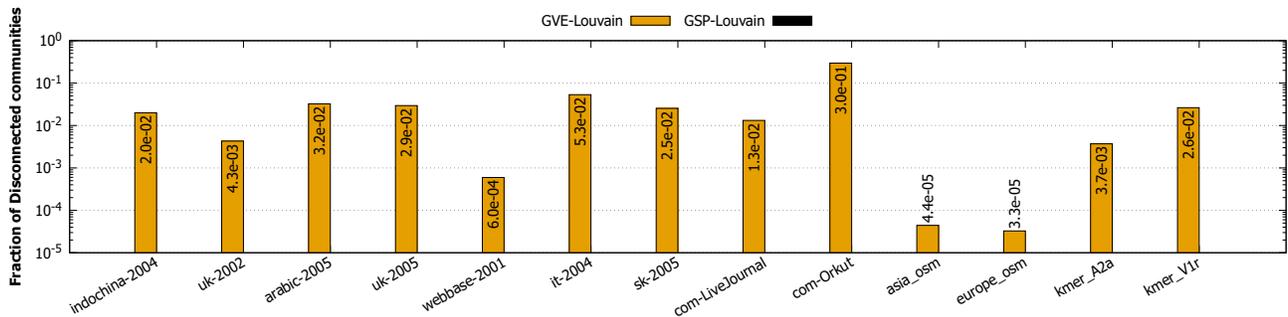
(a) Runtime in seconds (logarithmic scale) with *GVE-Louvain* and *GSP-Louvain*



(b) Speedup of *GSP-Louvain* with respect to *GVE-Louvain*.



(c) Modularity of communities obtained with *GVE-Louvain* and *GSP-Louvain*.



(d) Fraction of disconnected communities (logarithmic scale) with *GVE-Louvain* and *GSP-Louvain*.

**Figure 7: Runtime in seconds (log-scale), speedup, modularity, and fraction of disconnected communities (log-scale) with *GVE-Louvain* and *GSP-Louvain* for each graph in the dataset.**