

# COMP5112/MSBD5009 Parallel Programming

## Assignment 2: Pthread Programming

Due on 8 April 2021 at 5:00 pm

### Instructions

- This assignment counts for 10 points.
- This is an individual assignment. You can discuss with others and search online resources, but your submission should be your own code.
- Add your name, student id and email at the first line of comments in your submission.
- Your submission will be compiled and tested on Azure machines through remote terminals.
- Submit your assignment through Canvas before the deadline.
- **No late submission will be accepted!**

### Assignment Description

Graph structural clustering is a common data analysis task to cluster vertices by their edge connections in the graph. [SCAN \(Structural Clustering Algorithm for Networks\)](#) [1] is such an algorithm that clusters vertices based on a structural similarity measure. The algorithm is efficient in both computation and memory.

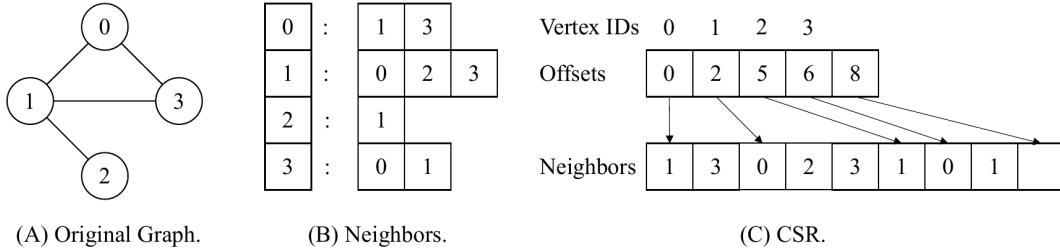


Figure 1: Illustration of data structure.

In our implementation, a graph  $G = (V, E)$  is stored in a CSR structure as shown in Fig. 1: where one array **Neighbors** stores the neighbors of all vertices, and the other array **Offsets** stores the offset of each vertex into the **Neighbors** array.

Please note that: (1) The length of the array **Offsets** is  $|V| + 1$ , with the last element of this array points to the end of **Neighbors**. (2) Accordingly, the length of **Neighbors** is the number of edges plus 1 ( $|E| + 1$ ).

The SCAN algorithm, shown in Algorithm 1, can be divided into two stages: (1) find pivot vertices using the structural similarity measure, (2) expand clusters starting from the pivot vertices in the depth-first order.

**Algorithm 1:** The SCAN Algorithm

**Input** : a graph  $G = (V, E)$ , the structural similarity threshold  $\epsilon$ , and the minimal cluster size  $\mu$

**Output:** the total number of clusters in  $G$ , and the cluster ID of each vertex

```

1 main()
2   read files and initialization
   // Stage 1: find pivot vertices each of which has at least  $\mu$  neighbors
   // whose similarity with the vertex exceeds  $\epsilon$ 
3   foreach  $v \in V$  do
4       foreach  $w \in \text{Neighbors}(v)$  do
5            $\Gamma(v) \leftarrow \text{Neighbors}(v) \cup \{v\}$ ,  $\Gamma(w) \leftarrow \text{Neighbors}(w) \cup \{w\}$ 
6           if  $\text{similarity}(v, w) \leftarrow \frac{|\Gamma(v) \cap \Gamma(w)|}{\sqrt{|\Gamma(v)| \times |\Gamma(w)|}} > \epsilon$  then
7                $\text{Neighbors}_\epsilon(v) \leftarrow \text{Neighbors}_\epsilon(v) \cup \{w\}$ 
8           end
9       end
10      if  $|\text{Neighbors}_\epsilon(v)| > \mu$  then
11           $\text{pivots}[v] \leftarrow \text{true}$ 
12      else
13      end
   // Stage 2: expand clusters from pivot vertices
14  foreach  $v \in V$  do
15      if  $\text{pivots}[v]$  and  $\text{!visited}[v]$  then
16           $\text{visited}[v] \leftarrow \text{true}$ 
17           $\text{cluster\_result}[v] \leftarrow v$ 
18           $\text{expansion}(v, v)$ 
19           $\text{num\_clusters} \leftarrow \text{num\_clusters} + 1$ 
20      else
21      end
22  output  $\text{num\_clusters}$  and  $\text{cluster\_result}$ 
23 expansion( $v$ , label)
24  foreach  $w \in \text{Neighbors}_\epsilon(v)$  do
25      if  $\text{pivots}[w]$  and  $\text{!visited}[w]$  then
26           $\text{visited}[w] \leftarrow \text{true}$ 
27           $\text{cluster\_result}[w] \leftarrow \text{label}$ 
28           $\text{expansion}(w, \text{label})$ 
29      end
30  end

```

## Input and output

In this assignment, you will implement a **pthread** version of the SCAN algorithm.

In the assignment folder:

- `clustering_sequential.cpp` is the sequential version for your reference. You can check the clustering result and compare its running time.
- `clustering.h` provides some utility functions.
- `main.cpp` contains the main function of the pthread version.
- `clustering_pthread_skeleton.cpp` is the code skeleton for your work. **Your task is to complete** the following function which returns the array `cluster_result`:  

```
int *scan(float epsilon, int mu, int num_threads, int num_vs, int num_es, int *nbr_offs, int *nbrs);
```

- The `datasets` and `results` folders contain the datasets and clustering results respectively.

Here is the output of executing the algorithm on `./sequential ../../data/test1/1.txt 0.7 3`:

**Screen output:**

```
Elapsed Time: 0.000016516 s
Number of clusters: 2
```

**Result file output (sequential.txt):**

```
2
-1 -1 -1 -1 -1 -1 6 6 -1 -1 -1 -1 -1 13 -1 13 -1 -1
```

2 in the first line is the number of clusters. The second line 6 13 is the cluster IDs (-1 if not in any cluster) of all vertices in order. Please note that in the parallel setting, we set the cluster ID be the lowest vertex ID of the pivots in the cluster. Make sure your file output follows the same format.

## Submission

1. You only need to complete and submit the `clustering_pthread_skeleton.cpp` to Canvas. Make sure your name information is added as comments in the first line. You can add or adjust any helper functions and variables as you wish in the skeleton file, but keep the other files unchanged.
2. We will use different input data (`num_vs`  $\leq 1,200,000$  and `num_es`  $\leq 50,000,000$ ) and specify different numbers of threads ( $0 < \text{num\_threads} \leq 8$ ) to test your program.
3. The correctness, running time and speedup of your program will be considered in grading.
4. We will perform code similarity checks. In case a submission has code similarity issues, we may request clarification and deduct partial marks or full marks on a case-by-case basis.

## References

- [1] Xu, Xiaowei, et al. "Scan: a structural clustering algorithm for networks." Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining. 2007.