



UNIVERSITY OF THE PHILIPPINES

BACHELOR OF SCIENCE IN COMPUTER SCIENCE

ALIYA AHLANNA C. MIRANDA

**W-CPMd:  
A Fast Clique Percolation Method for Directed  
Networks**

Adviser:  
HENRY ADORNA, PH.D.  
Department of Computer Science  
University of the Philippines

May 2020

Thesis Classification:  
**F**  
This Undergraduate Research Paper is available to the public.

16 *“I hereby grant the University of the Philippines a non-exclusive, worldwide,*  
17 *royalty-free license to reproduce, publish and publicly distribute copies of this un-*  
18 *dergraduate research paper in whatever form subject to the provisions of applicable*  
19 *laws, the provisions of the UP IPR policy and any contractual obligations, as well*  
20 *as more specific permission marking on the Title Page.”*

21

22 *“Specifically, I grant the following rights to the University:*

- 23 *a) To upload a copy of the work in the theses database of the college/school/institute/*  
24 *department and in any other databases available on the public internet;*
- 25 *b) To publish the work in the college/school/institute/department journal, both*  
26 *in print and electronic or digital format and online; and*
- 27 *c) To give open access to above-mentioned work, thus allowing “fair use” of the*  
28 *work in accordance with the provisions of the Intellectual Property Code of*  
29 *the Philippines (Republic Act No. 8293), especially for teaching, scholarly*  
30 *and research purposes.”*

31

ALIYA AHLANNA C. MIRANDA

32

Date of Submission: May 2020

33  
34  
35  
36

Department of Computer Science  
College of Engineering  
University of the Philippines  
Diliman, Quezon City

37

## ENDORSEMENT

38 This undergraduate research paper hereto attached, entitled **W-CPMd:**  
39 **A Fast Clique Percolation Method for Directed Networks**, prepared and  
40 submitted by **Aliya Ahlanna C. Miranda**, in partial fulfillment of the require-  
41 ment for the degree of **Bachelor of Science in Computer Science**, is hereby  
42 accepted.

43

HENRY ADORNA, PH.D.  
Adviser

44 This undergraduate research paper is hereby officially accepted and approved as  
45 partial fulfillment of the requirements for the degree of **Bachelor of Science in**  
46 **Computer Science**.

47

JAN MICHAEL YAP, PH.D.  
Chair  
Department of Computer Science

48 I hereby declare that I have created this work completely on I own and used no  
49 other sources or tools than the ones listed, and that I have marked any citations  
50 accordingly.  
51

52 ALIYA AHLANNA C. MIRANDA  
53 Diliman, Quezon City  
54 May 2020

# Abstract

55 **W-CPMd:**  
56 **A Fast Clique Percolation Method for Directed Networks**

57 Aliya Ahlanna C. Miranda                      Adviser:  
University of the Philippines, 2020              Henry Adorna, Ph.D.

58        Community detection algorithms are necessary for finding the relationships  
59 within a network. Real-time networks contains large amounts of data with a de-  
60 fined structure. W-CPM is a good algorithm for getting overlapping communities  
61 in large-scale networks due to its fast computational process but it is limited to  
62 undirected networks.

63

64        In this study, an algorithm that uses similar concepts to W-CPM but applicable  
65 to directed networks is proposed, called W-CPMd. If successful, this will help in  
66 distinguishing different communities in a vast variety of networks.

# Table of Contents

67	<b>Abstract</b> . . . . .	<b>v</b>
68	<b>Acknowledgments</b> . . . . .	<b>vii</b>
69	<b>List of Tables</b> . . . . .	<b>viii</b>
70	<b>List of Figures</b> . . . . .	<b>ix</b>
71	<b>1 Introduction</b> . . . . .	<b>1</b>
72	<b>2 Theoretical Framework and Related Literature</b> . . . . .	<b>2</b>
73	2.1 Theoretical Framework . . . . .	2
74	2.2 Related Work . . . . .	3
75	<b>3 Problem Statement</b> . . . . .	<b>5</b>
76	<b>4 Objectives of the Study</b> . . . . .	<b>6</b>
77	<b>5 Methodology</b> . . . . .	<b>7</b>
78	5.1 Directed Weak Cliques . . . . .	7
79	5.2 Node Priority . . . . .	7
80	5.3 Sequential Algorithm . . . . .	8
81	<b>6 Testing</b> . . . . .	<b>9</b>
82	6.1 Normalized Mutual Information(NMI)[7] . . . . .	9
83	<b>7 Results and Discussion</b> . . . . .	<b>10</b>
84	<b>8 Results and Discussion</b> . . . . .	<b>11</b>
85	<b>9 Conclusion</b> . . . . .	<b>12</b>
86	<b>10 Recommendations</b> . . . . .	<b>13</b>
87	<b>List of References</b> . . . . .	<b>14</b>
88	<b>11 Appendix</b> . . . . .	<b>15</b>

# Acknowledgments

<sup>89</sup> To be written.

# List of Tables

90	5.1	Pseudocode: Setting the Node Priority . . . . .	7
91	5.2	Pseudocode: Getting the Referral Count of a node . . . . .	8



# List of Figures

92	2.1	The bottom three triangles show little information circulating while	
93		the upper four contains 3-cycles . . . . .	2

# Chapter 1

## Introduction

Community detection is the process of getting groups of interacting vertices, known as clusters or communities[4]. These clusters share common properties and specific roles in a network. By knowing these similarities, services can be improved like efficient storage data and navigation of queries[3]. It is also good for determining the structure of networks, like hierarchies. It is useful for getting target audiences or market in social networks and e-commerce businesses[4].

There are different types of communities. Disjoint communities are clusters that does not have common vertices, while overlapping communities have nodes sharing with multiple communities[4]. There are also directed and undirected networks, where in directed, relationships between vertices are not reciprocal, while it does in undirected. Creating algorithms for directed networks are a difficult task since it is characterized by asymmetrical matrices, thus its analysis is more complex[3].

In real networks, vertices commonly belong to more than one group, and highly organized with a hierarchical structure[3]. Clique percolation method, or CPM, is created by the concept that every vertex is connected to other vertex and forming a clique[4]. It is focused on undirected graphs by finding its k-cliques. This is extended for directed graphs in CPMd by formulating directed k-cliques[9]. Since CPM is an NP-Complete problem, W-CPM is created to lessen the processing time to polynomial time[10]. W-CPM is limited to undirected graphs thus creating a similar algorithm that works for directed graphs will open more possible networks to be classified to their communities.

## Chapter 2

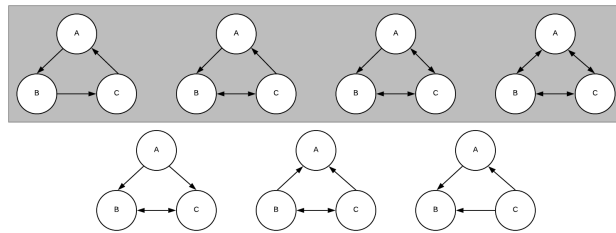
# Theoretical Framework and Related Literature

## 2.1 Theoretical Framework

There is no universally accepted definition for a community. A common description for one is that there are more edges inside the community than the rest of the network. A community should have a connectedness inside it such that there is a pair between each vertex and running only inside the cluster. The intra-cluster density, which are the edges within a cluster, should be larger than the average link density, while the inter-cluster density, which are the edges going to the rest of the graph, is much smaller.

The basic composition of a clique is a triangle. The amount of information propagated is dependent on the direction of the edges. There are seven types of triangles. The ones with a 3-cycle provides the most information since it circulates on all three vertices[5].

There are four types of criterion for communities based on social media analysis - complete mutuality, reachability, vertex degree, and comparison of internal and external cohesion. Small subgroups consists mostly of maximal subgraphs or cliques. But in large networks, larger cliques are not frequent and hierarchies



**Figure 2.1:** The bottom three triangles show little information circulating while the upper four contains 3-cycles

138 exist. Since a community does not consist of a single large clique, one can define  
 139 subgroups as clique-like subgraphs. In reachability, one can define a community  
 140 based on the existence and lengths of paths between vertices. Another criteria is  
 141 vertex degree, where it roots from the idea that a vertex is adjacent to a mini-  
 142 mum number of other vertices in the subgraph. And finally, the comparison of  
 143 internal and external cohesion, where a strong community is defined based on its  
 144 intra-cluster and inter-cluster densities such that the former is larger than the  
 145 latter[3].

146 CPM is based on complete mutuality criterion of a community in its strictest  
 147 sense[8]. CPMd redefined a clique such that it takes in consideration the hierar-  
 148 chical structures of networks[9]. W-CPM added a cohesion property by merging  
 149 triangular cliques by the number of links connected between them and gets the  
 150 priority of a node based its vertex degree[10].

## 151 2.2 Related Work

152 There are a lot of different algorithms in identifying communities. CPM is rooted  
 153 to the concept that real-life communities have cliques inside them. To be more  
 154 precise, CPM defines a k-clique community as a maximal complete subgraph.  
 155 These cliques are first located, then the communities are identified using a standard  
 156 component analysis of the clique-clique overlap matrix[8].

157 Real-time network relationships are not necessarily reciprocal. Thus, G. Palla  
 158 extended his work on directed networks called CPMd. Here, the k-cliques are  
 159 replaced by directed k-cliques, defined as complete sub-graphs of size k where  
 160 it can be any pair of nodes such that there is a directed link pointing from a  
 161 higher order towards a lower one. For these directed k-cliques, double-links are  
 162 not allowed, no directed loops and the restricted out degree of each vertex in the  
 163 clique should differ[9].

164 Communities do not necessarily form complete cliques. In SCP, a link is cre-  
 165 ated between two nodes to for small k-cliques based on the amount of common  
 166 neighbors two nodes have and creates an temporary link for these two nodes to  
 167 create them.[6]. An algorithm proposed by X. Zhang (WCPM) also uses the con-  
 168 cept of weak cliques to determine communities, which is based on the smallest  
 169 cliques formed - triangles.[10]

170     Computing for the identification of  $k$ -cliques and grouping for communities is  
171     an NP-complete problem. To solve this, WCPM only gets the weak cliques from  
172     two adjacent nodes and these cliques are then merged by priority and similarity  
173     to form communities. This solved CPM's problem of getting restricted to  $k$ -clique  
174     communities while having a computational difficulty of  $O(dm)$ [10]. J. Kumpula  
175     had a similar algorithm that searches for nodes with at least  $(k-1)$ -degrees to get  
176      $k$ -cliques and uses merges these cliques in a sequential manner to minimize the  
177     time taken in determining the communities.[6]

## 178 Chapter 3

# 179 Problem Statement

180 The study aims to create a directed CPM using the concepts provided in W-CPM  
181 and have it run in polynomial time. This will aid in classifying large and real-time  
182 networks in less processing time compared to CPMd.

183 This study aims to answer the following questions:

- 184 • Is there a way to create a directed CPM using W-CPM concepts?
- 185 • Can an algorithm be made such that it is comparable in accuracy and  
186 speed on current popular directed network algorithms?

## 187 Chapter 4

# 188 Objectives of the Study

189 This study will focus only on finding alternatives on weak cliques and merging  
190 mechanisms that will work on directed networks.

191 The main objectives of this study are as follows:

- 192 • Create a directed network method using the main concepts of W-CPM
- 193 • Have the algorithm run in polynomial time

## 194 Chapter 5

# 195 Methodology

## 196 5.1 Directed Weak Cliques

197 Definition 1 (Directed Weak Cliques Determined by an Adjacent Source Node and  
198 Drain Node): Given a directed network  $G = (V, E)$ , where  $V$  is the set of nodes  
199 and  $E$  is a set of links. Let  $u$  and  $v$  be two adjacent nodes in  $G$  with  $u$  having a  
200 minimum of one out-degree and  $v$  having a minimum of one in-degree. A directed  
201 weak clique determined by  $u$  and  $v$  is defined as

$$G_{uv} = (V_{uv}, E_{uv})$$

202 where  $V_{uv} = \{(u, v) \cup (N_u \cap N_v)\}$ ,  $E_{uv} = \{(x \rightarrow y) \in E | x, y \in V_{uv}\}$ ,  $N_u = \{x | (x \rightarrow$   
203  $u) \in E\}$ ,  $N_v = \{x | (v \rightarrow x) \in E\}$ .

## 204 5.2 Node Priority

205 Definition 2 (Node Strength[2]): Given a weighted network  $G = (V, E)$ , where  $V$   
206 is the set of nodes and  $E$  is a set of links. Let  $u$  and  $v$  be two adjacent nodes  
207 in  $G$ . The weight of edge  $e_{uv}$  is  $w_{uv}$  where  $w_{uv} = 0$  when nodes  $u$  and  $v$  are not  
208 connected by an edge. The node strength  $k_u$  is defined as

$$k_u = \sum_{v \in V} w_{uv}$$

209

210 Definition 3 (Common Neighbors[1]):

$$s_{xy} = |N_u \cap N_v|$$

Input: Graph $G$ , Node $n$
Output: number
if $G$ is directed, return number of incoming edges of $n$
if $G$ is weighted, return sum of weights of incoming edges of $n$

**Table 5.1:** Pseudocode: Setting the Node Priority



Input: Graph G, Node u, v
Output: number
$in_u \leftarrow$ incoming nodes of u
$out_v \leftarrow$ outgoing nodes of v
return $ in_u \cap out_v $

**Table 5.2:** Pseudocode: Getting the Referral Count of a node

211 Definition 3.5 (Referral Count):

$$s_{xy} = \text{outgoing link nodes } N_u \cup \text{incoming link nodes } N_v$$

212

### 213 5.3 Sequential Algorithm

214 Definition 4 (SCP Phase II[6]): Given a set of k-cliques, each component (k-1)-  
 215 cliques are extracted. If the (k-1)-cliques is a component of a different k-clique,  
 216 then merge. Repeat for the entire set.

217 Definition 4.5 (Directed Weak Clique Merge): Given a set of directed weak cliques,  
 218 its directed links are extracted. If the link is a component of a different weak clique,  
 219 then merge. Repeat for the entire set.

## 220 Chapter 6

## 221 Testing

### 222 6.1 Normalized Mutual Information(NMI)[7]

## <sup>223</sup> Chapter 7

## <sup>224</sup> Results and Discussion

## <sup>225</sup> Chapter 8

# <sup>226</sup> Results and Discussion

<sup>227</sup> **Chapter 9**

<sup>228</sup> **Conclusion**

<sup>229</sup> **Chapter 10**

<sup>230</sup> **Recommendations**

# List of References

- [1] BOJANOWSKI, M., AND CHROL, B. Proximity-based methods for link prediction in graphs with r package linkprediction. <https://cran.r-project.org/>, 2018. [Online; accessed 27 Feb 2020].
- [2] CHEN, D., SHANG, M., LV, Z., AND FU, Y. Detecting overlapping communities of weighted networks via a local algorithm. *Physica A: Statistical Mechanics and Its Applications* (2010).
- [3] FORTUNATO, S. Community detection in graphs. *Elsevier Physics Reports* (2009).
- [4] JAVED, M. A., YOUNIS, M. S., LATIF, S., QADIR, J., AND BAIG, A. Community detection in networks: A multidisciplinary review. *Journal of Network and Computer Applications* (2018).
- [5] KLYMKO, C., GLEICH, D., AND KOLDA, T. G. Using triangles to improve community detection in directed networks.
- [6] KUMPULA, J. M., KIVELA, M., KASKI, K., AND SARMAKI, J. Sequential algorithm for fast clique percolation. *Physical Review* (2008).
- [7] MCDAID, A. F., GREENE, D., AND HURLEY, N. Normalized mutual information to evaluate overlapping community finding algorithms.
- [8] PALLA, G., FARKA, I. J., DERENYI, I., AND VICSEK, T. Uncovering the overlapping community structure of complex networks in nature and society. *Nature Letters* 435 (2005).
- [9] PALLA, G., FARKA, I. J., POLLNER, , DERENYI, I., AND VICSEK, T. Directed network modules. *New Journal of Physics* (2007).
- [10] ZHANG, X., WANG, C., SU, Y., PAN, L., AND ZHANG, H.-F. A fast overlapping community detection algorithm based on weak cliques for large-scale networks. *IEEE Transactions on Computational Social Systems* 4 (December 2017).

<sup>257</sup> **Chapter 11**

<sup>258</sup> **Appendix**