

Cohesive Subgraph Detection on Large Graphs

by

Bohua Yang

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF
THE REQUIREMENTS FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

Australian Artificial Intelligence Institute (AII)
Faculty of Engineering and Information Technology (FEIT)
University of Technology Sydney (UTS)

April, 2021

CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Bohua Yang declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Signature: Production Note:
Signature removed prior to publication.

Date: 09/04/2021

ACKNOWLEDGEMENTS

First of all, I would like to deliver my sincere gratitude to my principal supervisor A/Prof. Lu Qin. Thanks for his selfless help during my PhD life. He is a good mentor and friend to me. He is very professional; he is a perfect guide who leads me to the beautiful world of the graph. He is very patient; he is always ready to help me whenever I meet problems in research or life. I could not finish my thesis without his illuminating instructions.

Secondly, I would like to thank my co-supervisor Prof. Ying Zhang. He gives me much accurate and valuable advice on my research topics. Besides, he often encourages and guides me to optimize my career development plan.

Thirdly, I would like to thank Dr. Dong Wen. All the research works presented in this thesis are conducted together with him. He patiently guides me in acquiring research skills, such as thinking about problems, writing papers, and preparing presentations. I learn a lot from the collaboration with him.

Fourthly, I would like to thank Prof. Xuemin Lin and Dr. Lijun Chang for supporting the works in this thesis. I would like to thank Prof. Lin for offering an interesting but rigorous research environment. I would like to thank Dr. Chang for his insightful advice on my research works.

Fifthly, I am very grateful to A/Prof. Xin Zhao, my mentor at the Renmin

University of China, for stimulating my research interest and providing helpful advice for my academic career. Furthermore, I would like to thank Prof. Jirong Wen, A/Prof. Hui Sun, Prof. Qing Zhu, and A/Prof. Yan Yu at the Renmin University of China for their kind help.

I would also like to appreciate Prof. Jeffrey Xu Yu, Prof. Wenjie Zhang, Prof. Ronghua Li, A/Prof. Ling Chen, Dr. Yulei Sui, Dr. Xin Cao, A/Prof. Zengfeng Huang, Dr. Shiyu Yang, A/Prof. Yixiang Fang, Prof. Weiwei Liu, and Prof. Xiaoyang Wang for sharing brilliant thoughts and experiences. Thanks to Dr. Dian Ouyang, Dr. Wentao Li, Dr. Conggai Li, Dr. Mingjie Li, Dr. Adi Lin, Dr. Yuan Liang, Mr. Yuxuan Qiu, Mr. Junhua Zhang, Dr. Fan Zhang, Dr. Wanqi Liu, Mr. Hanchen Wang, Mr. Yilun Huang, Dr. Xubo Wang, Dr. Longbin Lai, Dr. Long Yuan, Dr. Xing Feng, Dr. Haida Zhang, Dr. Wei Li, Dr. Chen Zhang, Dr. Yang Yang, Dr. Kai Wang, Dr. You Peng, Dr. Boge Liu, Mr. Yuanhang Yu, Mr. Peilun Yang, Ms. Xiaoshuang Chen, Mr. Xuefeng Chen, Mr. Yuren Mao, Mr. Yixing Yang, Mr. Zhengyi Yang, Mr. Qingyuan Linghu, Mr. Michael Ruisi Yu, Mr. Chenji Huang, Mr. Yu Hao, Mr. Kongzhang Hao, Dr. Peng Zhang, Mr. Xunxiang Yao, Dr. Binbin Huang, Dr. Yarui Chen, Ms. Kun Wang, Mr. Yijun Li, and Ms. Ruilin Liu. The days we spent together bear unforgettable memories.

Finally, I would like to thank my grandfather Mr. Zhencai Yang, and my grandmother Ms. Yuling Zhao, who give me selfless love and endless encouragement. I would like to thank the support from other relatives and friends.

ABSTRACT

Graphs have been widely used to model sophisticated relationships between different entities due to their strong representative properties. Social networks, traffic networks, and biological networks are among the applications that benefit from being expressed as graphs. The cohesive subgraph is an essential structure for understanding the organization of many real-world networks, and cohesive subgraph detection is a crucial problem in network analysis. There are many cohesive subgraph models, such as k -core, strongly connected component, and maximum density subgraph.

Uncertain graph management and analysis have attracted much research attention. Among them, computing k -cores in uncertain graphs (aka, (k, η) -cores) is an important problem and has emerged in many applications. However, the existing algorithms for computing (k, η) -cores heavily depend on the two input parameters k and η . In addition, computing and updating the η -degree for each vertex is the costliest component in the algorithm, and the cost is high.

To overcome these drawbacks, we propose an index-based solution for computing (k, η) -cores. The index size is well-bounded by $O(m)$, where m is the number of edges in the graph. Based on the index, queries for any k and η can be answered in optimal time. We propose an algorithm for index construction with several different optimizations.

We also discuss the (k, η) -core computation when graphs cannot be entirely stored in memory. We adopt the semi-external setting, which allows $O(n)$ mem-

ory usage, where n is the number of vertices in the graph. This assumption is reasonable in practice, and it has been widely adopted in massive graph analysis. We design an index-based solution for I/O efficient (k, η) -core computation.

Given the frequent updates in many real-world graphs, detecting strongly connected components (SCC) in dynamic graphs is a very complicated problem. In the thesis, we study the fully dynamic depth-first search (DFS) problem in directed graphs, which is a crucial basis of dynamic SCC detection. In the literature, most works focus on the dynamic DFS problem in undirected graphs and directed acyclic graphs. However, their methods cannot easily be applied in the case of general directed graphs. Motivated by this, we propose a framework and corresponding algorithms for both edge insertion and deletion in general directed graphs. We further give several optimizations to speed up the algorithms.

We conduct extensive experiments on several large real-world graphs to practically evaluate the performance of all proposed algorithms.

PUBLICATIONS

- **Bohua Yang**, Dong Wen, Lu Qin, Ying Zhang, Lijun Chang, and Rong-Hua Li. *Index-Based Optimal Algorithm for Computing K-Cores in Large Uncertain Graphs*. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 64-75. IEEE, 2019.
- **Bohua Yang**, Dong Wen, Lu Qin, Ying Zhang, Xubo Wang, and Xuemin Lin. *Fully Dynamic Depth-First Search in Directed Graphs*. *Proceedings of the VLDB Endowment*, 13(2):142-154, 2019.
- Dong Wen, **Bohua Yang**, Lu Qin, Ying Zhang, Lijun Chang, and Rong-Hua Li. *Computing K-Cores in Large Uncertain Graphs: An Index-based Optimal Approach*. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2020.
- Dong Wen, **Bohua Yang**, Ying Zhang, Lu Qin, Dawei Cheng, and Wenjie Zhang. *Span-Reachability Querying in Large Temporal Graphs*. *The VLDB Journal*, Revision Submitted.

TABLE OF CONTENT

| | |
|---|-----------|
| CERTIFICATE OF AUTHORSHIP/ORGINALITY | ii |
| ACKNOWLEDGEMENTS | iii |
| ABSTRACT | v |
| PUBLICATIONS | vii |
| TABLE OF CONTENT | viii |
| LIST OF FIGURES | xi |
| LIST OF TABLES | xiii |
| Chapter 1 INTRODUCTION | 1 |
| 1.1 Core Computation in Large Uncertain Graphs | 2 |
| 1.2 Fully Dynamic DFS in Large Directed Graphs | 7 |
| 1.3 Graph Model | 12 |
| 1.4 Roadmap | 13 |
| Chapter 2 LITERATURE REVIEW | 15 |
| 2.1 k -Core and Uncertain Graphs | 15 |
| 2.2 SCC and DFS | 18 |
| 2.3 Other Cohesive Subgraph Models | 21 |
| Chapter 3 INTERNAL MEMORY CORE COMPUTATION IN LARGE UNCERTAIN GRAPHS | 25 |
| 3.1 Overview | 25 |
| 3.2 Preliminary | 26 |
| 3.3 Online (k, η) -Cores Computation | 28 |
| 3.3.1 An Existing Solution for η -Core Decomposition | 28 |
| 3.3.2 Our Approach to Compute (k, η) -Cores | 31 |
| 3.4 An Index-based Approach | 34 |

| | | |
|--|--|-----------|
| 3.4.1 | The Index Structure | 34 |
| 3.4.2 | Query Processing | 36 |
| 3.4.3 | Index Construction | 37 |
| 3.5 | Making Query Processing Optimal | 40 |
| 3.5.1 | Forest-based Index Structure | 41 |
| 3.5.2 | Optimal Query Processing | 43 |
| 3.5.3 | Optimizations for the Index Construction Algorithm | 46 |
| 3.6 | Experiments | 51 |
| 3.6.1 | Performance of Query Processing | 53 |
| 3.6.2 | Performance of Index Construction | 55 |
| 3.7 | Chapter Summary | 60 |
| Chapter 4 SEMI-EXTERNAL MEMORY CORE COMPUTATION IN LARGE UNCERTAIN GRAPHS | | 61 |
| 4.1 | Overview | 61 |
| 4.2 | UCF-Index in External Memory | 62 |
| 4.3 | Local η -Threshold Computation | 63 |
| 4.4 | Semi-external Algorithms | 65 |
| 4.5 | Further Optimizations | 68 |
| 4.5.1 | Reducing η -Threshold Estimations | 68 |
| 4.5.2 | Partial Neighbor Loading | 68 |
| 4.5.3 | Vertex Ordering | 69 |
| 4.6 | Experiments | 70 |
| 4.6.1 | Performance of Semi-external Query Processing | 70 |
| 4.6.2 | Performance of Semi-external Index Construction | 70 |
| 4.7 | Chapter Summary | 75 |
| Chapter 5 FULLY DYNAMIC DEPTH-FIRST SEARCH IN LARGE DIRECTED GRAPHS | | 77 |
| 5.1 | Overview | 77 |
| 5.2 | Preliminary | 78 |
| 5.3 | A Flexible Framework | 80 |
| 5.3.1 | Efficient Validity Check | 80 |
| 5.3.2 | The Framework | 82 |
| 5.3.3 | Framework Analysis | 83 |
| 5.4 | Implementations | 87 |
| 5.4.1 | Edge Insertion | 87 |
| 5.4.2 | Edge Deletion | 89 |
| 5.5 | The Improved Approaches | 90 |
| 5.5.1 | Edge Insertion | 91 |
| 5.5.2 | Edge Deletion | 95 |

TABLE OF CONTENT

| | | |
|---------------------------|--|------------|
| 5.5.3 | Batch Update | 102 |
| 5.6 | Experiments | 103 |
| 5.6.1 | Overall Efficiency | 106 |
| 5.6.2 | Effectiveness of Optimizations | 109 |
| 5.6.3 | Scalability Test | 110 |
| 5.6.4 | Memory Usage | 112 |
| 5.7 | Chapter Summary | 112 |
| Chapter 6 EPILOGUE | | 115 |
| BIBLIOGRAPHY | | 117 |

LIST OF FIGURES

| | | |
|-----|---|-----|
| 1.1 | The (k, η) -cores of \mathcal{G} for $k = 2$ and $\eta = 0.3$ | 3 |
| 1.2 | An example graph G and its DFS-Tree \mathcal{T} (γ is a virtual root connecting all vertices in G) | 8 |
| 3.1 | The <i>UCO-Index</i> of \mathcal{G} | 36 |
| 3.2 | The η -tree of \mathcal{G} for $k = 2$ | 43 |
| 3.3 | The optimization of core-based ordering | 47 |
| 3.4 | Query time for different k ($\eta = 0.5$) | 54 |
| 3.5 | Query time for different η ($k = 15$) | 56 |
| 3.6 | Query time on different datasets | 56 |
| 3.7 | Index size for different datasets | 57 |
| 3.8 | Time cost for index construction | 58 |
| 3.9 | Scalability of index construction | 59 |
| 4.1 | The <i>UCEF-Index</i> of \mathcal{G} for $k = 2$ | 62 |
| 4.2 | Query time on different datasets | 71 |
| 4.3 | I/O cost of external query processing | 71 |
| 4.4 | Memory usage for external index construction | 72 |
| 4.5 | Time cost for external index construction | 73 |
| 4.6 | I/O cost for external index construction | 73 |
| 4.7 | Speedup for external index construction | 74 |
| 4.8 | I/O reduction for external index construction | 75 |
| 4.9 | Scalability of external index construction on Orkut | 76 |
| 5.1 | The non-tree edges and time intervals of the DFS-Tree \mathcal{T} | 80 |
| 5.2 | The updated DFS-Tree for the inserted edge (v_8, v_{13}) in the graph G | 88 |
| 5.3 | The updated DFS-Tree for the deleted edge (v_5, v_6) in the graph G | 90 |
| 5.4 | Running time for edge insertion | 106 |
| 5.5 | Running time for edge deletion | 107 |
| 5.6 | Running time for tree updates | 108 |
| 5.7 | Percentage of vertices performing graph search or tree search | 111 |
| 5.8 | Scalability for edge insertion | 113 |

LIST OF FIGURES

5.9 Scalability for edge deletion 114
5.10 Memory usage 114

LIST OF TABLES

- 1.1 Notations 14
- 3.1 Network statistics 53
- 5.1 Network statistics 105
- 5.2 Percentage of forward-cross edge insertions 109
- 5.3 Percentage of tree edge deletions 109

LIST OF TABLES
