

Graph Pattern Mining: consolidating models, systems, and abstractions

Vinícius Dias², Dorgival Guedes¹

¹Departamento de Ciência da Computação – Universidade Federal de Minas Gerais (UFMG)
Belo Horizonte – MG – Brazil

²Departamento de Ciência da Computação – Universidade Federal de Lavras (UFLA)
Lavras – MG – Brazil

viniciusdias@ufla.br, dorgival@dcc.ufmg.br

Abstract. *This text summarizes the key contributions of the dissertation entitled “Graph Pattern Mining: consolidating models, systems, and abstractions”¹, approved in the Graduate Program in Computer Science of the Federal University of Minas Gerais (DCC/UFMG) Graph Pattern Mining (GPM) refers to a class of problems involving the processing of subgraphs extracted from larger graphs. Applications to GPM algorithms include querying subgraphs with given properties of interest, identifying motif structures in biological networks, among others. GPM algorithms are challenging to develop and thus, general-purpose GPM systems emerge as a solution to improve the user experience with such algorithms. In this dissertation we propose a primitive-based model for representing GPM algorithms, a distributed system implementing this model, and an extensive experimental study of popular algorithms used in GPM systems. We demonstrate empirically the effectiveness of our model by showing competitive performance without sacrificing the expressiveness of algorithms.*

Resumo. *Este texto resume as contribuições da tese intitulada “Mineração de Padrões em Grafos: consolidando modelos, sistemas e abstrações”, aprovada no Programa de Pós-Graduação em Ciência da Computação da Universidade Federal de Minas Gerais (DCC/UFMG). Mineração de Padrões em Grafos (Graph Pattern Mining, GPM) refere-se a uma classe de problemas que envolve o processamento de subgrafos extraídos de grafos maiores. As aplicações para algoritmos de GPM incluem a consulta de subgrafos com propriedades específicas de interesse, a identificação de estruturas de motivos em redes biológicas, entre outras. Os algoritmos de GPM são desafiadores de desenvolver e, assim, sistemas de GPM de uso geral surgem como uma solução para melhorar a experiência do usuário com tais algoritmos. Nesta tese, propomos um modelo baseado em primitivas para representar algoritmos de GPM, um sistema distribuído que implementa este modelo, e um extenso estudo experimental de algoritmos populares usados em sistemas de GPM. Demonstramos empiricamente a eficácia do modelo proposto, mostrando um desempenho competitivo em relação aos concorrentes, sem sacrificar a expressividade dos algoritmos.*

¹Full dissertation text available: <https://repositorio.ufmg.br/handle/1843/51806>

1. Dissertation defense information

- Data da Defesa: 24/03/2023
- Graduate program: Graduate Program in Computer Science of the Federal University of Minas Gerais (DCC/UFMG)
- Category: Doctorate
- Members of committee: Srinivasan Parthasarathy (Ohio State University), Arlei Lopes da Silva (Rice University), Ítalo Fernando Scotá Cunha (Universidade Federal de Minas Gerais), Vinícius Fernandes dos Santos (Universidade Federal de Minas Gerais), Wagner Meira Júnior (Universidade Federal de Minas Gerais)

2. Context and problem

Graphs are widely used to model problems in various areas, including web applications, social media, biological networks, brain networks, conceptual graphs, among others. In this work we navigate the trade-off between abstractions and system performance in the context of Graph Pattern Mining (GPM): a class of problems marked by the processing of subgraphs extracted from larger graphs. The relevance of GPM computation is multidisciplinary, including applications such as motif extraction from biological networks [Agrawal et al. 2018], frequent subgraph mining [Elseidy et al. 2014], subgraph searching over semantic data [Elbassuoni and Blanco 2011], social media network characterization [Ugander et al. 2013], community discovery [Benson et al. 2016], periodic community discovery [Qin et al. 2019], temporal hotspot identification [Yang et al. 2016], identification of dense subgraphs in social networks [Hooi et al. 2020], link spam detection [Leon-Suematsu et al. 2011], recommendation systems [Zhao et al. 2019], graph learning [Meng et al. 2018], to cite a few.

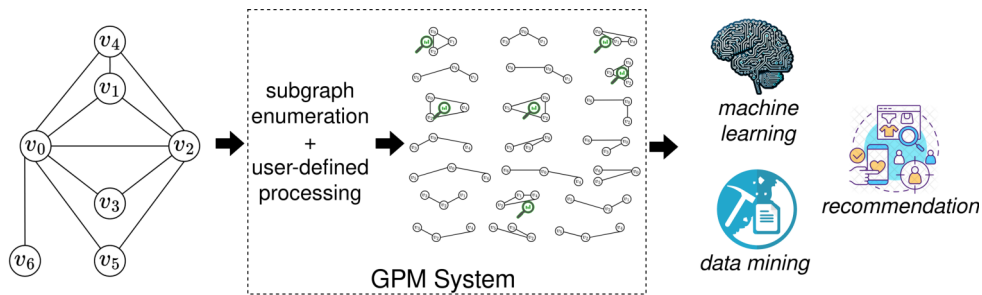


Figure 1. High-level operation of a Graph Pattern Mining (GPM) system.

In this context, general-purpose graph pattern mining systems (Fig. 1) emerge as an alternative for programming and for maintaining parallel GPM applications. Next we highlight the main challenges concerning GPM processing. *[The lack of a standard algorithm model]:* The space of existing general-purpose GPM systems is diverse and little effort has been made to model GPM algorithms in such a way that is independent of system implementation details. This limits the generalization, the proper evaluation, and the extensibility of existing GPM solutions. *[The need for an efficient, productive, and integrated GPM system]:* GPM tasks are computational intensive, irregular in terms of load balancing and memory access. They are also complex to develop from scratch since it often include non-trivial concepts from graph theory and mathematics (e.g., isomorphism and combinatorics) and often used as a pre-processing step in data analytics pipelines. It

is not trivial to address all these aspects altogether. *[The need for a fair and informative evaluation of GPM paradigms]*: Existing GPM systems are not ideal for a wide experimental characterization of GPM paradigms since implementation details are too merged into application design, this challenges the fairness of performance comparisons and the identification of opportunities for research directions.

3. Objectives

The thesis statement of this work is that GPM systems can benefit from a strong, well-defined model for algorithms that is independent of implementation details such as system architecture, programming language, and parallelization strategies. Specific objectives: (1) *Propose* a simple and expressive algorithm model for general-purpose GPM; (2) *Design and implement* a GPM system that adopts the proposed model and that deals with system challenges concerning the efficiency and the programming productivity of GPM applications. (3) *Present* an evaluation study to consolidate collective knowledge about GPM processing and to identify promising future work.

4. Contributions

We highlight three contributions of these work (Fig. 2). First, we propose a model for representing GPM algorithms, unveiling important building blocks for standardized application design, which besides improving productivity also allows a more consistent access to performance diagnostics and optimizations. Second, we provide the design and the implementation of *Fractal*, a general-purpose distributed and parallel system for GPM. *Fractal* offers an expressive and compositional programming interface, bounds memory demand via a stateless subgraph enumeration algorithm, and includes an adaptive and dynamic load balancing layer via work stealing. Third, we leverage our well defined model and system to provide an extensive experimental study of GPM workloads, including a wide range of application scenarios considering multiple algorithms and over real-world datasets.

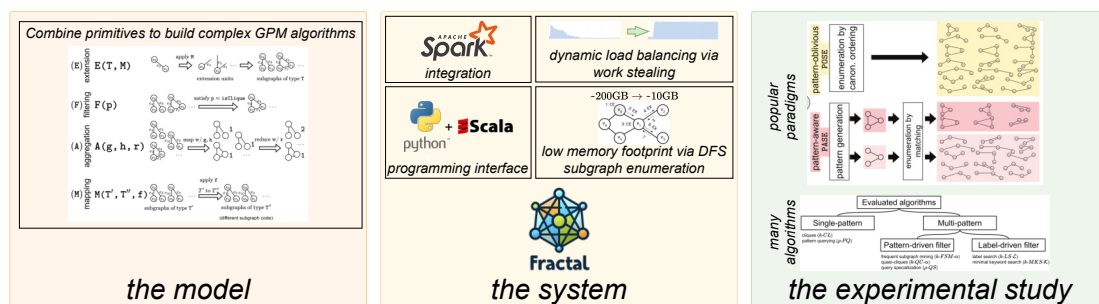


Figure 2. Overview of our contributions and solutions.

5. Advancement in the state-of-the-art

Overall, existing systems has its very particular application model and/or makes strong assumptions about the GPM tasks supported. In this work we focus on a system design that not only can be easily integrated to existing data analysis pipelines but also exhibit optimal and competitive performance for various tasks – none of the existing related work accomplishes these requirements altogether nor are adequate for

a fair experimental study of GPM algorithms. For instance, first-generation system Arabesque [Teixeira et al. 2015] suffers from intractable memory demands. Peregrine [Jamshidi et al. 2020] and Automine [Mawhirter and Wu 2019] are hand-tuned approaches that do not handle distributed environments nor offer integrated solution capabilities. G2Miner [Chen and Arvind 2022] (GPU) and G-Miner [Chen et al. 2018] (distributed) offer an efficient system design but not easily integrated into data processing tools. Tesseract [Bindschaedler et al. 2021] handles distributed and dynamic graphs but lacks the flexibility for supporting multiple paradigms in its subgraph enumeration engine.

Our model (Fig. 2-left) is based on a set of concise primitives that can be combined in such a way that abstracts all the low level implementation details and parallel/distributed deployment. Our system (Fig. 2-center) implements the primitive based model on top of Spark framework. We extended Spark’s execution model to accommodate our optimizations: load balancing via dynamic work stealing, depth-first subgraph enumeration for bounded memory, and integration with Resilient Distributed Datasets (RDD) abstraction. Our experimental study (Fig. 2-right) include multiple GPM algorithms (categorized by semantics) and a wide range of real-world graphs. We reimplemented the algorithms under the same conditions to enable a fair comparison and an informative set of conclusions.

6. Evaluation

Fractal is efficient and competitive with existing baselines: We evaluated Fractal’s performance against specialized baselines (able to solve a single problem) and other general-purpose GPM systems over real-world graphs. Fractal outperforms some baselines and stays competitive against hand-optimized ones – although Fractal is not always more efficient, it allows solving various problems with a reduced programming effort. Our evaluation also unveils that system optimizations proposed are able to enhance resource utilization by near-perfect load balancing and orders of magnitude improvement in memory demand. *Experimental study enhance knowledge about trade-offs between GPM paradigms:* We implemented and evaluated popular GPM algorithms using our solutions, enforcing comparison fairness and unveiling challenges and opportunities for the area. Our main conclusion is that there is no silver bullet in terms of GPM paradigms, contrary to existing claims in the literature.

7. Scientific and Technical Production, and Awards

The following are scientific and technical products of this dissertation:

1.[publication] Full paper (Qualis A1) [Dias et al. 2019]: *Dias, V., Teixeira, C. H. C., Guedes, D., Meira Jr., W., and Parthasarathy, S. (2019). Fractal: A general-purpose graph pattern mining system. In Proceedings of the 2019 International Conference on Management of Data (SIGMOD).*

Link:

<https://doi.org/10.1145/3299869.3319875>

2.[publication] Full paper with best paper nomination (honors) (Qualis A3) [Dias et al. 2023]: *Dias, V., Ferraz, S., Vadlamani, A., Erfanian, M., Teixeira,*

C. H., Guedes, D., Meira, W., and Parthasarathy, S. (2023). *Graph pattern mining paradigms: Consolidation and renewed bearing*. In *2023 31st IEEE International Conf. on High Performance Computing, Data, and Analytics (HiPC)*.

Link:

<https://doi.org/10.1109/HiPC58850.2023.00040>

3.[software] Fractal project is publicly available, including reproducibility artifacts and guides on how to deploy and to use our system to produce new integrated solutions to custom GPM problems.

Code link (including reproducibility of above papers):

<https://github.com/dccspeed/fractal>

Data link:

<https://drive.google.com/drive/folders/1ViLAlQt45hFDtqTCJnOfqk4WZ71E3IUN>

4.[short course and book chapter] Short course and hence, also a book chapter, accepted to this year’s 39th Brazilian Database Symposium (SBBD ’24): “Practical Graph Pattern Mining: Systems, Applications, and Challenges”. It is intended to share the knowledge produced with this dissertation with the database community through practical examples with Fractal.

References

- Agrawal, M., Zitnik, M., and Leskovec, J. (2018). Large-scale analysis of disease pathways in the human interactome. *Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing*, 23:111–122.
- Benson, A. R., Gleich, D. F., and Leskovec, J. (2016). Higher-order organization of complex networks. *Science*.
- Bindschaedler, L., Malicevic, J., Lepers, B., Goel, A., and Zwaenepoel, W. (2021). *Tesseract: Distributed, General Graph Pattern Mining on Evolving Graphs*, page 458–473. Association for Computing Machinery, New York, NY, USA.
- Chen, H., Liu, M., Zhao, Y., Yan, X., Yan, D., and Cheng, J. (2018). G-miner: An efficient task-oriented graph mining system. In *Proceedings of the Thirteenth EuroSys Conference*, EuroSys ’18.
- Chen, X. and Arvind (2022). Efficient and scalable graph pattern mining on GPUs. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, pages 857–877, Carlsbad, CA. USENIX Association.
- Dias, V., Ferraz, S., Vadlamani, A., Erfanian, M., Teixeira, C. H., Guedes, D., Meira, W., and Parthasarathy, S. (2023). Graph pattern mining paradigms: Consolidation and renewed bearing. In *2023 31st IEEE International Conference on High Performance Computing, Data, and Analytics (HiPC)*.
- Dias, V., Teixeira, C. H. C., Guedes, D., Meira Jr., W., and Parthasarathy, S. (2019). Fractal: A general-purpose graph pattern mining system. In *Proceedings of the 2019 International Conference on Management of Data (SIGMOD)*.

- Elbassuoni, S. and Blanco, R. (2011). Keyword search over rdf graphs. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management, CIKM '11*, pages 237–242, New York, NY, USA. ACM.
- Elseidy, M., Abdelhamid, E., Skiadopoulou, S., and Kalnis, P. (2014). Grami: Frequent subgraph and pattern mining in a single large graph. *Proc. VLDB Endow.*, 7(7):517–528.
- Hooi, B., Shin, K., Lamba, H., and Faloutsos, C. (2020). Telltail: Fast scoring and detection of dense subgraphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):4150–4157.
- Jamshidi, K., Mahadasa, R., and Vora, K. (2020). Peregrine: A pattern-aware graph mining system. In *Proceedings of the Fifteenth European Conference on Computer Systems, EuroSys '20*, New York, NY, USA. Association for Computing Machinery.
- Leon-Suematsu, Y. I., Inui, K., Kurohashi, S., and Kidawara, Y. (2011). Web Spam Detection by Exploring Densely Connected Subgraphs. In *2011 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology*, volume 1, pages 124–129.
- Mawhirter, D. and Wu, B. (2019). Automine: Harmonizing high-level abstraction and high performance for graph mining. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP '19*, pages 509–523, New York, NY, USA. ACM.
- Meng, C., Mouli, S. C., Ribeiro, B., and Neville, J. (2018). Subgraph pattern neural networks for high-order graph evolution prediction.
- Qin, H., Li, R.-H., Wang, G., Qin, L., Cheng, Y., and Yuan, Y. (2019). Mining periodic cliques in temporal networks. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 1130–1141.
- Teixeira, C. H. C., Fonseca, A. J., Serafini, M., Siganos, G., Zaki, M. J., and Aboulmaga, A. (2015). Arabesque: A system for distributed graph mining. In *Proceedings of the 25th Symposium on Operating Systems Principles, SOSP '15*, pages 425–440.
- Ugander, J., Backstrom, L., and Kleinberg, J. (2013). Subgraph frequencies: mapping the empirical and extremal geography of large graph collections. In *WWW*.
- Yang, Y., Yan, D., Wu, H., Cheng, J., Zhou, S., and Lui, J. C. (2016). Diversified temporal subgraph pattern mining. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, page 1965–1974, New York, NY, USA. Association for Computing Machinery.
- Zhao, H., Zhou, Y., Song, Y., and Lee, D. L. (2019). Motif enhanced recommendation over heterogeneous information network. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM '19*, pages 2189–2192, New York, NY, USA. ACM.