Master Thesis

# Collaboration networks in open-source software development

Jozsef Csepanyi

Date of Birth: 06.09.1996
Student ID: 11927479

**Subject Area:** Information Systems

**Studienkennzahl:** h11927479

**Supervisor:** Johannes Wachs

**Date of Submission:** 31.March 2021

# Contents

# List of Figures

# List of Tables

## Abstract

Aenean commodo ligula eget dolor. Aenean massa. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Donec quam felis, ultricies nec, pellentesque eu, pretium quis, sem. Nulla consequat massa quis enim. Donec pede justo, fringilla vel, aliquet nec, vulputate eget, arcu. In enim justo, rhoncus ut, imperdiet a, venenatis vitae, justo. Nullam dictum felis eu pede mollis pretium. Integer tincidunt. Cras dapibus. Vivamus elementum semper nisi. Aenean vulputate eleifend tellus. Aenean leo ligula, porttitor eu, consequat vitae, eleifend ac, enim. Aliquam lorem ante, dapibus in, viverra quis, feugiat a, tellus. . . .

# 1 Introduction

In recent years open source software solutions have become widely popular and frequently used in both scientific and enterprise use, which can be attributed to a number of factors, most importantly the ease of development and deployment of IT projects, improved cybersecurity and enhanced scalability [3]. This increases the contribution to open source projects from enterprises and individuals alike. Due to its nature, open source software projects are driven by community contributions, and depend heavily on active participation in all phases of the project. Because there is a high dependency on the community in open source software projects, by understanding how contributions are included and what patterns emerge we can gain valuable insight into the project's current state and its trajectory.

. . .

## 1.1 Literature review

Software development in a corporate environment usually follows a strict hierarchial structure, where each participant is given a precise position and responsibility, like project manager, scrum master, senior or junior developer, and employees do not tend to work outside of their assigned tasks and territories. The main purpose of maintaining software development structures is for the company to ensure that the outcome of the project is in accordance with the business objective, adheres to the pre-set quality criteria and it is completed in a given timeframe; in other words to asses the risks associated with the business objective of the software project [5]. This is achieved by breaking down the developed software into smaller, less complex components, and grouping the developers into managable teams, where the communication is moderated between teams [1].

On the other hand, Free/Libre Open Source Software (FLOSS) projects usually do not follow an organizational hierarchy, and are usually self-organizing and dynamic [1]. Issues, bugs and progress are tracked openly, and everyone is encouraged to contribute based on the current topics and expertise, but purely on a volunteering basis. The lack of access restriction to certain modules allows for much more spontaneous interaction between developers, which generate large, complex networks [2]. These complex networks can be seen as large social networks of developers based on collaboration.

The collaboration networks of open source software (OSS) have been a subject of many academic research. Raymond [4] has defined collaboration based on bug report interaction, and observed the collaboration network of 124 large-scale SourceForge projects. The generated networks have widely different centralization properties, but it was observed that larger sized projects tend to be more decentralized. The broad community roles contributors tend to take have been also identified in [4], which have been coined as the *onion model*.

- Open-source software development properties

  - centralized vs decentralized
  - No collocation
  - Enterprise support
  - Version control, issue tracking

- Relevant social aspects of OS projects

- State of the art

  - Collaboration by coediting files
  - Contributors form dynamic social networks
  - Problem of analysing changes over time in a network
  - Other studies in this field...

- Preliminary analysis results (pandas, networkx, ...)

## 1.2 Motivation of research problem and research question

- Importance of OS project analysis based on lit rew

- Analysing effects of large events within the lifecycle of the OS project in order to improve them or adapt

  - Planned, foreseeable changes (e.g. upcoming major release)
  - Unforeseeable changes (e.g. end of support, pandemic)

- Research questions

– What social patterns emerge within large-scale open-source software projects?

    ∗ Are there smaller "core" collaborator networks connected with weak links or do they form one large interconnected network?

    ∗ Are there usually key contributors, who are central to the project and collaborate with most contributors, or is it completely decentralized?

    ∗ How does the size of the project change these properties?

– How does the structure of OS software development collaboration change over time?

    ∗ Are there any major changes over the natural project lifecycle? Are they visible in the collaboration network? (e.g. planning, developing, bugfixing, sunset?)

    ∗ How does a sudden major event change the participation and development?

# 2 Proposed research method

- Developing a tool, that can extract the collaboration information from any OS project (from GitHub/git repository)

- Data cleaning - method to merge authors, excluding common folders, etc. . .

- Qualitative research

  – Observing collaboration statistics and networks in order to discover patterns: connected components, centrality, changes over time

- Quantitative research

  – Composing a large set of repositories (different sizes, properties)

  – Detecting past changes automatically based on changes in measured statistics

# 3 Outline of thesis

- Literature review

- – Network analysis, relevant metrics
- – Properties of social collaboration networks

- Used repositories, selection criteria

- Data cleaning - files, authors, max modifications

- Implementation

  - – . . .

- Qualitative analysis

- Quantitative analysis

- Conclusion

## 3.1 (Preliminary literature list - in references)

## 3.2 Work plan including milestones

- Data cleaning - files, authors, max modifications

- Implementation

  - – . . .

- Qualitative analysis

- Quantitative analysis

- Conclusion

# References

[1] Christian Bird, David Pattison, Raissa D'Souza, Vladimir Filkov, and Premkumar Devanbu. Latent social structure in open source projects. In *Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, SIGSOFT '08/FSE-16, pages 24–35, New York, NY, USA, November 2008. Association for Computing Machinery.

[2] Juan Martinez-Romo, Gregorio Robles, Jesus M. Gonzalez-Barahona, and Miguel Ortuño-Perez. Using Social Network Analysis Techniques to Study Collaboration between a FLOSS Community and a Company. In Barbara Russo, Ernesto Damiani, Scott Hissam, Björn Lundell, and Giancarlo Succi, editors, *Open Source Development, Communities and Quality*, volume 275, pages 171–186. Springer US, Boston, MA, 2008.

[3] PwC. Leading benefits of open-source software among enterprises worldwide as of 2016. *Statista*, 2016.

[4] Eric S. Raymond. The cathedral and the bazaar. https://firstmonday.org/ojs/index.php/fm/article/download/578/499?inline=1, March 1998.

[5] Ashish Sureka, Atul Goyal, and Ayushi Rastogi. Using social network analysis for mining collaboration data in a defect tracking system for risk and vulnerability analysis. In *Proceedings of the 4th India Software Engineering Conference*, ISEC '11, pages 195–204, New York, NY, USA, February 2011. Association for Computing Machinery.