

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: 02. April 2021

Department of Information Systems and Operations, Vienna University of Economics and Business, Welthandelsplatz 1, 1020 Vienna, Austria



**DEPARTMENT FÜR INFORMATIONS-
VERARBEITUNG UND PROZESS-
MANAGEMENT** DEPARTMENT
OF INFORMATION SYSTEMS AND
OPERATIONS

Contents

1	Introduction	6
2	Background and rationale	7
2.1	Collaboration in FLOSS projects	7
2.2	Social network of Open-source projects	8
2.3	OSS project success	9
3	Motivation of research problem and research questions	9
3.1	Research methodology	11
4	Gitminer implementation	12
4.1	git2net miner	12
4.2	repo_tools miner	13
4.3	Data preprocessing	13
4.4	Collaboration networks	14
4.4.1	Temporal bipartite network	14
4.4.2	Static networks	15
4.5	Core and periphery, centralization	18
4.5.1	Degree centrality	18
4.5.2	Degree centralization	19
4.5.3	Clustering coefficient	20
4.5.4	Hierarchy	21
4.6	Project measures	22
4.6.1	Release and release measures	22
4.6.2	Project issues and measures	25
4.7	Time window	27
5	Collaboration pattern analysis	27
5.1	Observed projects and events	27
5.2	SNA metrics analysis	27
5.2.1	K-cores	27
5.3	Results	28
6	Quantitative analysis of projects during crunch time	28
6.1	Collaboration network changes	28
6.2	Prediction of outcome based on collaboration changes	28
7	Discussion and results	28
8	Conclusion and future work	28

List of Figures

1	Onion model of collaboration types in FLOSS projects [9]. . .	7
2	Sequential snapshots of the networkx collaboration network with a moving time-window of 30 days and 7-day steps.	15
3	A bipartite network of authors (green) and edited files (light blue) in the pandas project within the timeframe 31/01/2021 and 15/02/2021	16
4	Weighted Jaccard similarity collaboration network of pandas generated from the bipartite network in Figure 3.	17
5	Degree centrality within the pandas and curl projects' collaboration networks.	19
6	The local clustering coefficient demonstrated on an unweighted network of 4 vertices [15].	20
7	Hierarchical network and its corresponding degree number vs clustering coefficient plot.	22
8	Semantic versioning. Figure source: [1].	23
9	Survival curves of the pandas library.	26

List of Tables

1	Releases collected information.	24
2	Keywords and drop words for <i>bug</i> and <i>feature</i> categorization. .	27

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 [26]. 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.

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 [29]. 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 [7].

As opposed to commercial software development, Free/Libre Open Source Software (FLOSS) projects usually do not follow an organizational hierarchy, and are usually self-organizing and dynamic [7]. 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 [19]. These complex networks can be seen as large social networks of developers based on collaboration.

Because contribution to FLOSS projects are voluntary, participants have a different motivation for taking part than in commercial software development. According to El Asir et al. [11], FLOSS participation can be motivated by internal and external factors. Internal factors include self-improvement,

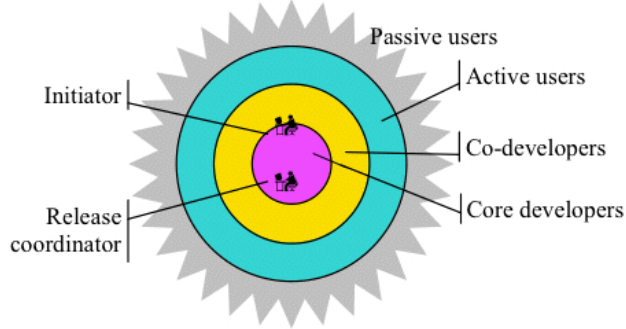


Figure 1: Onion model of collaboration types in FLOSS projects [9].

learning and contribution as a hobby or pass-time activity [2, 31], whereas external factors are motivated by marketing and demonstrating certain skills, thus increasing and improving employability [2].

2 Background and rationale

2.1 Collaboration in FLOSS projects

Collaboration networks of open source software (OSS) have been a subject of many academic research. Raymond [9] 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 [9], which have been coined as the *onion model* in [19] (Figure 1).

The onion model describes the types of participants in an OSS project as layers. The center represents the small group of core developers, who are responsible for the majority of contributions to the software. They are surrounded by a larger group of co-developers, whose main contributions are usually bug fixes reported by the active users. The passive users are usually the largest in numbers, who do not contribute or report any bugs. In a healthy FLOSS project, each layer of contributors are about one magnitude larger in numbers than the preceeding inner layer [22].

El Asir et al. [11] used a K-means classification to categorise project

participants into a similar core-periphery structure (core, gray in-between area, and periphery) based on SNA metrics with a monthly timeframe, and analysed how and why contributors transition between groups. They found that technical contributions like code commits and lines added have a much heavier impact on becoming a core developer as opposed to other activities, such as testing, reviewing and commenting.

A literature review conducted by McClean et al. [21] systematically analysed the state-of-the-art research of 46 scientific papers in the field of FLOSS social networks, and categorised them into three groups based on topic: structure, lifecycle and communication. They conclude, that the existence of core-periphery structure in OSS projects is well established in the field, which is also an indicator of a healthy FLOSS software. Regarding the lifecycle, generally the core development team does not change significantly over time, however, the project becomes more decentralised and distributed as it matures. A lack of research regarding temporal analyses were identified in the most current knowledge, which was suggested as a future research area in this field.

2.2 Social network of Open-source projects

In a larger FLOSS project, developers usually cannot understand every part of the project, therefore collaboration is required with each other. The network created by the collaborators can be considered as a social network, because collaboration requires some kind of social interaction with each other. Social network theory describes how social interaction patterns affect the individual behaviour [20]. We can model an OSS project's social network as a graph, where nodes represent collaborators (developers, bug reporters, etc...) and edges represent the social interaction between them. In mathematical terms...

The type of the social interaction determines the created network, therefore choosing the basis of collaboration can have a significant impact on the network structure. The common types of developer social networks (DSNs) are Version Control System-based (VCS-DSN), Bug Tracking System-based (BTS-DSN) networks and DNSs, that are purely based on social elements [3]. The VCS-DSN take the version control application as a source for network generation by recording collaboration based on co-edits of the same module, file or code section. Choosing the granularity can impact the precision of true collaborations represented in the network. Co-edits by multiple

developers to a single module or file does not necessarily mean actual collaboration was required from the authors, as the parts edited could work functionally independent from each other. By increasing the granularity to file sections (classes or functions within a single file) or even lines, we can be more certain, that coordination was required, but we risk leaving out semantically connected parts of the project [17]. In contrast to VCS-DSNs' purely technical approach, the BTS-DSNs use semi-technical bases for connecting participants, such as comments on issues, bugs or reviews [11]. These artifacts, although being tightly related to specific sections of the source code, allow for taking into account conversational elements as contribution. For example, participants, who do not contribute directly to the software source code, but actively review and comment, are also considered. Lastly, social networks of developers can be constructed on project participation, following, starring or through communication means like mailing lists. The technical aspect of collaboration is minimized in such DSNs, and they are more fit for project organization and communication analyses in FLOSS projects (mailing list vs file file edits vs line edits)

2.3 OSS project success

community maturity ([18] in [3])

"Successful projects will likely have modular structure from the start or after refactoring as the source code grows larger and more unwieldy" [4]

success factors: Average Time Efforts, Number of Developers, Comments, Total Code Lines, Comment Ratio, Number of Rater [30] Truck Factor [5]

3 Motivation of research problem and research questions

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. As stated before, SNA analysis of OSS have been extensively studied, but there is a lack of research regarding temporal models analyzing the lifecycle of a FLOSS project.

The goal of this paper is to fill in this gap by examining OSS project collaboration networks over time using SNA metrics. More specifically, one part of the research will focus on the evolution of such collaboration networks and comparing and contrasting these networks with the software outcome. The second part will focus on events during a project, and how it affects the developer collaboration. The research questions, which are broken down into subquestions, are as follows:

1. How does the temporal lifecycle information of a project influence its success?

- (a) *Based on temporal models of collaboration, is it possible to predict the outcome of the project?* Since it has been proven that the core collaborators do not change much over the course of the OS software development, our assumption is that any sudden or long-term change, that is not consistent with the other observed projects, can have a significant impact on the outcome (negative or positive alike).
- (b) *Can stages of a FLOSS project with a maturity model be observed?* As most OS software starts with a small collaborator basis and grows over time, it can be assumed, that each project goes through the same steps of open source maturity levels. On the other hand, it is also possible that due to the uniqueness of each project, no such stages are observable.

2. How do major events in the project lifecycle change the collaboration network of the project?

- (a) *Do planned or foreseeable events change the collaboration structure?* Major software version releases can be considered foreseeable events of the project lifecycle, which could have an effect on the developer collaboration. For example, there might be a higher rate of interaction between contributors just before a new version is released to clear up the backlog of tasks. But it is also possible, that commit and change rates drop during this time, because the focus shifts to stability and testing instead of new features.
- (b) *How unforeseeable internal or external events affect FLOSS collaboration?* Sudden shocks to the project, such as an announcement of disinterest from major users of the software, discontinued enterprise support of the project, large-scale global events like the

pandemic, or sudden employee firings can have significant effect on the core and periphery collaborators alike. By analysing the collaboration network before, during and after such changes, we might be able to recognise patterns, that regularly occur around these events.

3.1 Research methodology

To find answers to the research questions above, first we build a repository analyzer tool, which mines collaboration data from FLOSS projects, generates static snapshot collaboration networks at each given time interval and calculates SNA metrics for each snapshot. Then these metrics can be aggregated over time, or plotted against time to discover changes in the network. The `git2net`¹ [12] Python library provides the necessary tools to mine any project repository that uses git version control. It also incorporates temporal network generation capability, which can be used as a source for creating static collaboration networks aggregated over a given period of time.

We apply a hybrid methodology of qualitative and quantitative research. First, as part of the qualitative research, we choose a small number of repositories to be analyzed. We observe the number of connected components, centrality, number of nodes and mean degree SNA metrics in order to discover the core and peripheral collaborators over the project lifecycles. The basis of collaboration, due to the unavailability of other means of communication, is coediting files. Based on the state of the art research in this field, file coediting proves to be an effective and easy way to represent collaboration between developers.

After discovering the collaboration structure over time, we will match the breakpoints and unexpected spikes or troughs to events within the lifespan of the project. We expect that the key SNA metrics will show a periodicity around planned releases and other reoccurring events (e.g. holiday season). Outstanding values without reoccurrence, on the other hand, are more likely to be consequences of unexpected events. In these cases, it should be observed whether the network is capable of reorganizing itself, or does the event leave a permanent mark on the collaboration structure. A categorization of unexpected events and the level of impact each category has should be observed.

For the quantitative research to be conducted, we will gather a large set

¹<https://github.com/gotec/git2net>

of repositories along with major events in its lifecycles. We will then run the miner for all repositories, and with the findings of the qualitative research, we will try to detect all major events and their type (planned or unexpected). We will utilize the `ruptures`² library to detect changes in the continuous SNA metrics. If the model is capable to accurately recognise events, then we can also apply it on any repository to detect changes, which will allow us to discover changes in the collaboration network that are not related to publicly known events or releases.

4 Gitminer implementation

To find answers to the research questions, we implement an analysis tool to mine and analyze project repositories, which allows us to generate collaboration networks and network metrics for the analysed projects.

4.1 git2net miner

The process begins with the project mining. After cloning the repository, the `git2net` [12] library is used to collect data related to commits. Specifically, who is the author of each commit, which files were modified (created, edited, deleted) with the commit, and when was the commit created. Additionally, the lines edited by the author within each commit are collected separately, allowing for a more fine-grained collaboration network generation if necessary. The results are collected into an `sqlite`³ database file's *commits* and *edits* tables.

The `git2net` mining process by default collects all the commits throughout the project's lifecycle. However, the processing time of each commit differs based on the number of edits, the affected number of files and the file types as well, which makes collecting certain commits very resource-intensive and time-costly. Therefore, we exclude every commit, which contains more than 100 file modifications, during each repository mining using the *max_modifications* parameter. As observed by Gote et al [12], this exclusion criteria does not affect significantly the generated network, because they are mostly merge commits or project restructurings, which do not mark any

²<https://github.com/deepcharles/ruptures>

³<https://www.sqlite.org/index.html>

true collaboration effort between developers. During the data mining in certain repositories, we encountered commits, that were not mineable with this method and the mining process halted, presumably due to processing error because of binary file changes in these commits. We also excluded these commits from our data mining process.

This exclusion criteria resulted in an average of 3% of commits excluded in all repositories subject to our analyses, with the highest excluded commit rate being 20%.

4.2 `repo_tools` miner

We use the `repo_tools` ⁴ Python library to query the Github API for additional repository data extraction, such as:

- Releases
- Tags
- Issues
- Stars and followers

The mining output is also stored in a `sqlite` relational database, which is queried later on during the analysis.

4.3 Data preprocessing

The collaboration networks with the `git2net` library connect the authors to their edited files using only file and author names instead of IDs. This creates an issue when generating the networks, because authors with the same name will show up as one node, and they will be connected to the files they touched combined. Furthermore, authors that change their displayed name ('author_name' field in the mining database) or log in from different accounts, where they have different names, will show up as multiple nodes instead of a single vertex.

We utilize the `gambit` [13] rule-based disambiguation tool to resolve the author names. Furthermore, the created networks have issues when the node names contain special characters or spaces. Therefore, after disambiguation,

⁴https://github.com/wschuell/repo_tools/

we replace every unique author name with its ID number.

As the files are also labelled by their filename property in the network outputs, the same filenames but in different folders are also displayed as single nodes. In order not to create false collaborations, we simply remove the files from the network with filenames, that occur more than once in all the repository subdirectories. We argue that this does not remove any significant collaboration data, since most files sharing their name with other files are technical files, like `__init__.py` for a Python project.

4.4 Collaboration networks

When creating a DSN from the mined data, we have multiple methods at hand. The `git2net` library provides its own co-editing network function, which returns a temporal network of collaborators. This uses the co-authorship algorithm developed by Gote et. al. [12], however, we would like to have more control over the network generation method, such as simple file-based co-authorship in order to customize the network for our needs, like weighing each relation or generating undirected graphs.

4.4.1 Temporal bipartite network

As a first step, we generate a temporal bipartite network of authors and their edited files with the `git2net` built-in `get_bipartite_network` method. A temporal network is a `pathpy`⁵ graph object, which contains a collection of timestamped graphs of a single network at each point in time within the observed timeframe. Such a snapshot $S_t = (U, V, E_t)$, where U is the set of authors, V is the set of files and E_t is the set of file edits as edges at t timestamp. By connecting the authors, who touched the same files, and removing the nodes representing the edited files (converting the bipartite network to a regular network), we can observe the evolution of the collaboration over time, represented in Figure 2.

Although a temporal network preserves the time aspect of the graph by the edges being tied to the time dimension of the graph, calculating network metrics like centrality on such networks is infeasible. Visualization also proves to be difficult in representations where animation is not possible.

⁵<https://www.pathpy.net/>

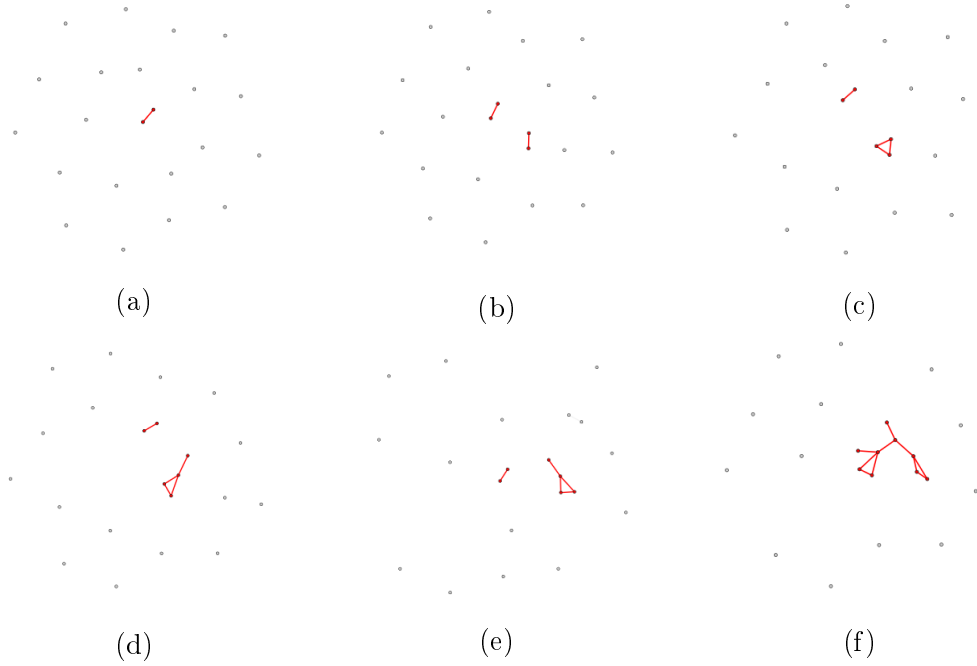


Figure 2: Sequential snapshots of the `networkx` collaboration network with a moving time-window of 30 days and 7-day steps.

Therefore, we aggregate the bipartite network over a given timeframe into a static network. All nodes within the temporal net are preserved, and all directed edges are added to the network with the edge weight representing how many times that author edited the file.

4.4.2 Static networks

The generated static weighed bipartite network loses its time-varying component, but now we are able to manipulate and calculate complex statistics over it. Figure 3 is an example of such a network. As a next step, we convert the the bipartite network into an authors' network by removing the nodes representing files.

We have multiple methods to convert the directed and weighted bipartite network into a projection of authors. We could simply remove the files and connect each author, that worked on the same file, however, the end result would be an unweighted graph. This would falsely show, that all collaborations are weighed equally, which is clearly not the case, as multiple continuous

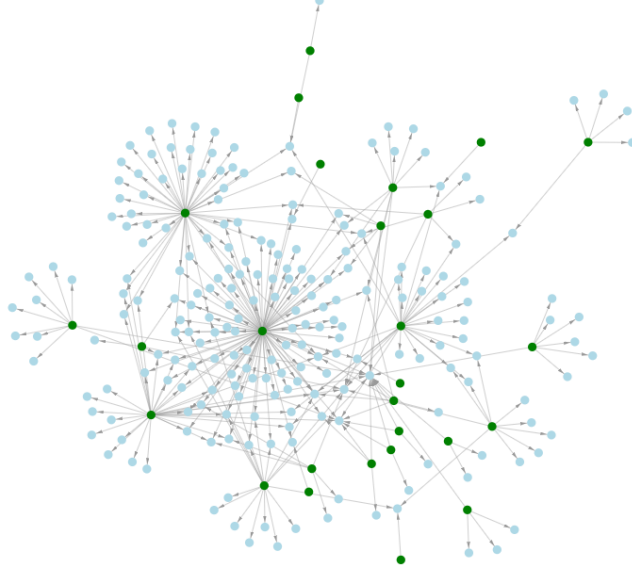


Figure 3: A bipartite network of authors (green) and edited files (light blue) in the `pandas` project within the timeframe 31/01/2021 and 15/02/2021

edits on the same file from both parties should represent a stronger collaborative connection. Therefore, firstly we implement the Weighted One-Mode Projection (WOMP) method [28]. The WOMP method converts the bipartite network $G(A, F, E)$, where E is the edge list containing tuples (a_i, f_i, w_{ij}) , and $w_{ij} \in E$ is the weight between author $a_i \in A$ and file $f_i \in F$. With this notation, a weighted directed edge can be calculated for any $a_a, a_b \in A$ as follows:

$$w_{ab}^{A \rightarrow A} = \sum_{j=1}^m \frac{w_{aj}}{W_a^F},$$

where W_a^F is the sum of all outgoing edge weights from author a to all files F denoted as $W_a^F = \sum_{i=1}^n w_{ai}$. This creates a bidirectional weighted collaboration network between authors a_1 and a_2 , where the weight w_{12} represents the relative collaboration effort of a_1 towards a_2 compared to all the other developers a_1 has collaborated with. Consequently, every edge is in the range $[0, 1]$ in the resulting WOMP network.

A disadvantage of the WOMP method is, that the generated collaboration network is bidirectional, meaning if there was any common authored files between a_1 and a_2 , then there will be both w_{12} and w_{21} connecting them.

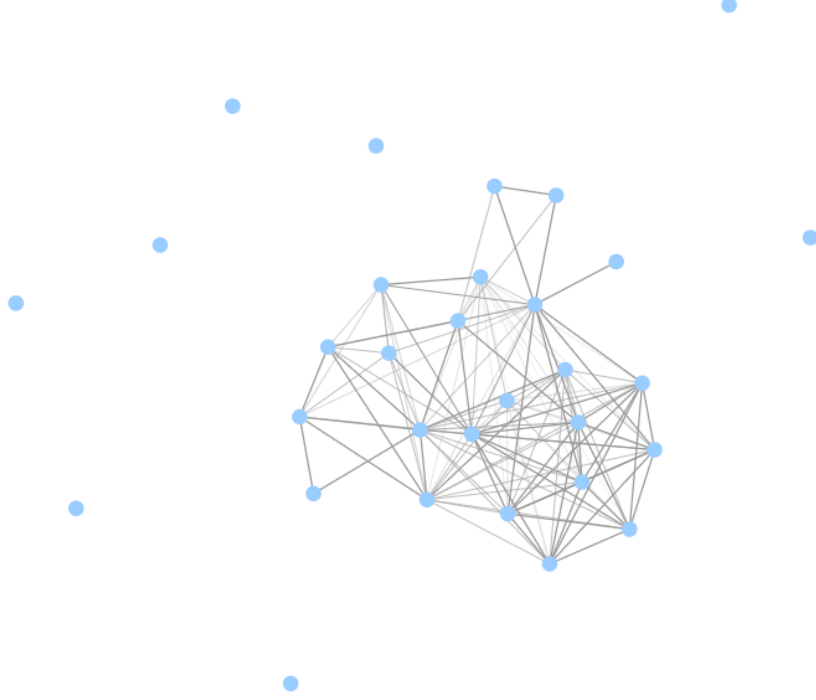


Figure 4: Weighted Jaccard similarity collaboration network of `pandas` generated from the bipartite network in Figure 3.

To simplify the network, we want to generate an authors network, where the edges are undirected. For this, we are using the weighted Jaccard method on the files-authors bipartite network:

$$w_{ab} = \frac{\sum_{f \in F} \min(f_a, f_b)}{\sum_{f \in F} \max(f_a, f_b)}.$$

For each file f that a_1 and a_2 authors touch, we sum up the minimum and maximum weights the authors have towards each file, then we divide the sum of minimums with the sum of maximums. This results in the undirected author-author network with edge weights in range $[0, 1]$. By default, this method removes isolated contributors, who do not collaborate with each other, but are actively editing the files. We add these nodes manually. Figure 4 shows the final author network.

4.5 Core and periphery, centralization

A critical part of the OSS software projects is the existence of core and periphery developers. It has been observed, that in each FLOSS project there are a small number of developers, who provide the vast majority of development effort into the project. It has been also established, that the members of core developers do not change substantially during the project’s lifecycle. However, there was no effort on whether there is a change in the collaboration pattern, especially before, during or after a major lifecycle event. Therefore, we make efforts identifying the core developer network to observe these changes.

4.5.1 Degree centrality

We use the degree centrality of each node (i.e. developer) to identify the core members. Degree centrality of a node is the fraction of all possible nodes it is connected to. We can calculate it by dividing the degree with $n - 1$, where $n = |G|$ the number of nodes within the network. Since the core developers contribute the majority of commits and edits of the project, they are expected to be connected with more nodes. Joblin et. al. [16, 17] have also identified degree centrality as the best predictor of core developers. In cases, where binary classification of core or periphery is needed, we assign developers to the core network if their degree centrality score is in the top 20th percentile, otherwise they are considered as periphery. We also take note, that this method does not consider the weighted edges, only the number of edges (degree) a node has. Although this method could be refined to consider the node degree weighted with the edges, we argue that this could lead to invalidity. In case of two developers, who only contributed to one file, they will be represented with a strong connection and would receive a high weighted degree value, whereas a core contributor, who edits many files, can have many weak connections but these might not add up to one strong connection of the two isolated developers when weighted with the edge weights. It is clear, that a developer with many connections, regardless of the strength of the collaboration, should be considered core. Figure 5 shows two examples for degree centrality within a collaboration network. In the figure, darker colors represent a higher degree centrality value. The highlighted nodes are in the highest 20th percentile of degree centrality, classifying as members of the core developers. We can observe in both one-month periods, that `pandas` is much more decentralized, whereas `curl` is largely dependent on one developer.

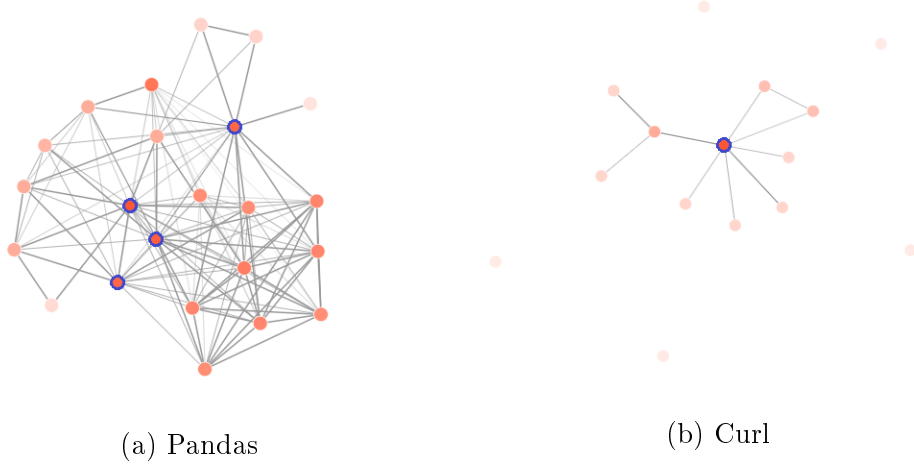


Figure 5: Degree centrality within the `pandas` and `curl` projects' collaboration networks.

4.5.2 Degree centralization

The degree *centrality* can be calculated for every node, however, through our analysis we would also like to measure a global *centralization* metric, which is applicable to the whole network. As suggested by Crowston and Howison [10], we calculate the degree *centralization* by summing the differences between the maximum and each node's degree *centrality*.

$$C_D(A) = \frac{\sum_{i=1}^n (C_d(a^*) - C_d(a_i))}{H},$$

where $C_d(a)$ is the degree *centrality* of an author a , a^* is the author with the highest degree *centrality* value, and n is the number of authors in the collaboration network A . The value H is for normalizing the sum by dividing by the theoretical maximum *centralization*. Since the *centrality* values are already in the range $[0, 1]$, we only need to normalize for the network's size. We get the highest centrality score with a star graph, where each node is only connected to a single central node, which has exactly one edge to all other nodes. The central node has a centrality of 1 in this case, whereas all the other $n - 1$ nodes have $C_d(a) = \frac{1}{n-1}$. This means that in case of a star graph:

$$H = (n - 1) \left(1 - \frac{1}{n - 1}\right) = n - 2.$$

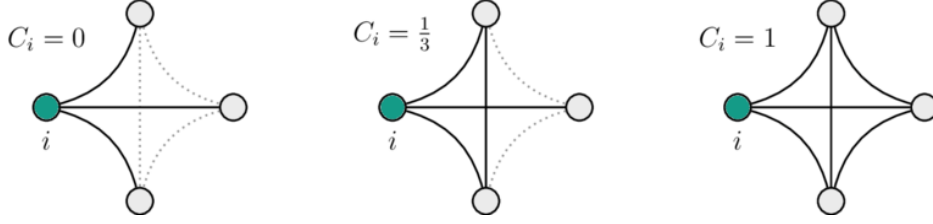


Figure 6: The local clustering coefficient demonstrated on an unweighted network of 4 vertices [15].

Within certain timeframes, when a project is inactive, it could happen, that the network contains 2 nodes or less. We define $C_D(A) = 0$ if $|A| = n \leq 2$. The resulting output will always have a value in $[0, 1]$, where 1 means a completely centralized network (star graph) and 0 means a completely decentralized network. It is important to emphasize that a centralization score of 0 does not necessarily mean that there is no collaboration and every developer is isolated. Rather it means that each developer is co-authoring with just as many authors, as the others do.

4.5.3 Clustering coefficient

While centralization helps us describe the centralness of the network and how much it is centered around a single, or a small number of developers, it does not help us describing the structure of the network in more detail. Our goal is to gain an understanding of also the modularity of our network, meaning how much developers tend to cluster together [17]. We expect that authors form smaller clusters, which are more tightly connected together, and these clusters have somewhat weaker ties to other clusters. This builds on the assumption that the social network of the software follows the modules which build up the software itself, thus authors of a specific function should also cluster together within the network [8, 17]. To measure this "clusteredness", we calculate the *local clustering coefficient* for each node.

The *local clustering coefficient* quantifies on a scale $[0, 1]$ how likely it is that a node's neighbours are also neighbours. We use the number of how many triangles (also called clique, triplet) is every node a part of. This is illustrated for unweighted networks on Figure 6 with the formula:

$$C_i = \frac{2T(i)}{\deg(i)(\deg(i) - 1)},$$

where $T(i)$ is the number of triangles through node i and $\deg(i)$ is the degree of i . However, in a weighted network we also have to consider the edge weights, since it is easy to see that a clustering coefficient of 1 with also the maximum weighted edges in a triplet does not represent the same clustering as being connected with a weak links. We expect weaker links connecting larger clusters, whereas stronger links within each cluster. Therefore, we use geometric averaging of the subgraph edge weights (as implemented by the `networkx`⁶ library [24]):

$$c_i = \frac{\sum_{jk} (\hat{w}_{ij} \hat{w}_{ik} \hat{w}_{jk})^{1/3}}{\deg(i)(\deg(i) - 1)}.$$

The \hat{w}_{ij} represents the normalized weight of edge e_{ij} over the maximum weight in the network.

4.5.4 Hierarchy

The degree centrality and the clustering coefficient are in themselves able to express meaningful aspects of the developer social network, however, by combining the two metrics, we can also assess how hierarchical the network is. In a scale-free social hierarchical network (such as the collaboration network), nodes tend to cluster around a single or a few hubs, which are more likely to have weak connections to other hubs [27, 17]. The nodes within these formed groups are relatively stronger than the connections connecting the hubs, but they are less likely to connect to nodes outside of their group. Therefore in a hierarchical network, the hubs have a high degree number and a low clustering coefficient, whereas the group members clustering around the hubs have a high clustering coefficient, but low degree numbers.

We can visualize the degree of hierarchy by plotting each node's clustering coefficient against the number of degrees, shown in Figure 7. In hierarchical networks, the plotted linear regression trendline will decrease steeply, as there is a negative correlation between the degree and clustering coefficient. In networks, where this cannot be observed, the trendline stays flat, meaning these two metrics are independent from each other and the structure is not hierarchical. To measure the hierarchical level numerically within a network, we take the trendline's slope, which is β_1 in the $y = \beta_1 x + \beta_0$ general linear regression equation.

⁶<https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.cluster.clustering.html>

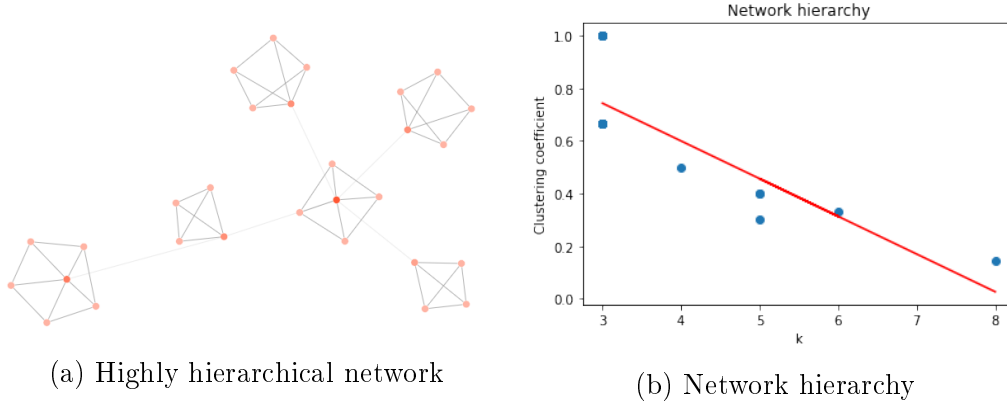


Figure 7: Hierarchical network and its corresponding degree number vs clustering coefficient plot.

In every network isolated authors can be observed, who are only working on files that noone else has edited (in the given timeframe). These nodes can skew the hierarchy score, because an isolated node’s degree and clustering coefficient are by definition 0, which disproportionately makes the trendline much flatter. Therefore, we remove the isolated nodes from the network. If the network only contains isolated nodes, and there is no linear regression to be calculated, we set the hierarchy value to 0.

4.6 Project measures

One of our assumptions is that the collaboration network structure changes depending on the project’s lifecycle. In order to discover cause and effect relationships between the network structure and project lifecycle, we gather basic project metrics to pinpoint the time and date of events, as well as the effort required within the project. We achieve this by gathering the dates of each release, and the quality and relative stress is measured by the issues within the project (create and close times).

4.6.1 Release and release measures

To measure the network changes around releases, first we collect the list of releases for the given project. Our goal is to gather the version number of each release, as well as the dates they were released. There are two relevant APIs regarding the release version numbers: GitHub releases and Git tags. Tags are marked and annotated commits supported by Git, therefore other projects, that are not on GitHub can also have tags. *Tags* are most

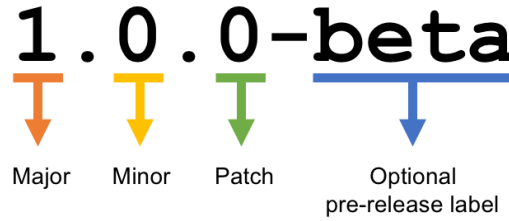


Figure 8: Semantic versioning. Figure source: [1].

commonly used to mark new release versions, but they can also be used to annotate other information, such as release editions or milestones. On the other hand, a *Release* is a high-level GitHub concept, which allows the project organizers to announce Git tags as project releases by adding a version number, release notes and binary artifacts [23]. A commit, that represents a new GitHub release, must be tagged, but a tagged commit does not necessarily need to be a new GitHub release.

Throughout our analysis, we use the tags as indicators of a new release, because there could be new versions of a software, that are not released publicly, and therefore do not have a GitHub release. Furthermore, GitHub only announced the Releases workflow in 2013, when the release versioning by Git tags was already a common practice. This means, that releases before 2013 can only be analysed via the repository tags.

Most large-scale OSS projects follow the semantic versioning convention, but only to a certain extent. The majority of tags follow the major-minor-patch (also known as breaking-feature-fix) semantic version naming convention in the format of X.Y.Z, where X is the major number, Y is the minor number, and Z is the patch number [25]. At the end, optionally pre-release and build data can be marked, for example: *1.11.6-pre*. The major number signifies an API-breaking change compared to the previous release, meaning backwards compatibility is not guaranteed for depending applications. Minor releases add new features, but compatibility is ensured with the older version. Patch releases usually handle bugs and security updates within the package.

During the network analysis, we expect the network measures to change around a release, but it is also expected, that a larger release, that required more collaboration, will have a greater impact, while smaller changes have smaller or no impact at all. The issue with the semantic versioning and re-

Index	Name	Tag name	Created	Type	Modifications	Lines added	Lines removed	Total change
0	Pandas v0.13.0	v0.13.0	2013-12-30 17:02:51	unknown	20031	1.665357e+09	1.624343e+09	3.289700e+09
1	Pandas v0.13.1	v0.13.1	2014-02-03 04:52:01	patch	836	2.231875e+06	2.349833e+06	4.581708e+06
2	Pandas v0.14rc1	v0.14.0rc1	2014-05-16 22:28:09	minor	1861	2.251189e+07	2.050846e+07	4.302035e+07
3	v0.14.0 final	v0.14.0	2014-05-30 11:47:40	unknown	318	3.220804e+06	5.314470e+05	3.752251e+06
4	v0.14.1 final	v0.14.1	2014-07-10 23:46:19	patch	720	1.660133e+06	3.796690e+05	2.039802e+06
5	v0.15.0 Pre-release	v0.15pre	2014-09-07 12:52:01	unknown	826	1.028779e+07	2.458664e+06	1.274645e+07
...
66	Pandas 1.1.5	v1.1.5	2020-12-07 11:42:10	patch	2167	4.251057e+06	4.611365e+06	8.862422e+06
67	Pandas 1.2.0rc0	v1.2.0rc0	2020-12-08 12:31:44	minor	41	1.060500e+04	8.850000e+02	1.149000e+04
68	Pandas 1.2.0	v1.2.0	2020-12-26 13:47:00	unknown	683	5.782570e+05	2.989310e+05	8.771880e+05
69	Pandas 1.2.1	v1.2.1	2021-01-20 11:21:02	patch	1306	1.864636e+07	9.228629e+07	1.109326e+08
70	Pandas 1.2.2	v1.2.2	2021-02-09 10:55:19	patch	844	6.783283e+06	1.827142e+07	2.505470e+07
71	Pandas 1.2.3	v1.2.3	2021-03-02 09:43:36	patch	959	1.002189e+07	5.783745e+06	1.580563e+07
72	Pandas 1.2.4	v1.2.4	2021-04-12 15:59:13	patch	289	3.402820e+05	1.704960e+05	5.107780e+05

Table 1: Releases collected information.

lease names is that there are no constraints, which would enforce a strict version naming, and it is entirely up to the developers to set the tag names. This leads to inconveniences when measuring collaboration effort through release version number, because a patch might require more collaboration effort than a minor release, and a minor release within one project could require significantly more teamwork than in another. Furthermore, tag names can be inconsistent, and there could be version numbers, that do not adhere to the major-minor-patch naming convention at all (e.g. test releases, or releases like 'latest-release-v11'). The possibility to add extra information at the end of tag names, like `-beta` or `rc` further complicates tracking the collaboration effort, because the majority of collaboration effort might happen before the `rc` version, or it might happen after it. As it can also be seen in Table 1, the `v0.14.0rc1` tag contains more modifications than the final `v0.14.0` version, whereas the `v1.2.0rc0` only contains a small fraction of modifications compared to `v1.2.0`.

In order to have a more fine-grained measure of how much effort a release required (besides the semantic versioning), we measure the number of lines added and lines removed in that release. This is calculated by adding up each commit's 'total lines added' and 'total lines removed', which was authored after the previous release but up to and including the current release tag's commit. The total change of lines for a release is simply the sum of lines added and lines removed, which are provided by the `git2net` miner. The miner also provides the number of modifications for each commit. A modification is a section of the source code modified, which can mean multiple lines added and deleted at the same part of the document. For example, if a new function is added to the project, which requires 20 new lines and removes 2 lines (e.g. empty space that was there before), then it will be considered as 1 modification, but 22 total line change. Changes to binary files are due to generated artifacts, which do not carry any collaboration effort, therefore

they are excluded, and commits, that do not have a hash are also removed.

The release type contains the semantic version of the release, which was gathered from the tag name with a Regular Expression matching the conventional versioning X.Y.Z. We also capture the version number in tags, that contain additional notations such as `beta` or `rc`, then we compare the current release to the previous to identify whether the release is a major, minor or patch release. When the found version numbers in the previous and in the current release are the same, we leave that as unknown, as this is mostly the case in pre-releases. The first release doesn't have a preceeding tag, therefore we consider every modification before the tag as part of the release. This leads to the first tag seeming to have significantly more edits than the rest of the releases, whereas in fact this is just the result of not tracking from the beginning (see example in Table 1). Therefore, in our analyses we remove the first release, as this would lead to falsely weighing the network results.

4.6.2 Project issues and measures

To measure the productivity within the project, we collect some basic information regarding the GitHub issues. Issues keep track of bugs, beatures and tasks, contributors can comment and discuss the task at hand within an issue, and it can be assigned to users and milestones can be set. We collect the following information of issues:

- Issue title
- Issue number
- Created at
- Closed at
- Open for
- Bug or feature

The *issue title* is a short description of the issue. Most projects create their own convention of naming and tagging issues, for example each issue starts with a 3-letter abbreviation of a category, e.g. `BUG` or `DOC`. *Issue number* is a unique number for each issue, that is increased sequentially. *Created at* is the date and time of the issue being created, and *closed at* is the time when it was closed. If the issue was still open on the day of data mining (May

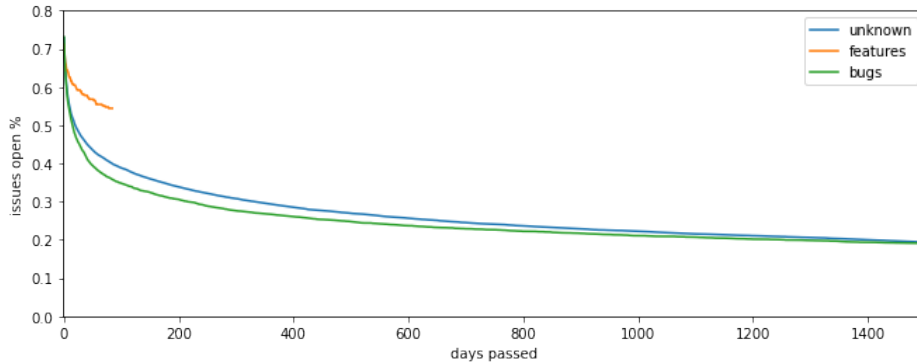


Figure 9: Survival curves of the `pandas` library.

9 2021), this field is empty. The *open for* field is the difference between the *closed at* and *created at* dates, or it is empty if the issue was never closed.

The last measure is *bug or feature*. Because issues can cover a wide range of possible topics, from discussions to performance, it is worth categorizing them to analyse the differences. However, categorization in practice proves to be difficult due to the different conventions each project uses, and because GitHub does not apply any constraints to the text of the title. We are searching for specific keywords in the *issue title* to categorize each issue into either a bug or a feature. We also drop some of the words from the text, because they would give us 'false positive' matches. If keywords for both or neither categories are found, the issue is marked as unknown. The dropped words, as well as the bug and feature keywords can be seen in Table 2. Although this simple method can already classify 10-15 percent of issues, it would require extensive manual effort to further improve this ratio. In future research, sentiment analysis of issue titles and issue descriptions could achieve much higher percentages.

Statistics for issues like average close time also prove to be difficult, because there are continuously open issues, which inevitably leads to the fact that a large portion of issues were still open and continue to be open. For these issues, we cannot know the issue close time, and if we try to calculate global measures by aggregating all (or portion) of issues, we have to consider the *survival bias* [14]. Figure 9 shows the issues survival curve of the `pandas` library. It is clearly visible, that a large portion (20 percent) of issues are not closed, and if we left them out of aggregate values such as average open time, our data will be skewed. Nevertheless, we can see in the example, that the categorization of issues makes a difference, as bugs get closed a bit earlier

Drop words	Bug keywords	Feature keywords
debug debugger	bug defect incorrect unexpected error missing warning problem	feature enhancement improvement suggestion wishlist wish list

Table 2: Keywords and drop words for *bug* and *feature* categorization.

than an average (unknown) issue, because after the same amount of days passign for both categories, more bugs tend to be closed in this example. In contrast, features have a high life expectancy, as they are more likely to be still open than a bug, if the same number of days pass. The reason could be that features take longer to develop and plan than bugs, which is consistent with the findings of Jarczky et. al. [14].

4.7 Time window

5 Collaboration pattern analysis

5.1 Observed projects and events

5.2 SNA metrics analysis

5.2.1 K-cores

Besides taking the 20th percentile of the top clustering coefficients to identify the number of core developers in the network, we [6]

5.3 Results

6 Quantitative analysis of projects during crunch time

6.1 Collaboration network changes

6.2 Prediction of outcome based on collaboration changes

7 Discussion and results

8 Conclusion and future work

References

- [1] Semantic versioning for UT. <https://forums.ubports.com/topic/1822/semantic-versioning-for-ut>, October 2018.
- [2] Shaosong Ou Alexander Hars. Working for Free? Motivations for Participating in Open-Source Projects. *International Journal of Electronic Commerce*, 6(3):25–39, April 2002.
- [3] Mohammed Aljemabi and Zhongjie Wang. Empirical Study on the Evolution of Developer Social Networks. *IEEE Access*, PP:1–1, September 2018.
- [4] M. Antwerp. Evolution of open source software networks. pages 25–39, January 2010.
- [5] Guilherme Avelino, Leonardo Passos, Andre Hora, and Marco Tulio Valente. A Novel Approach for Estimating Truck Factors. *2016 IEEE 24th International Conference on Program Comprehension (ICPC)*, pages 1–10, May 2016.
- [6] V. Batagelj and M. Zaversnik. An $O(m)$ Algorithm for Cores Decomposition of Networks. *arXiv:cs/0310049*, October 2003.
- [7] 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.
- [8] Melvin E Conway. How Do Committees Invent? page 4, 1968.
 - [9] Kevin Crowston and James Howison. The social structure of free and open source software development. <https://firstmonday.org/ojs/index.php/fm/article/download/1478/1393?inline=1>, February 2005.
 - [10] Kevin Crowston and James Howison. Hierarchy and centralization in free and open source software team communications. *Knowledge, Technology & Policy*, 18(4):65–85, December 2006.
 - [11] Ikram El Asri, Nouredine Kerzazi, Lamia Benhiba, and Mohammed Janati. From Periphery to Core: A Temporal Analysis of GitHub Contributors' Collaboration Network. In Luis M. Camarinha-Matos, Hamideh Afsarmanesh, and Rosanna Fornasiero, editors, *Collaboration in a Data-Rich World*, IFIP Advances in Information and Communication Technology, pages 217–229, Cham, 2017. Springer International Publishing.
 - [12] Christoph Gote, Ingo Scholtes, and Frank Schweitzer. Analysing Time-Stamped Co-Editing Networks in Software Development Teams using git2net. *arXiv:1911.09484 [physics]*, November 2019.
 - [13] Christoph Gote and Christian Zingg. Gambit – An Open Source Name Disambiguation Tool for Version Control Systems. *arXiv:2103.05666 [physics]*, March 2021.
 - [14] Oskar Jarczyk, Szymon Jaroszewicz, Adam Wierzbicki, Kamil Pawlak, and Michal Jankowski-Lorek. Surgical teams on GitHub: Modeling performance of GitHub project development processes. *Information and Software Technology*, 100:32–46, August 2018.
 - [15] Arkadiusz Jędrzejewski. *The Role of Complex Networks in Agent-Based Computational Economics*. PhD thesis, June 2016.
 - [16] Mitchell Joblin, Sven Apel, Claus Hunsen, and Wolfgang Maurer. Classifying Developers into Core and Peripheral: An Empirical Study on Count and Network Metrics. *arXiv:1604.00830 [cs]*, April 2016.
 - [17] Mitchell Joblin, Sven Apel, and Wolfgang Maurer. Evolutionary trends of developer coordination: A network approach. *Empirical Software Engineering*, 22(4):2050–2094, August 2017.

- [18] Yu-Ru Lin, Hari Sundaram, Yun Chi, Junichi Tatemura, and Belle Tseng. Blog Community Discovery and Evolution Based on Mutual Awareness Expansion. In *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence*, pages 48–56, November 2007.
- [19] 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.
- [20] M. R. Martínez-Torres. A genetic search of patterns of behaviour in OSS communities. *Expert Systems with Applications*, 39(18):13182–13192, December 2012.
- [21] Kelvin McClean, Des Greer, and Anna Jurek-Loughrey. Social network analysis of open source software: A review and categorisation. *Information and Software Technology*, 130:106442, February 2021.
- [22] Audris Mockus, Roy T. Fielding, and James D. Herbsleb. Two case studies of open source software development: Apache and Mozilla. *ACM Transactions on Software Engineering and Methodology*, 11(3):309–346, July 2002.
- [23] Rick Olson. Release Your Software, July 2013.
- [24] Jukka-Pekka Onnela, Jari Saramäki, János Kertész, and Kimmo Kaski. Intensity and coherence of motifs in weighted complex networks. *Physical Review E*, 71(6):065103, June 2005.
- [25] Tom Preston-Werner. Semantic Versioning 2.0.0. <https://semver.org/>.
- [26] PwC. Leading benefits of open-source software among enterprises worldwide as of 2016. *Statista*, 2016.
- [27] Erzsébet Ravasz and Albert-László Barabási. Hierarchical organization in complex networks. *Physical Review E*, 67(2):026112, February 2003.
- [28] Rotem Stram, Pascal Reuss, and Klaus-Dieter Althoff. Weighted One Mode Projection of a Bipartite Graph as a Local Similarity Measure. pages 375–389, June 2017.

- [29] 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.
- [30] Xuan Yang, Daning Hu, and Davison M. Robert. How Microblogging Networks Affect Project Success of Open Source Software Development. In *2013 46th Hawaii International Conference on System Sciences*, pages 3178–3186, January 2013.
- [31] Yunwen Ye and K. Kishida. Toward an understanding of the motivation of open source software developers. In *25th International Conference on Software Engineering, 2003. Proceedings.*, pages 419–429, May 2003.