



## AN ABSTRACT OF THE DISSERTATION OF

Amirhosein “Emerson” Azarbakht for the degree of Doctor of Philosophy in  
Computer Science presented on May 18, 2017.

Title: Longitudinal Analysis of Collaboration in Forked Open Source Software  
Development Projects

Abstract approved: \_\_\_\_\_

Carlos Jensen

Social interactions are a ubiquitous part of our lives, and the creation of online social communities has been a natural extension of this phenomena. Free and Open Source Software (FOSS) development efforts are prime examples of how communities can be leveraged in software development, where groups are formed around communities of interest, and depend on continued interest and involvement.

Not everything works smoothly all the time in open source projects. Problems arise for a variety of reasons, including collaboration and communication problems, which results in uncertainty about the operational health and survivability of the projects. Many stake-holders are affected by this uncertainty, including industry sponsors, individual contributors, corporate developers, and users, who all have decided to invest time and effort in the project, and will be affected if a project suffers from troubles.

Forking in FOSS, either as a non-friendly split or a friendly divide, affects the community. Such effects have been studied, shedding light on how forking

happens. However, most existing research on forking is post-hoc. In this study, we focus on the seldom-studied run-up to forking events.

We used the following two approaches to study the evolution and social dynamics of FOSS communities; 1) Time series analysis of the contents of the messages sent and received on the projects developers mailing list, for the time period of 10-month run-up to the fork was analyzed for anomalies, indicative of simmering conflicts. 2) Social network analysis using a developer-oriented approach to statistically model the changes a community goes through in the run-up to a fork, in which the model represents tie formation, tie breakage, and tie maintenance between developers. We estimated several model parameters that capture the variance in the changes the community goes through. We found that conflict-driven forks exhibited anomalies; time series analysis of sentiments showed the anomalies occurred before and close to the fork event. Whereas non-conflict-driven forks did not suffer from such pre-fork anomalies. The objective was to be able to evaluate the operational health of the project community, and intervene if need be. We suggest anomaly detection of the time series analysis may be used by the project stakeholders/investors as key indicators left in the record, that can be used to identify problems among developers, and intervene if need be.

We also found that in conflict-driven forks, (1) the developers maintained a preference for interacting with developers who had similar out-degrees, in contrast to the non-conflict-driven forks, where the developers did not require similar out-degrees. The interpretation may be that a project with non-conflict-driven forks had a more inclusive and classless core developer team. (2) The interactions were reciprocal, in contrast to the non-conflict-driven forks, where the interactions did not need to be reciprocal to happen. The interpretation

may be that the projects with non-conflict-driven forks were more open to interactions whether or not they would get something back in return from the other developer. (3) In conflict-driven forks, the more senior developers preferred to interact with other more senior developers, in contrast to the non-conflict-driven forks. The interpretation may be that the senior developers in projects with conflict-driven forks were less involved with junior developers than in projects with non-conflict-driven forks. (4) In non-conflict-driven forks, the developers with high source code contribution levels interacted more with other high source code contributors. (5) In non-conflict-driven forks, high levels of contribution to the source code brings you connections more rapidly, while high levels of contributions to the mailing list is not suggestive of this. This can be interpreted as a sign of meritocracy based on code, rather than talk, which captures a healthy dynamic in these projects.

©Copyright by Amirhosein “Emerson” Azarbakht

May 18, 2017

All Rights Reserved

Longitudinal Analysis of Collaboration in Forked Open Source Software  
Development Projects

by  
Amirhosein “Emerson” Azarbakht

A DISSERTATION

submitted to

Oregon State University

in partial fulfillment of  
the requirements for the  
degree of

Doctor of Philosophy

Presented May 18, 2017  
Commencement June 2017

Doctor of Philosophy dissertation of Amirhosein “Emerson” Azarbakht presented  
on May 18, 2017

APPROVED:

---

Major Professor, representing Computer Science

---

Director of the School of Electrical Engineering and Computer Science

---

Dean of the Graduate School

I understand that my dissertation will become part of the permanent collection of Oregon State University libraries. My signature below authorizes release of my dissertation to any reader upon request.

---

Amirhosein “Emerson” Azarbakht, Author

## ACKNOWLEDGMENTS

I would like to express sincere appreciation to family, friends, colleagues, and committee members.



## TABLE OF CONTENTS

	<u>Page</u>
1 Introduction.....	1
2 Related Work .....	7
3 Methodology and Results.....	11
3.1 Time series analysis of open source developers collabo- ration communication sentiments .....	12
3.2 Social network analysis of open source developers inter- action graphs .....	34
4 Discussion.....	65
5 Threats to Validity.....	69
Appendices .....	71
A Appendix A: List of projects forked for undesirable reasons ....	71
B Appendix B: Initial study: Temporal analysis using the network- specific measurement approach .....	73
B.1 Visualization .....	75
B.2 Initial study results and discussion.....	76
C Appendix C: Mathematical Definition of Effects in the Statistical Model.....	84
C.1 Structural Effects for the Objective Function .....	84
C.2 Behavior-related Effects.....	86
Bibliography .....	88

## LIST OF FIGURES

<u>Figure</u>	<u>Page</u>
1 Time Series of sentiments in developers' messages for projects (From top to bottom, left to right) Kamilio (Conflict-Driven), ffmpeg (Conflict-Driven), Amarok (Non-Conflict-Driven), Apache CouchDB (Non-Conflict-Driven), Pidgin (Non-Conflict-Driven), MPlayer (Non-Conflict-Driven). It is not easy to visually analyse these raw time series.....	19
2 Time Series of sentiments in developers' messages for projects (From top to bottom, left to right) Asterisk (Non-Conflict-Driven), freglut (Non-Conflict-Driven), Ceph (Not Forked), OpenStack Neutron (Not Forked), and GlusterFS (Not Forked). It is not easy to visually analyse these raw time series.....	20
3 Time series decomposed into trend, seasonality and residual time Series for projects (From top to bottom, left to right) Kamilio (Conflict-Driven), ffmpeg (Conflict-Driven), Amarok (Non-Conflict-Driven), Apache CouchDB (Non-Conflict-Driven).	21
4 Time series decomposed into trend, seasonality and residual time Series for projects (From top to bottom, left to right) Pidgin (Non-Conflict-Driven), MPlayer (Non-Conflict-Driven), Asterisk (Non-Conflict-Driven), and freglut (Non-Conflict-Driven).	22
5 Time series decomposed into trend, seasonality and residual time Series for projects (From top to bottom, left to right) Ceph (Not Forked), OpenStack Neutron (Not Forked), and GlusterFS (Not Forked). ....	23
6 Time series with outliers detected for project Kamilio (conflict-driven).....	24
7 Time series with outliers detected for project ffmpeg (conflict-driven).....	25
8 Time series with outliers detected for project Amarok (non-conflict-driven). ....	26
9 Time series with outliers detected for project Pidgin (non-conflict-driven). ....	27
10 Time series with outliers detected for project MPlayer (non-conflict-driven). ....	28

## LIST OF FIGURES (Continued)

<u>Figure</u>	<u>Page</u>
11 Time series with outliers detected for project Asterisk (non-conflict-driven). . . . .	29
12 Time series with outliers detected for project freeglut (non-conflict-driven). . . . .	30
13 Time series with outliers detected for project Ceph (Not forked). . . . .	31
14 Time series with outliers detected for project OpenStack Neutron (Not forked). . . . .	32
15 Time series with outliers detected for project GlusterFS (Not forked). . . . .	33
16 The methodology in a glance . . . . .	35
17 Heat-map color-coded examples Degree, Closeness, Betweenness, Eigenvector centralities . . . . .	74
18 The number of nodes and edges over time, as network-specific measurements . . . . .	77
19 The network diameter change over time, as network-specific measurements . . . . .	78
20 Snapshots from video visualization of Kamilio project's graph . . . . .	79
21 Kamilio top contributors' betweenness centralities and network diameter over time . . . . .	79
22 Amarak project's top contributors' betweenness centralities and network diameter over time . . . . .	81
23 Asterisk project's top contributors' betweenness centralities and network diameter over time . . . . .	82

## LIST OF TABLES

<u>Table</u>	<u>Page</u>
1 The main reasons for forking .....	4
2 The behavioral measures used by Asur et al. [1] .....	9
3 The measures of diversity used by Kunegis et al. [31] .....	10
4 Projects in conflict-driven forks and non-conflict driven forks categories selected and analyzed based on the criteria described in section 3.2.1 taken from [47] .....	13
5 The number of messages sentiment-analyzed for each project. ..	15
6 Anomalies detected in the time series of sentiments for project Kamailio (Conflict-Driven). Type column refers to the differ- ent anomalies as described in section 3.1.3. Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 128.8: log likeli- hood = -1170.2, aic = 2356.41) .....	24
7 Anomalies detected in the time series of sentiments for project ffmpeg (Conflict-Driven). Note that all listed outliers are statis- tically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 57.33: log likelihood = -1043.78, aic = 2101.57) .....	25
8 Anomalies detected in the time series of sentiments for project Amarok (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 375.6: log likelihood = -1326.61, aic = 2661.22) .....	26
9 Anomalies detected in the time series of sentiments for project Pidgin (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 312.3: log likelihood = -1307.25, aic = 2622.49) .....	27
10 Anomalies detected in the time series of sentiments for project MPlayer (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 225.5: log likelihood = -1251.94, aic = 2517.88) .....	28

## LIST OF TABLES (Continued)

<u>Table</u>	<u>Page</u>
11 Anomalies detected in the time series of sentiments for project Asterisk (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 358.2: log likelihood = -1227.85, aic = 2463.7 ) .....	29
12 Anomalies detected in the time series of sentiments for project freeglut (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$ estimated as 348.6: log likelihood = -1307.6, aic = 2619.2) .....	30
13 Anomalies detected in the time series of sentiments for project Ceph (Not forked). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. .	31
14 Anomalies detected in the time series of sentiments for project OpenStack Neutron (Not forked). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ....	32
15 Anomalies detected in the time series of sentiments for project GlusterFS (Not forked). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value.....	33
16 Estimates parameters in the model for Project Kamilio (Conflict-Driven). The outdegree density, transitive triplets, out-out-degree assortativity, and developers' seniority effects are statistically significant. The interpretations are discussed in section 4.....	55
17 Estimates parameters in the model for Project ffmpeg (Conflict-Driven). The outdegree density, transitive triplets, reciprocity, 3-cycles and out-out-degree assortativity are statistically significant. The interpretations are discussed in section 4. ....	56
18 Estimates parameters in the model for Project Amarok (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4. ....	57

## LIST OF TABLES (Continued)

<u>Table</u>	<u>Page</u>
19 Estimates parameters in the model for Project Apache CouchDB (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4. ....	58
20 Estimates parameters in the model for Project Pidgin (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4. ....	59
21 Estimates parameters in the model for Project MPlayer (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4. ....	60
22 Estimates parameters in the model for Project freeglut (Non-Conflict-Driven). The outdegree density, transitive triplets effects are statistically significant. The interpretations are discussed in section 4. ....	61
23 Estimates parameters in the model for Project Ceph (Not forked). The outdegree density, reciprocity, transitive triplets, 3-cycles, out-out degree assortativity and developers' alter source code activity effects are statistically significant. The interpretations are discussed in section 4. ....	62
24 Estimates parameters in the model for Project OpenStack Neutron (Not forked). The outdegree density, reciprocity, transitive triplets, 3-cycles, out-out degree assortativity and developers' seniority effects are statistically significant. The interpretations are discussed in section 4. ....	63
25 Estimates parameters in the model for Project GlusterFS (Not forked). The outdegree density, transitive triplets, out-out degree assortativity and developers' alter seniority effects are statistically significant. The interpretations are discussed in section 4. ....	64
26 Conflict-Driven vs. Non-Conflict-Driven parameter effects .....	65
27 List of all projects forked because of “ <i>personal differences among the developer team</i> ” conflict-driven .....	71

## LIST OF TABLES (Continued)

<u>Table</u>	<u>Page</u>
28 List of all projects forked because of the need for “ <i>more community-driven development</i> ” conflict-driven .....	72

# Longitudinal Analysis of Collaboration in Forked Open Source Software Development Projects

## 1 Introduction

Social networks are a ubiquitous part of our social lives, and the creation of online social communities has been a natural extension of this phenomena. Social media plays an important role in software engineering, as software developers use them to communicate, learn, collaborate and coordinate with others [56]. Free and Open Source Software (FOSS) development efforts are prime examples of how community can be leveraged in software development, where groups are formed around communities of interest, and depend on continued interest and involvement to stay alive [39].

Community splits in free and open source software development are referred to as forks, and are relatively common. Robles et al. [47] define forking as “when a part of a development community (or a third party not related to the project) starts a completely independent line of development based on the source code basis of the project.”

Although the bulk of collaboration and communication in FOSS communities occurs online and is publicly accessible by researchers, there are still many open questions about the social dynamics in FOSS communities. Projects may go through a metamorphosis when faced with an influx of new developers or the involvement of an outside organization. Conflicts between developers’ divergent



visions about the future of the project may lead to forking of the project and dilution of the community. Forking, either as an acrimonious split when there is a conflict, or as a friendly divide when new features are experimentally added, affect the community [10].

Previous research on forking ranges from the study by Robles et al. [47] that identified 220 significant FOSS projects that have forked over the past 30 years, and compiled a comprehensive list of the dates and reasons for forking (listed in Table 1, and depicted by frequency in Figure 1), to the study by Baishakhi et al. [8] on post-forking porting of new features or bug fixes from peer projects. It encompasses works of Nyman on developers’ opinions about forking [41], developers motivations for performing forks [36], the necessity of code forking as tool for sustainability [40], and Syeed’s work on sociotechnical dependencies in the BSD projects family [57].

Most existing research on forking, however, is post-hoc. It looks at the forking events in retrospect and tries to find the outcome of the fork; what happened after the fork occurred; what was the cause of forking, and such. The run-up to the forking events are seldom studied. This leaves several questions unanswered: Was it a long-term trend? Was the community polarized before the forking happened? Was there a shift in influence? Did the center of gravity of the community change? What was the tipping point? Was it predictable? Is it ever predictable? We are missing that context.

Additionally, studies of FOSS communities tend to suffer from an important limitation. They treat community as a static structure rather than a dynamic process. Longitudinal studies on open source forking are rare. To better understand and measure the evolution, social dynamics of forked FOSS projects, and integral components to understanding their evolution and direction, we need

new and better tools. Before making such new tools, we need to gain a better understanding of the context. With this knowledge and these tools, we could help projects reflect on their actions, and help community leaders make informed decisions about possible changes or interventions. This could also help potential sponsors make informed decisions when investing in a project, and throughout their involvement to ensure a sustainable engagement.

Identification is the first step to rectify an undesirable dynamic before the damage is done. A community that does not manage growing pains may end up stagnating or dissolving. Managing growing pains is especially important in the case of FOSS projects, where near half the project contributors are volunteers [21]. Oh et al. [42] have argued that openness in FOSS is

*“[...] generally perceived as having a positive connotation, however, the term can also be interpreted as referring to some nonconstructive characteristics, such as unobstructed exit, susceptible, vulnerable, fragile, lacking effective regulation, and so on. The unobstructed exit and lack of regulatory force inherent in the FOSS community can result in a community’s susceptibility and vulnerability to herded exits by its participants. Commercial vendor intervention, an alternative project becoming available, and licensing issues can result in some original core members ceasing to provide their loyal service for the community, which can prompt their coworkers to leave as well” [42].*

Identification of recipes for success or stagnation, sustainability or fragmentation may lead to a set of best practices and pitfalls.

In this research, we use time series analysis and longitudinal social network analysis to study the evolution and social dynamics of FOSS communities. Specifically, we use anomaly detection algorithms to analyze the time series of sentiments in communications of open source developers, in the run up to the

fork, and we also use longitudinal social network analysis to investigate the driving forces in formation and dissolution of communities. With these techniques we aim to identify measures associated with unhealthy group dynamics, for example a simmering conflict, in addition to early indicators of major events in the lifespan of a community. One set of dynamics we were especially interested in, are those that are associated with FOSS projects that have forked.

TABLE 1: Robles and Gonzalez-Barahona [47] defined the following reasons for forking.

<b>Reason for forking</b>	<b>Example forks</b>
Technical (Addition of functionality)	Amarok & Clementine Player
More community-driven development	Asterisk & Callweaver
Differences among developer team	Kamailio & OpenSIPS
Discontinuation of the original project	Apache web server
Commercial strategy forks	LibreOffice & OpenOffice.org
Experimental	GCC & EGCS
Legal issues	X.Org & XFree

## Research Goals

Social interactions reflect the changes the community goes through, and so, it can be used to describe the context surrounding a forking event. Social interactions in FOSS can happen, for example, in the form of mailing list email correspondence, and source code co-authoring.

For the rest of the discussion, we define conflict-driven forks as projects that have forked for personal differences, and non-conflict driven forks as projects that have forked for technical experimentation, and were not forked for personal differences.

We defined conflict-driven (undesirable) as the projects forked because of “Personal differences among developers team”, or because of the need for “more community-driven development”. These situations are undesirable because they imply an increase in cost of maintenance, redundant or wasted efforts, and the loss of shared values. One-fifth (20.5%) of the 220 forked projects studied by [47] fall into this category. Examples include Kamailio and OpenSIPS projects. We defined non-conflict-driven as projects forked for “technical differences (Addition of functionality)”, and were not forked for personal differences. More than a quarter (27.3%) of the 220 forked projects studied by [47] fall into this category. Examples include AmaroK and Clementine Player projects. Conflict-driven forks are especially important, as they are destructive, lead to resentment, and the dilution of the community workforce. In contrast, the non-conflict driven forks are friendly and often temporary splits, where developers experimentally modify the code and/or add new features, often with the intention of returning back to the original project and merging the new features back into the original project. In summary, conflict-driven forks are like violent exits with no intention of coming back, whereas non-conflict-driven forks are temporary trips.

In this study, we analyze, quantify and visualize how the community is structured, how it evolves, and the degree to which community involvement changes over time.

Specifically, our overall research objective is to identify these traces/social patterns associated with different types of conflict-driven forking.

**Do forks leave traces in the collaboration artifacts of open source projects in the period leading up to the fork?**

To study the properties of possible social patterns, we need to verify their existence. More specifically, we need to check whether the possible social patterns are manifested in the the collaboration artifacts of open source projects, e.g., mailing list communications data, source code data. This is going to be accomplished as explained in section 3.

**Do conflict-driven forks and non-conflict-driven forks leave different types of traces?**

If forks leave traces in the collaboration artifacts, do forks exhibit different social patterns? Are there patterns that exemplify these categories? For example, is there a distinct “conflict-driven” fork collaboration pattern? If so, do different forking reasons have distinctly different social patterns associated with them? We are going to investigate this by statistical modeling of the interaction graphs, as explained in detail in section 3.2.2.

**What are the key indicators that let us distinguish between conflict-driven and non-conflict-driven forks?**

What quantitative measure(s) can be used as an early warning sign of an inflection point (fork)? Are there metrics that can be used to monitor the odds of change, (e.g. forking-related patterns), ahead of time? This will be accomplished by time series analysis and statistical modeling of developer interactions as explained in more detail in section 3.

This dissertation is organized as follows: Section 2 presents related literature on open source social communities, the gap in the literature, and discusses why the issue needs to be studied. Then, in section 3 presents our methodology, including time series and statistical modeling. In section 3.2.5 we present the results of our analysis. In Section 4 we discuss the findings and their interpretations and as well their implications. Lastly, in section 5 the threats to validity are discussed.

## 2 Related Work

The free and open source software development communities have been studied extensively. Researchers have studied the social structure and dynamics of team communications [11][23][27][28][35], identifying knowledge brokers and associated activities [53], project sustainability [35][40], forking [39], requirement satisfaction [18], their topology [11], their demographic diversity [31], gender differences in the process of joining them [30], and the role of age and the core team in their communities [2][3][17][59]. Most of these studies have tended to look at community as a static structure rather than a dynamic process [16]. This makes it hard to determine cause and effect, or the exact impact of social changes.

Post-forking porting of new features or bug fixes from peer projects happens among forked projects [8]. A case study of the BSD family (i.e., FreeBSD, OpenBSD, and NetBSD, which evolved from the same code base) found that 10-15% of lines in BSD release patches consist of ported edits, and on average 26-58% of active developers take part in porting per release. Additionally, They found that over 50% of ported changes propagate to other projects within three

releases [8]. This shows the amount of redundant work developers need to do to synchronize and keep up with development in parallel projects.

Visual exploration of the collaboration networks in FOSS communities was the focus of a study that aimed to observe how key events in the mobile-device industry affected the WebKit collaboration network over its lifetime [58]. They found that *coopetition* (both competition and collaboration) exists in the open source community; moreover, they observed that the “firms that played a more central role in the WebKit project such as Google, Apple and Samsung were by 2013 the leaders of the mobile-devices industry. Whereas more peripheral firms such as RIM and Nokia lost market-share” [58].

The study of communities has grown in popularity in part thanks to advances in social network analysis. From the earliest works by Zachary [60] to the more recent works of Leskovec et al. [32][33], there is a growing body of quantitative research on online communities. The earliest works on communities was done with a focus on information diffusion in a community [60]. The study by Zachary investigated the fission of a community; the process of communities splitting into two or more parts. They found that fission could be predicted by applying the Ford-Fulkerson min-cut algorithm [20] on the group’s communication graph; “the unequal flow of sentiments across the ties” and discriminatory sharing of information lead to subcommunities with more internal stability than the community as a whole [60].

The dynamic behavior of a network and identifying key events was the aim of a study by Asur et al [1]. They studied three DBLP co-authorship networks and defined the evolution of these networks as following one of these paths: a) Continue, b) k-Merge, c) k-Split, d) Form, or e) Dissolve. They used the metrics listed in Table 2 and defined four possible transformation events for individual

members: 1) Appear, 2) Disappear, 3) Join, and 4) Leave. They compared groups extracted from consecutive snapshots, based on the size and overlap of every pair of groups. Then, they labeled groups with events, and used these identified events [1].

TABLE 2: The behavioral measures used by Asur et al. [1]

<b>Metrics</b>	<b>Meaning</b>
Stability	Tendency of a node to have interactions with the same nodes over time
Sociability	Tendency of a node to have different interactions
Influence	Number of followers a node has on a network and how its actions are copied and/or followed by other nodes. (e.g., when it joins/leaves a conversation, many other nodes join/leave the conversation, too)
Popularity	Number of nodes in a cluster (how crowded a sub-community is)

The communication patterns of free and open source software developers in a bug repository were examined by Howison et al. [27]. They calculated out-degree centrality as their metric. Out-degree centrality measures the proportion of times a node contacted other nodes (outgoing) over how many times it was contacted by other nodes (incoming). They calculated this centrality over time “in 90-day windows, moving the window forward 30 days at a time.” They found that “while change at the center of FOSS projects is relatively uncommon,” participation across the community is highly skewed, following a power-law distribution, where many participants appear for a short period of time, and a very small number of participants are at the center for long periods. Our proposed approach is similar to theirs in how we form collaboration graphs.



TABLE 3: The measures of diversity used by Kunegis et al. [31]

<b>Network property</b>	<b>Network is diverse when</b>	<b>Diversity Measures</b>
Paths between nodes	Paths are long	Effective diameter
Degrees of nodes	Degrees are equal	Gini coefficient of the degree distribution
Communities	Communities have similar sizes	Fractional rank of the adjacency matrix
Random walks	Random walks have high probability of return	Weighted spectral distribution
Control of nodes	Nodes are hard to control	Number of driver nodes

Our approach is different in terms of our project selection criteria, the metrics we examine, and our research questions.

The tension between diversity and homogeneity in a community was studied by Kunegis et al. [31]. They defined five network statistics, listed in Table 3, used to examine the evolution of large-scale networks over time. They found that except for the diameter, all other measures of diversity shrunk as the networks matured over their lifespan. Kunegis et al. [31] argued that one possible reason could be that the community structure consolidates as projects mature.

Community dynamics was the focus of a more recent study by Hanne-  
mann and Klamann [24] on three open source bioinformatics communities. They measured "age" of users, as starting from their first activity and found survival rates and two indicators for significant changes in the core of the community.

They identified a survival rate pattern of 20-40-90%, meaning that only 20% of the newcomers survived after their first year, 40% of these survivors made it through the second year, and 90% of the remaining ones, survived over the next years. As for the change in the core, they suggested that a falling maximum betweenness in combination with an increasing network diameter as an indicator for a significant change in the core, e.g., retirement of a central person in the community. Our initial network-specific study built on their findings, and the evolution of betweenness centralities and network diameters for the projects in our study are explained in the following sections.

### 3 Methodology and Results

We did the following two types of analysis:

- Time series analysis of open source developers collaboration communication sentiments
- Social network analysis of open source developers interaction graphs

In the following sections, the details of each method as well as the results of each method are described.

One of the challenges of studying the run-up to forking that you face is finding projects that have gone through a variety of different types of forking. We also need for those projects to have survived, and additionally for the records to have survived. Unfortunately, that leaves us with a small sample, which is typical of studies in this domain, even though our study features more projects that is typical of related work as discussed in section 2, which typically feature only 1-3 projects [8][58][24]. Table 4 lists all the attainable projects that meet

the criteria. Of the projects listed in Table 4, projects fall into three categories; conflict-driven, non-conflict-driven, and not forked. To find these projects, we looked at the list of all significant open source software forks in the past three decades as compiled by Robles and Gonzalez-Barahona [47]. Their study found the reasons behind each fork, listed in Table 1. We applied three selection criteria to the 220 forked projects on that list to find projects in conflict-driven and non-conflict-driven categories. A project was short-listed as either a conflict-driven or non-conflict-driven if the forking was relatively recent, i.e., happened after the year 2000, its data was existent and available to access and download online, or was made accessible to us after our requests; and the project had a sizable developer community, i.e., more than a dozen developers, which means it would be large enough that requires a meaningful statistical analysis. For the Not Forked category, we chose well-known, stable projects that had been around for a while (two years), and had large communities; and had not forked; and were similar in size of the development team. The preceding criteria resulted in the projects listed in Table 4.

### **3.1 Time series analysis of open source developers collaboration communication sentiments**

Time series analysis is a statistical technique used to understand the past, and predict the future. It is also used to do forecasting (predicting inference, a subset of statistical inference), which assumes that present trends continue. This assumption cannot be checked empirically, but, when we identify the likely causes for a trend, we can justify the forecasting (extrapolating it) for a few time-steps at least. Time series analysis is also used for anomaly detection, and classification (i.e. assigning a time series pattern to a specific category: e.g. gesture recognition of hand movements in sign language videos) and query by

TABLE 4: Projects in conflict-driven forks and non-conflict driven forks categories selected and analyzed based on the criteria described in section 3.2.1 taken from [47]

<b>Projects</b>	<b>Reason for forking</b>	<b>Year forked</b>
Kamailio & OpenSIPS	Conflict-driven	2008
ffmpeg & libav	Conflict-driven	2011
Asterisk & Callweaver	Non-conflict-driven	2007
freelut & OpenGLUT	Non-conflict-driven	2004
Amarok & Clementine Player	Non-conflict-driven	2010
Apache CouchDB & BigCouch	Non-conflict-driven	2010
Pidgin & Carrier	Non-conflict-driven	2008
MPlayer & MPlayerXP	Non-conflict-driven	2005
Ceph	Not forked	NA
OpenStack Neutron	Not forked	NA
GlusterFS	Not forked	NA

content, i.e. content-based image retrieval.

Time series analysis is a technique with applications in a variety of fields. Examples include census analysis, business forecasting, disease incidence tracking, understanding fluctuations in business sales, airline’s decision to buy airplanes because of passenger trends and decision to increase/maintain market share, and monitoring unemployment rate as an economic indicator used by decision makers.

Sentiment Analysis on the other hand, is a technique to determine the overall contextual polarity or emotional reaction to a document, interaction, or event, or the attitude of a writer. It uses natural language processing, statistics, or machine learning methods to extract, identify, or otherwise characterize the sentiment content of a text unit.

Sentiment analysis has applications in many fields. It is used for bias identification in news sources, identifying (in)appropriate content for ad placement, “Flame” detection, analyzing trends, targeting advertising/messages and gauging reactions, and identifying ideological bias.

### **3.1.1 Data**

A statistical time series analysis of the sentiments for the contents of the messages sent and received by the developers in the 10-month run-up to the fork was completed. The messages were collected from the projects’ mailing list messages, and contained all messages by all developers for that time period. These projects had on average 9,166 messages in this time period. Table 5 lists the number of messages for each project in this time period.

TABLE 5: The number of messages sentiment-analyzed for each project.

<b>Projects</b>	<b>Number of messages</b>	<b>Year forked</b>
Kamailio & OpenSIPS	5,450	2008
ffmpeg & libav	20,690	2011
Amarok & Clementine Player	1,642	2010
Apache CouchDB & BigCouch	3,746	2010
Pidgin & Carrier	4,095	2008
MPlayer & MPlayerXP	6,305	2005
Asterisk & Callweaver	28,012	2007
freeglut & OpenGLUT	2,267	2004
Ceph	7,232	NA
OpenStack Neutron	20,153	NA
GlusterFS	1,245	NA

### 3.1.2 Analysis

Our goal was to find out possible pre-fork negative anomalies happened, that could be an early indicator of problems in the community.

We used the R (Statistical Computing) [45] package SentimentAnalysis [19][44] to analyze the sentiments of the cleaned data, and each message was assigned a sentiment score between  $[-1, +1]$ , with negative values indicating negative sentiments, 0 for neutral, and positive value representing positive sentiments.

To analyze the raw time series data, we decomposed the raw series into the trend, seasonality, and residual components. The trend (a non-periodic systematic change in the time series) and residual series is the component we are interested in, especially the residual. The residual series is the part that can

give us insights into the longitudinal interdependency between the observations. The trend and seasonality are the long-time non-stationary process that needs to be removed before doing the time series analysis.

To detect outliers, we used the R (Statistical Computing) [45] package `tsoutliers` [34] for outlier detection by estimating outlier effects for four types of outliers. We found the outliers and fitted models to the series with the outlier effects removed. Failing to adjust for outliers can result in wrong models or biased parameter estimates, and increased forecasting error. The four models for outlier effect that we looked for were Additive outlier (AO), Level shift (LS), Temporary change (TC), and Innovational outlier (IO).

Figures 6-15 show the detected outliers in the sentiments time series, and includes both the original and adjusted time series, as well as the outlier effects.

### 3.1.3 Outlier Detection Models

We were interested in detecting anomalies that might have happened before the forking event, that could be used as indicators associated with forking. For this purpose, we used the following anomaly detection approach. All the following description from [14] describes outlier models.

For an ARIMA( $p, d, q$ ) process, we have

$$X_t = \frac{\theta(B)}{\alpha(B)\phi(B)} Z_t \quad (0.1)$$

Roots of  $\theta(B), \phi(B)$  outside unit circle

$$\alpha(B) = (1 - B)^d$$

$$Z_t \sim_{iid} \text{Normal}(0, \sigma^2)$$

The Observed series is represented as:

$$X_t^* = X_t + \text{outlier effect} \quad (0.2)$$

$$\text{AO:} \quad X_t^* = X_t + \omega I_t(t_1) \quad (0.3)$$

$$\text{LS:} \quad X_t^* = X_t + \frac{1}{1-B} \omega I_t(t_1) \quad (0.4)$$

$$\text{TC:} \quad X_t^* = X_t + \frac{1}{(1-\delta B)} \omega I_t(t_1) \quad (0.5)$$

$$\text{IO:} \quad X_t^* = X_t + \frac{\theta(B)}{\alpha(B)\phi(B)} \omega I_t(t_1) = \frac{\theta(B)}{\alpha(B)\phi(B)} [Z_t + \omega I_t(t_1)] \quad (0.6)$$

### 3.1.4 Outlier Estimation

The following Iterative procedure for detecting outliers, adjusting series, and fitting (seasonal) ARIMA model from [14] describes how the outliers are detected by estimation.

Obtain residuals  $\hat{e}_t$  from the observed series  $X_t^*$  by applying

$$\pi(B) = \frac{\alpha(B)\phi(B)}{\theta(B)} = 1 - \pi_1 B - \pi_2 B^2 - \pi_3 B^3 - \dots$$

(Remember  $X_t = \frac{\theta(B)}{\alpha(B)\phi(B)} Z_t$ )

If there were no outliers, result is white noise:  $\pi(B)X_t = Z_t$

When outlier present at  $t = t_1$ , residuals  $\hat{e}_t = \pi(B)X_t^*$  for  $t = t_1, \dots, n$  reveal outlier effect.

Least-squares estimate:

$$\hat{\omega} = \frac{\sum_{t=t_1}^n \hat{e}_t x_t}{\sum_{t=t_1}^n x_t^2}$$



Divide by standard error:

$$\hat{\tau} = \frac{\hat{\omega}}{\hat{\sigma} / \sqrt{\sum_{t=t_1}^n x_t^2}}$$

Approximately  $\sim \text{Normal}(0, 1)$

### 3.1.5 Outlier Detection

At each  $t = 1, \dots, n$ , for each outlier type (AO, LS, TC, IO), Estimate outlier effect  $\hat{\omega}$  and calculate  $\hat{\tau}$ . Large  $|\hat{\tau}|$  ( $> C$ ) indicates an outlier. Once outlier is detected, estimated effect can be subtracted to obtain adjusted series.[14]

### 3.1.6 Results

The time series shown in Figures 1 and 2 show the daily mean sentiment scores of the developers' messages.

Figures 3-5 show the decomposed time series into seasonal, trend and irregular components using moving averages, for all projects.

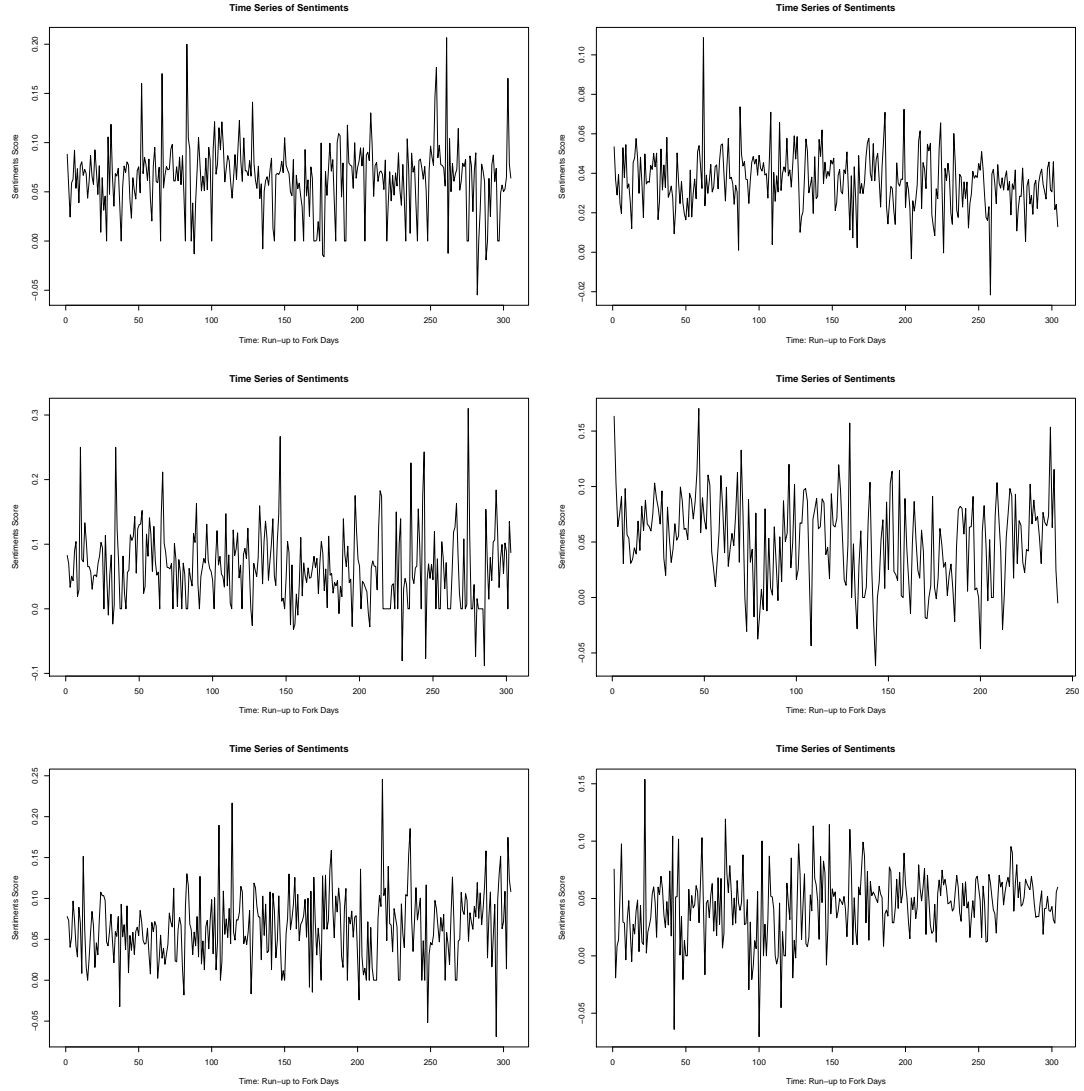


FIGURE 1: Time Series of sentiments in developers' messages for projects (From top to bottom, left to right) Kamailio (Conflict-Driven), ffmpeg (Conflict-Driven), Amarok (Non-Conflict-Driven), Apache CouchDB (Non-Conflict-Driven), Pidgin (Non-Conflict-Driven), MPlayer (Non-Conflict-Driven). It is not easy to visually analyse these raw time series.

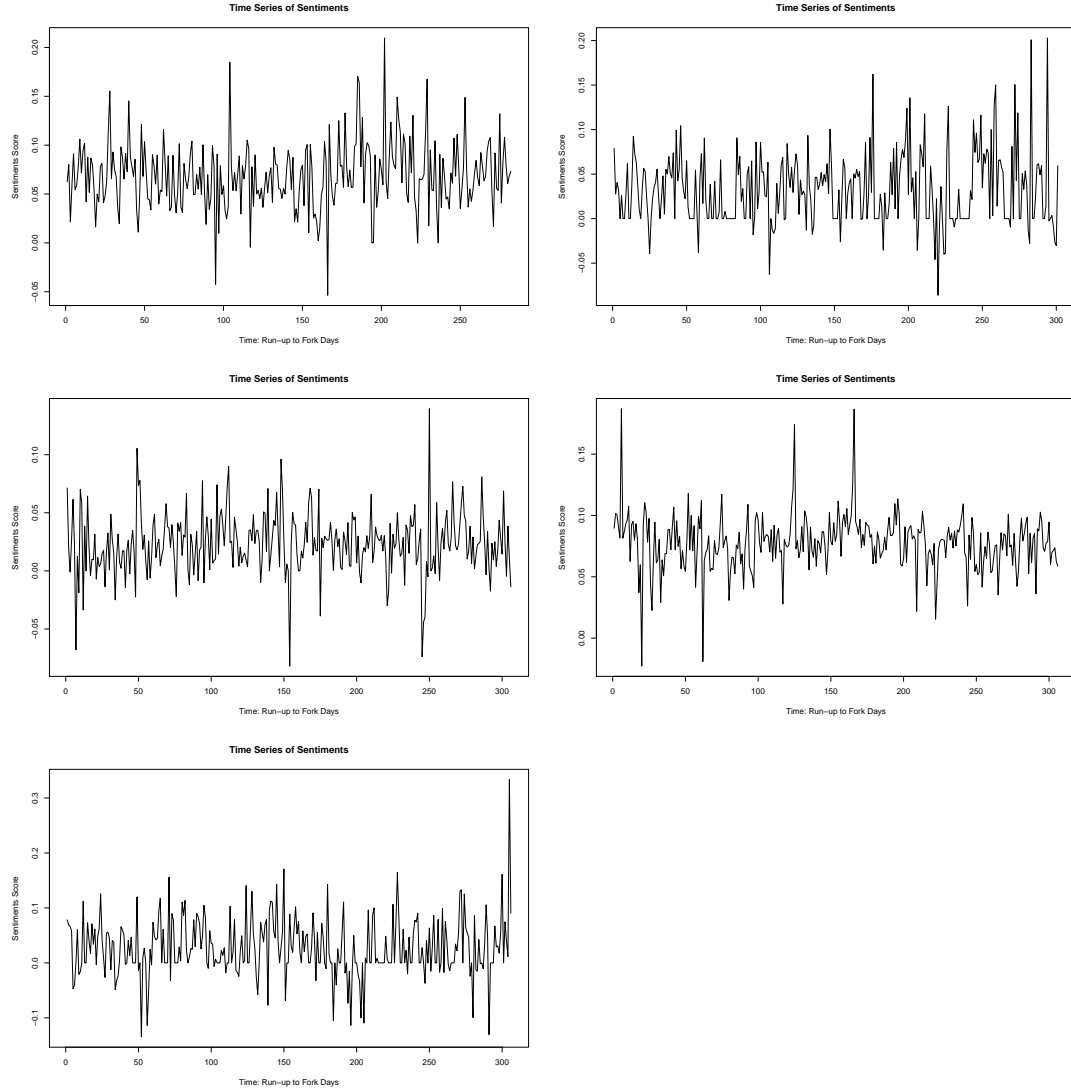


FIGURE 2: Time Series of sentiments in developers' messages for projects (From top to bottom, left to right) Asterisk (Non-Conflict-Driven), freeglut (Non-Conflict-Driven), Ceph (Not Forked), OpenStack Neutron (Not Forked), and GlusterFS (Not Forked). It is not easy to visually analyse these raw time series.

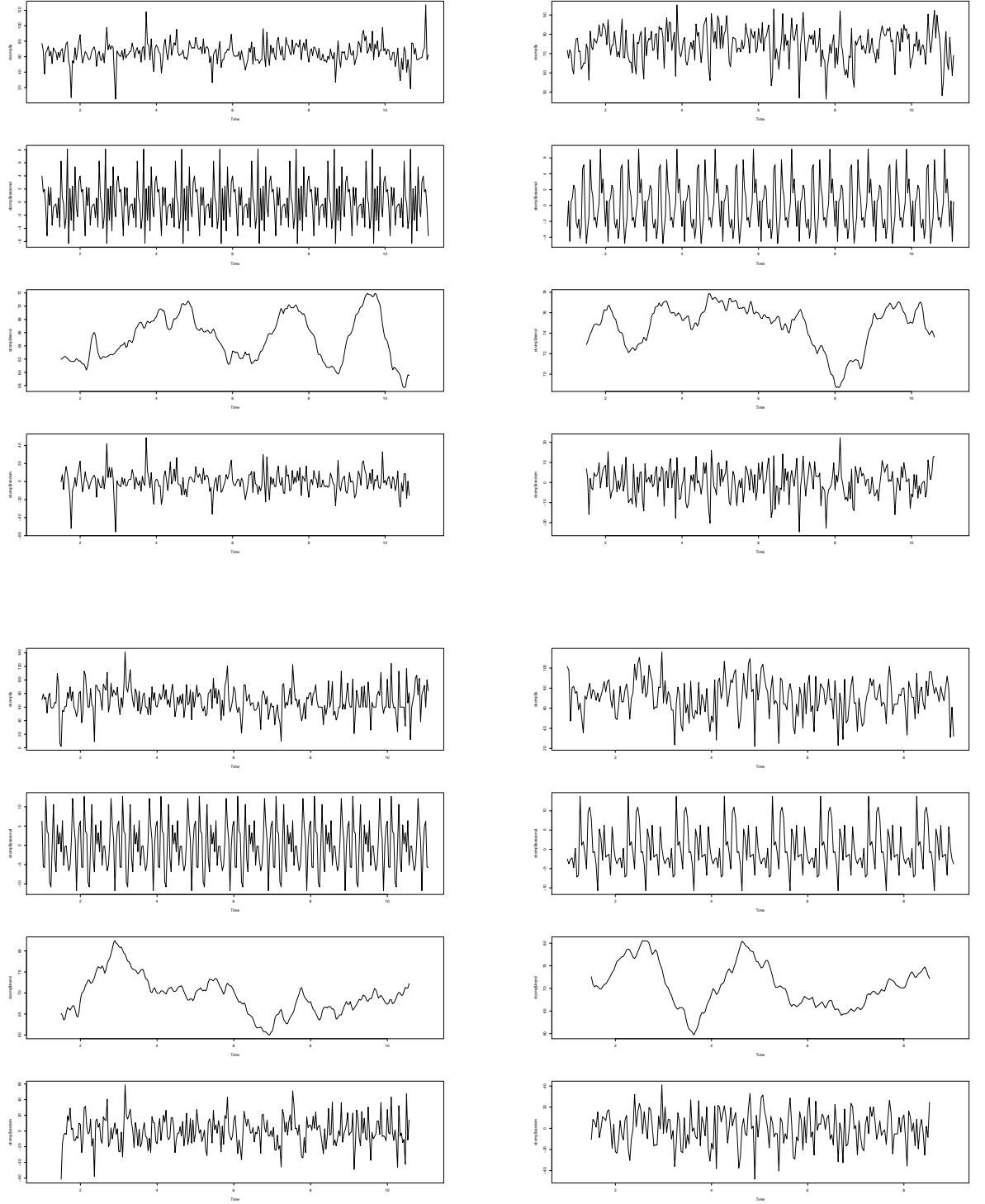


FIGURE 3: Time series decomposed into trend, seasonality and residual time Series for projects (From top to bottom, left to right) Kamailio (Conflict-Driven), ffmpeg (Conflict-Driven), Amarok (Non-Conflict-Driven), Apache CouchDB (Non-Conflict-Driven).

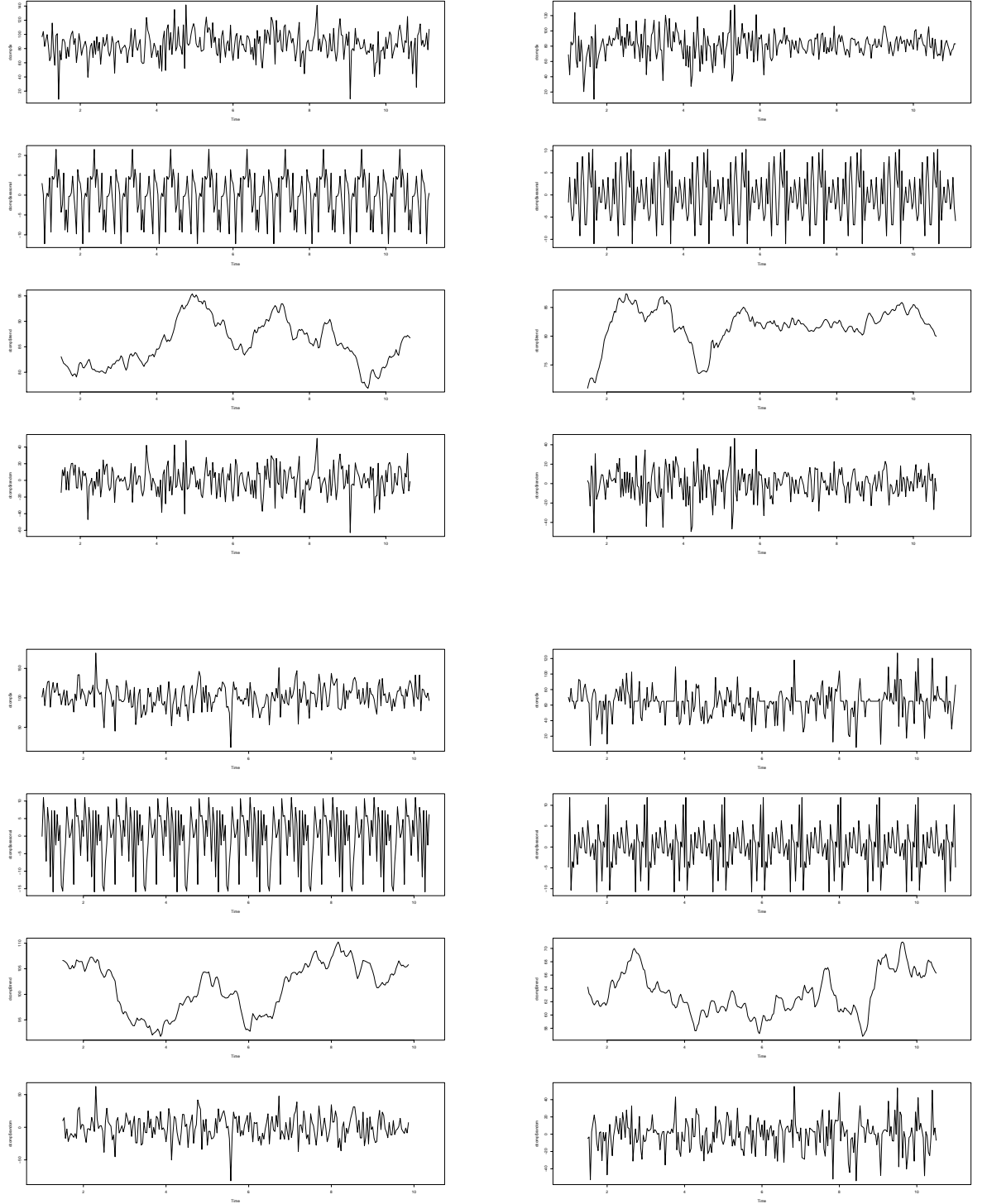


FIGURE 4: Time series decomposed into trend, seasonality and residual time Series for projects (From top to bottom, left to right) Pidgin (Non-Conflict-Driven), MPlayer (Non-Conflict-Driven), Asterisk (Non-Conflict-Driven), and freeglut (Non-Conflict-Driven).

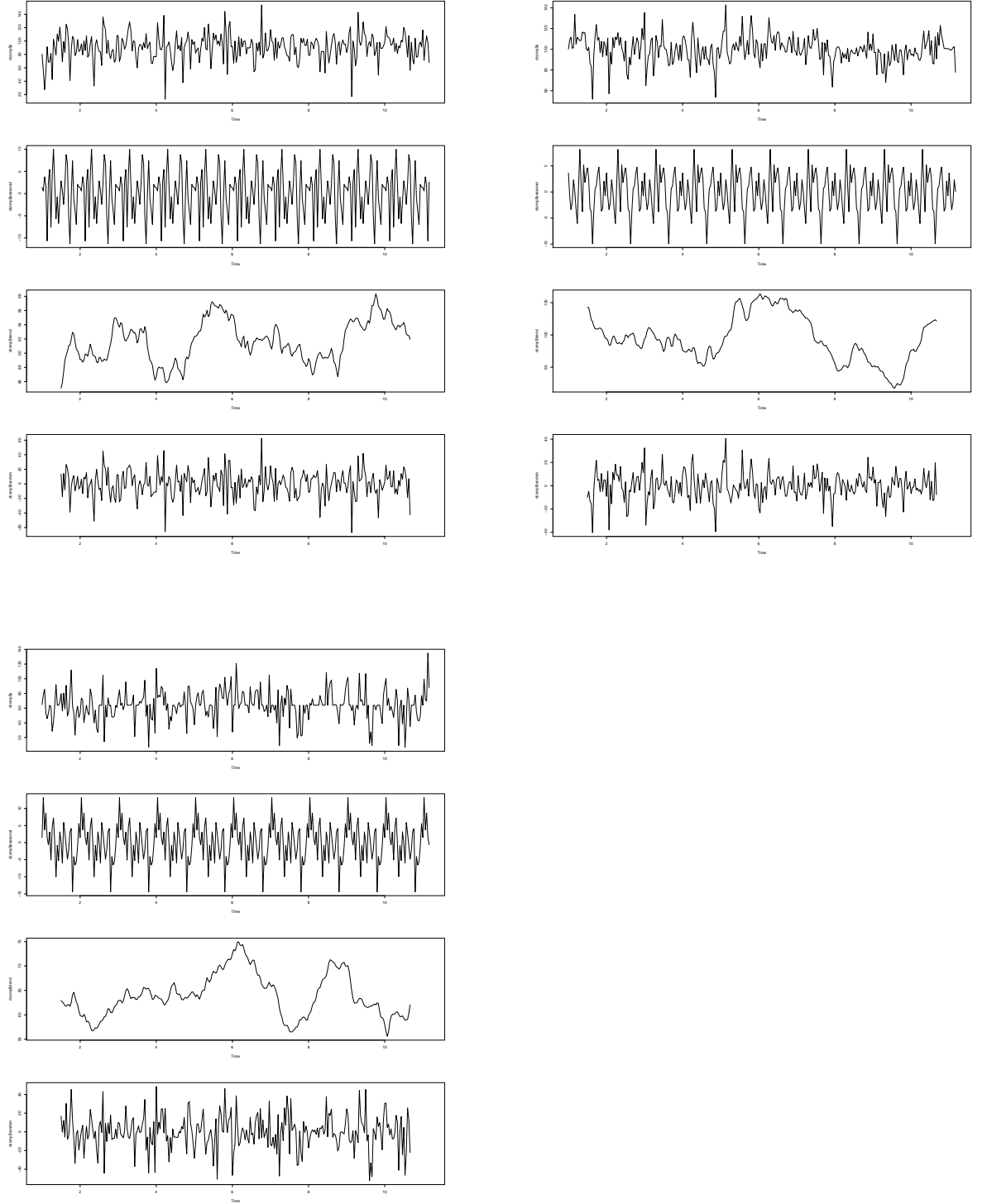


FIGURE 5: Time series decomposed into trend, seasonality and residual time Series for projects (From top to bottom, left to right) Ceph (Not Forked), Open-Stack Neutron (Not Forked), and GlusterFS (Not Forked).

TABLE 6: Anomalies detected in the time series of sentiments for project Kmailio (Conflict-Driven). Type column refers to the different anomalies as described in section 3.1.3. Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 128.8: log likelihood = -1170.2, aic = 2356.41)

	type	ind	coefhat	tstat	statistical significance (*)
1	AO	83	0.14	4.40	*
2	TC	170	-0.09	-4.12	*
3	TC	253	0.10	4.28	*
4	AO	261	0.14	4.44	*
5	TC	282	-0.10	-4.38	*

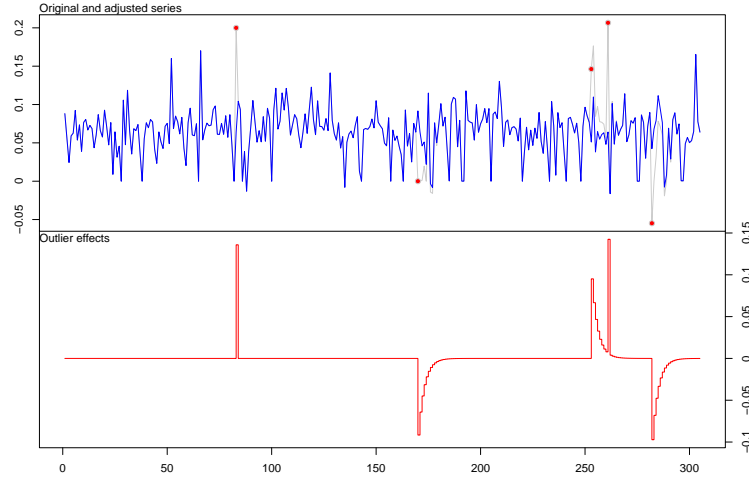


FIGURE 6: Time series with outliers detected for project Kmailio (conflict-driven).

TABLE 7: Anomalies detected in the time series of sentiments for project ffmpeg (Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 57.33: log likelihood = -1043.78, aic = 2101.57)

	type	ind	coefhat	tstat	statistical significance (*)
1	AO	62	0.07	5.80	*
2	AO	258	-0.05	-4.25	*

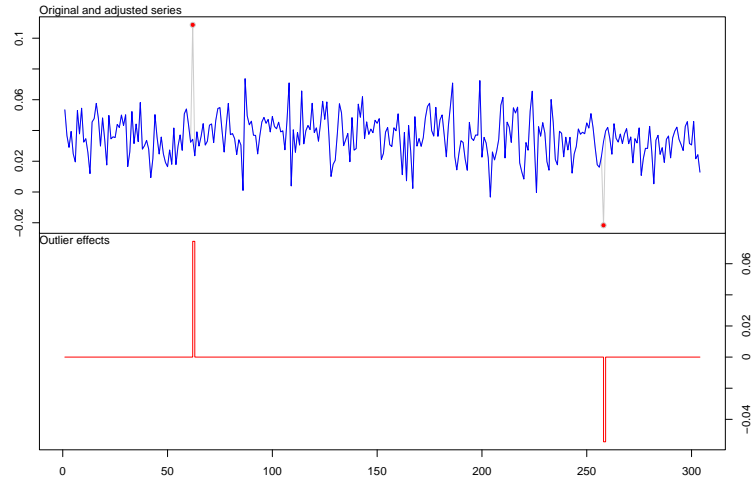


FIGURE 7: Time series with outliers detected for project ffmpeg (conflict-driven).



TABLE 8: Anomalies detected in the time series of sentiments for project Amarok (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 375.6: log likelihood = -1326.61, aic = 2661.22)

	type	ind	coefhat	tstat	statistical significance (*)
1	AO	10	0.19	3.87	*
2	AO	146	0.20	4.12	*
3	AO	235	0.18	3.79	*
4	AO	244	0.19	3.87	*
5	IO	274	0.20	4.62	*

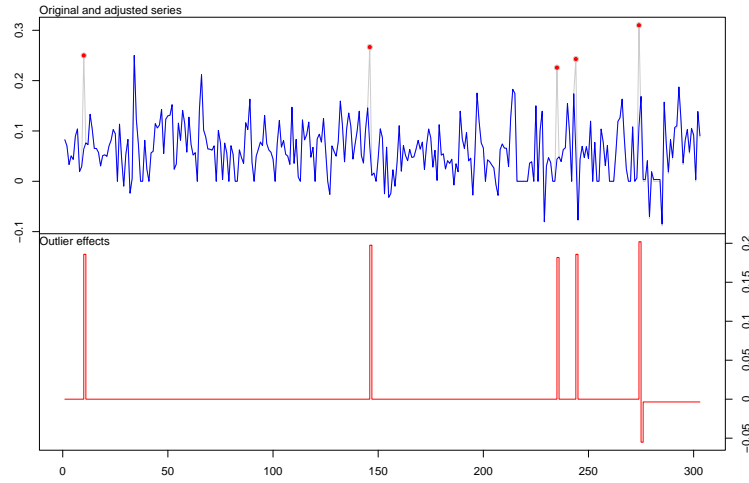


FIGURE 8: Time series with outliers detected for project Amarok (non-conflict-driven).

TABLE 9: Anomalies detected in the time series of sentiments for project Pidgin (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 312.3: log likelihood = -1307.25, aic = 2622.49)

	type	ind	coefhat	tstat	statistical significance (*)
1	IO	217	0.17	4.13	*

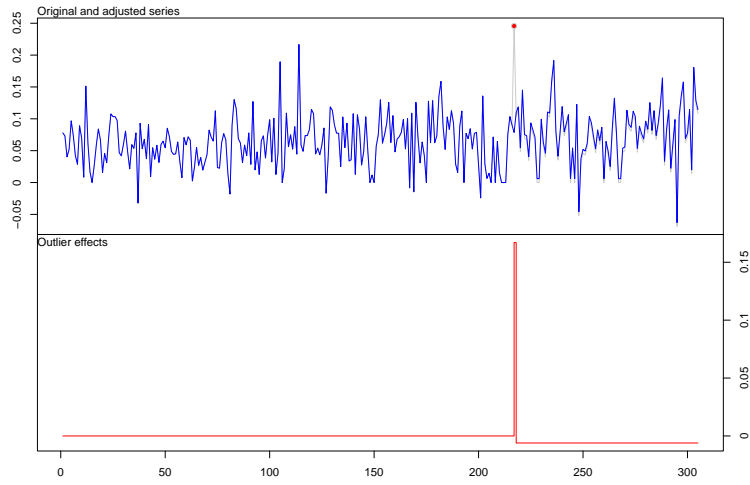


FIGURE 9: Time series with outliers detected for project Pidgin (non-conflict-driven).

TABLE 10: Anomalies detected in the time series of sentiments for project MPlayer (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 225.5: log likelihood = -1251.94, aic = 2517.88)

	type	ind	coefhat	tstat	statistical significance (*)
1	IO	22	0.11	4.75	*
2	AO	42	-0.11	-4.43	*
3	TC	93	-0.08	-3.91	*
4	AO	100	-0.11	-4.29	*
5	TC	115	-0.08	-3.81	*

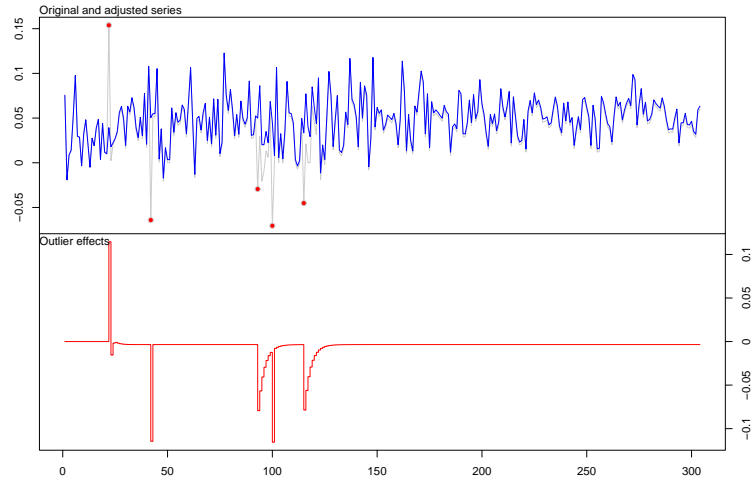


FIGURE 10: Time series with outliers detected for project MPlayer (non-conflict-driven).

TABLE 11: Anomalies detected in the time series of sentiments for project Asterisk (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 358.2: log likelihood = -1227.85, aic = 2463.7 )

	type	ind	coefhat	tstat	statistical significance (*)
1	AO	95	-0.11	-3.84	*
2	AO	104	0.12	4.20	*
3	AO	166	-0.12	-4.23	*
4	TC	185	0.10	5.06	*
5	AO	202	0.14	5.07	*
6	TC	210	0.08	4.07	*
7	AO	229	0.10	3.59	*

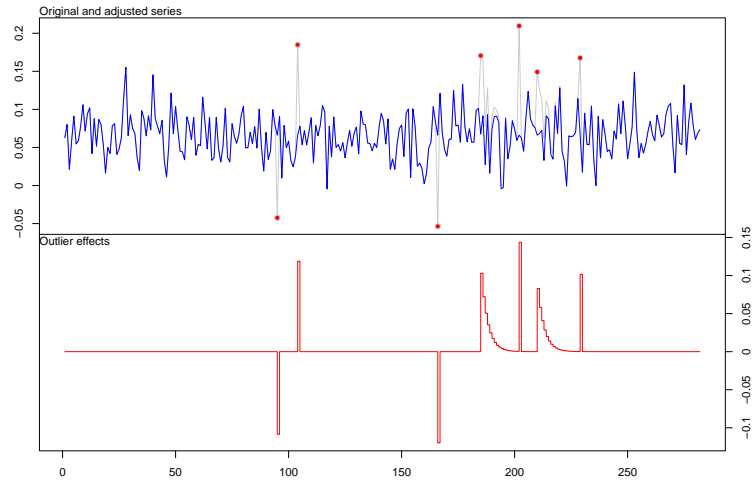


FIGURE 11: Time series with outliers detected for project Asterisk (non-conflict-driven).

TABLE 12: Anomalies detected in the time series of sentiments for project freeglut (Non-Conflict-Driven). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value. ( $\sigma^2$  estimated as 348.6: log likelihood = -1307.6, aic = 2619.2)

	type	ind	coefhat	tstat	statistical significance (*)
1	AO	176	0.13	3.80	*
2	AO	283	0.20	5.62	*
3	AO	294	0.19	5.46	*

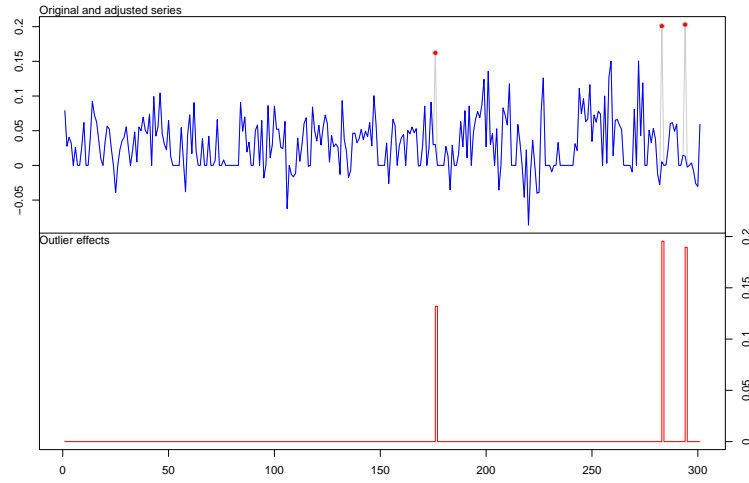


FIGURE 12: Time series with outliers detected for project freeglut (non-conflict-driven).

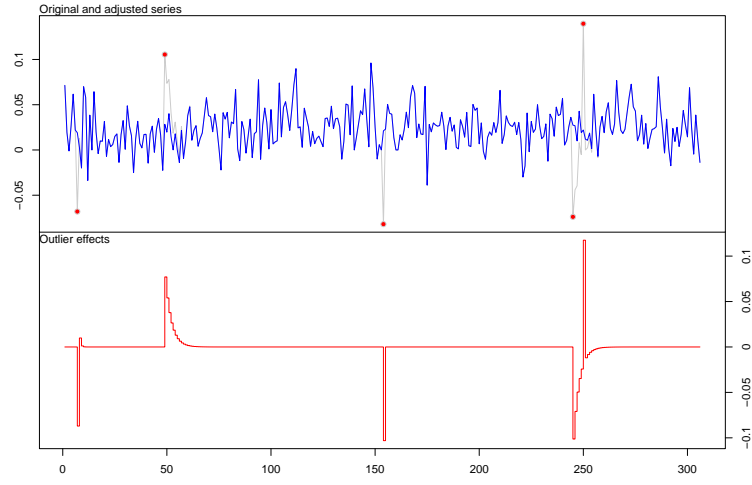


FIGURE 13: Time series with outliers detected for project Ceph (Not forked).

TABLE 13: Anomalies detected in the time series of sentiments for project Ceph (Not forked). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value.

	type	ind	coefhat	tstat	statistical significance (*)
1	IO	7	-0.09	-4.07	*
2	TC	49	0.08	4.47	*
3	AO	154	-0.10	-4.72	*
4	TC	245	-0.10	-5.91	*
5	AO	250	0.13	6.14	*

TABLE 14: Anomalies detected in the time series of sentiments for project OpenStack Neutron (Not forked). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value.

	type	ind	coefhat	tstat	statistical significance (*)
1	AO	6	0.11	6.34	*
2	AO	20	-0.10	-5.87	*
3	IO	62	-0.09	-5.65	*
4	AO	125	0.09	5.06	*
5	IO	166	0.08	5.01	*

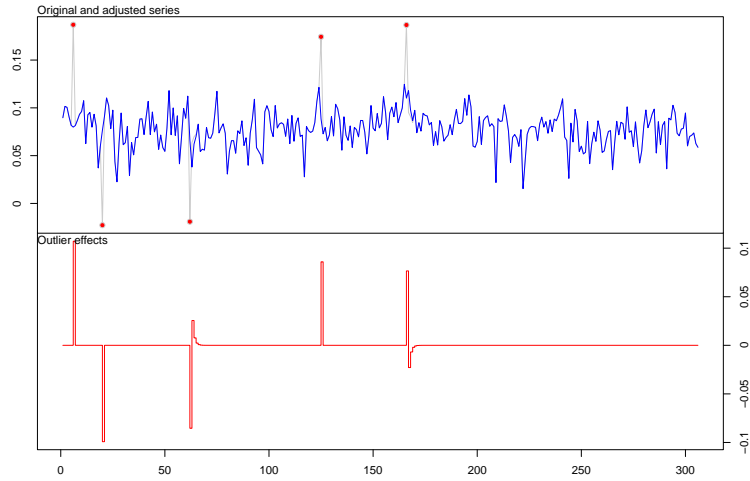


FIGURE 14: Time series with outliers detected for project OpenStack Neutron (Not forked).

TABLE 15: Anomalies detected in the time series of sentiments for project GlusterFS (Not forked). Note that all listed outliers are statistically significant with t-statistics greater than 1.96 in absolute value.

	type	ind	coefhat	tstat	statistical significance (*)
1	IO	305	0.28	5.77	*

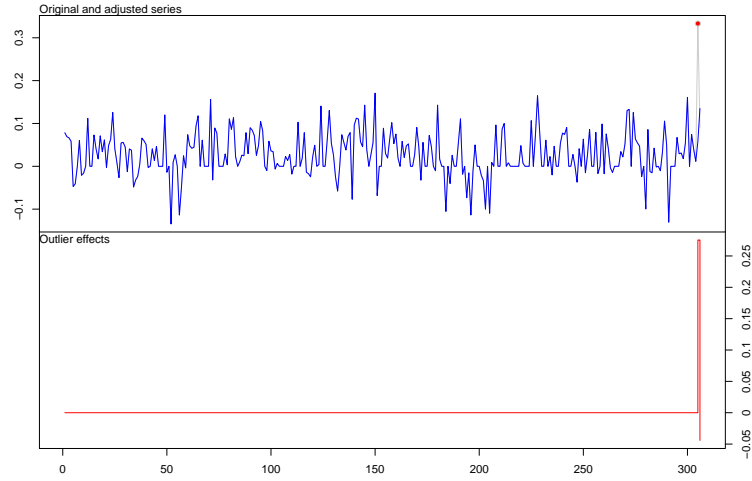


FIGURE 15: Time series with outliers detected for project GlusterFS (Not forked).

No anomalies were found in the time series of sentiments for project Apache CouchDB (Non-Conflict-Driven).



The description of these anomalies and how can they be used is discussed in section 4.

### **3.1.7 Summary**

Time series analysis of open source developers communication sentiments shows that in conflict-driven forks, we see negative dip anomalies pre-fork near the fork date, in contrast to the non-conflict-driven forks, where we see either no anomalies, or positive spikes anomalies pre-fork. We got close to determining when one vs. other has occurred, however, given the false positives and the scarcity of projects, we cannot be conclusive. So we turn to social network analysis 3.2, to figure out whether a detected anomaly is “a storm in a teacup”, or a situation needing intervention. What do these networks look like, who is involved, and what are the forces are shaping the community. This is discussed in the next section (3.2).

## **3.2 Social network analysis of open source developers interaction graphs**

We used a social network analysis developer-oriented approach to statistically model the changes a community goes through in the run-up to a fork [7]. The model represented tie formation, tie breakage, and tie maintenance between developers. It uses stochastic estimation methods to estimate several model parameters that capture the variance in the changes the community goes through. The model leads us to estimates of the significance of several parameters on this longitudinal change.

The process is described in detail in the following. Figure 16 show the overview of the methodology.

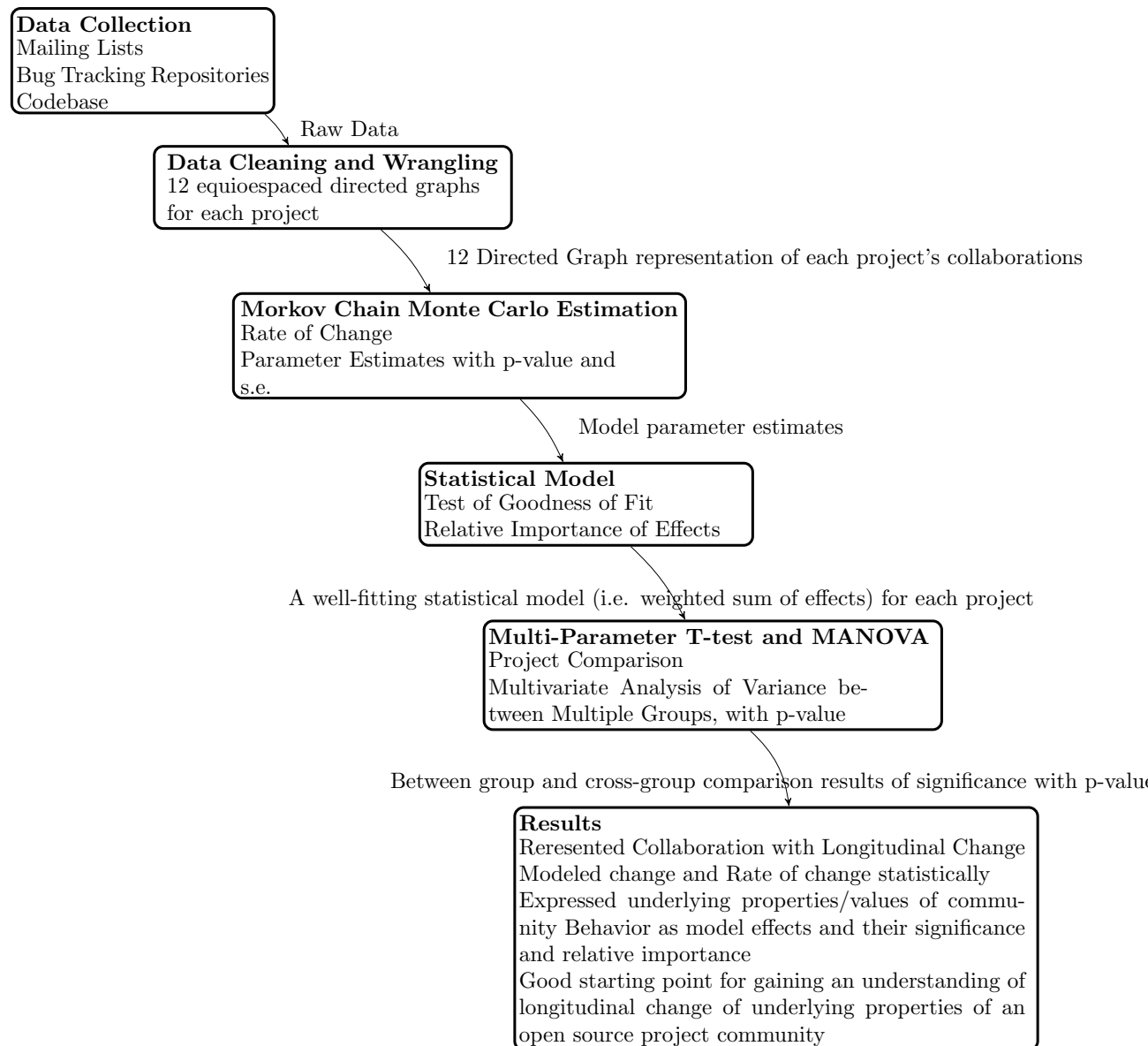


FIGURE 16: The methodology in a glance

### 3.2.1 Data

To find patterns uniquely associated with conflict-driven forks, we gathered data on projects from both conflict-driven and non-conflict-driven forking categories. The data consisted of developer interaction logs on the project’s developers mailing lists, where developers’ interact by sending and receiving emails and source-code repository contribution logs, where developers interact by modifying the code, and/or working on the same source files. The sociograms was formed based on interactions among developers and the contribution history was used as a covariate in the model. (Note that the content of the messages send and received by the top contributors of the project in the months leading to the forking events was sentiment-analyzed.) The time period for which data was collected is 10-months leading to when the fork happened. This should supposedly capture the social context prior, and at the time of the fork.

To find projects in conflict-driven and non-conflict-driven categories, we looked at the list of all significant open source software forks in the past three decades as compiled by Robles and Gonzalez-Barahona [47]. Their study found the reasons behind each fork, listed in Table 1. We applied three selection criteria to the 220 forked projects on that list to find projects in conflict-driven and non-conflict-driven categories. A project was short-listed as either a conflict-driven or non-conflict-driven if **a)** the forking was relatively recent, i.e., happened after the year 2000, **b)** its data was existent and available to access and download online, or was made accessible to us after our requests; and **c)** the project had a sizable developer community, i.e., more than a dozen developers, which means it would be large enough to make a sociogram for a meaningful statistical analysis. For the Not Forked category, we chose well-known, stable projects that had been around for a while (two years), and had large communities; and had not forked;

and were similar in size of the development team. Similarity in size is a constrain that our statistical method imposes for the results to be meaningful, as described in detail in section 3.2.2. The preceding criteria resulted in the projects listed in Table 4. The rest of the 220 forked projects were discarded, because they did not meet the described filtering criteria.

### 3.2.2 Analysis

#### Sociogram Formation and Statistical Study

Social connections and non-connections can be represented as graphs, in which the nodes represent actors (developers) and the edges represent the interaction(s) between actors or lack thereof. Such graphs can be a snapshot of a network – a static sociogram – or a changing network, also called a dynamic sociogram. In this phase, we process interactions data to form a communication sociogram of the community.

Two types of analysis can be done on sociograms: Either a *cross-sectional* study, in which only one snapshot of the network is looked at and analyzed; or a *longitudinal* study, in which several consecutive snapshots of the network are looked at and studied. We are interested in patterns in the run-up to forks, therefore, unlike most existing research on forking, we do a longitudinal study.

A longitudinal study can look at the sociograms in the following two distinct approaches:

1. A *network-specific* approach, which can be called a skin-deep *measurement-specific* look, in which, we focus on the observed networks, measure their properties, to describe the structure of the *observed networks* with numeric summaries/descriptors, e.g. see Appendix B Figures 18 and 19. Many studies measuring centralities as their only metrics, fall into this category.

2. A *population-processes* approach, in which we treat an observed network as one instance from a set of all possible networks with the same number of nodes, and with similar characteristics. In a population-processes approach, the observed network is only good to help us understand the social forces/processes that generated it, because we are interested in the forces/processes that underlie the structure and changes of the network; these reflect the values shared or not shared by the community members, and represent the group behavior, and behavior change trends that lead to a forking event.

Appendix B B Figures 18 and 19 are good examples of the expressiveness and limitations of a network-specific approach. Appendix B B Figure 18 shows the normalized change in the number of active nodes (developers) and edges (interaction between developers) for eight FOSS projects from conflict-driven and non-conflict-driven categories.

The normalized number of nodes and edges for the conflict-driven forks (i.e., second and third row) show a sudden sharp decrease around the month of fork. This can be because of dissolution of the development team, who either left the project, or stopped being active contributors as much as before. Appendix B B Figure 19 shows normalized change of diameter over time for the same projects. It is hard to make sense of normalized diameter changes across projects, and this demonstrates the limitations of network-specific approaches, as they may, or may not (often the case), help us understand the network dynamics.

As an example of expressiveness, Appendix B B Figure 18 shows that a measurements-only approach shows us the trends of change or no-change in the measured metrics. It also shows one limitation of such an approach: even though we can see trend changes, it is hard to make sense of such trends given the short

scope of the measurements. Additionally, it is hard to constitute a reference point and back it up mathematically or statistically.

The following subsection lists the reasons why we need to use a statistical model, used in the *population-processes* approach.

### **Why a statistical model is needed?**

One may wonder why we should look beyond the observed network. We can do measurements on the observed networks data, and get some descriptive statistics. This would be a superficial look at the data, and even though necessary, is not good enough and has several shortcomings. There is another way of looking at these networks, in a less-superficial way; namely, finding a model that fits the data and its longitudinal change. So, instead of a superficial look at the data, we can find a well-fitting statistical model of our observed interactions network. The following lists the reasons why we need to go beyond a superficial measurement-only approach.

1. For the *observed network*, a small observation error or sampling error, and the uncertainty involved with real-world communication, can result in large perturbation in the numeric descriptors used to describe static graphs.
2. The traditional *network-specific* approach assumes edges in sociograms are statistically independent, (and/or identically distributed). This can be misleading, as, social network data are relational. For example, in real-life human communication, the likelihood of forming ties with friends of a friend is higher than a stranger, as Balance Theory suggests [25].
3. In the population-specific approach, we try to identify the social forces that have formed the observed network, by simulation a population of similar

networks of the same characteristics. After finding the statistical distribution of network population, then we can compare the observed network to the population distribution of possible networks of that size, and see how significant and likely it is to observe such a graph, as compared to observing a randomly-generated graph. This is useful, because it gives us a reference point to compare our observed graph with, and to weed out the properties generated by random processes, and to find the statistically significant network statistics. In short, with a statistical model, we can draw inferences about whether certain network structures and substructures are more commonly observed in the observed network than might be expected by chance [46].

4. Stochastic models capture the regularities in the processes that caused the network ties form, as well as variability that are hard to model otherwise. A model that considers stochasticity allows us to understand the uncertainty associated with an observed network. It makes it possible to learn about the distribution of possible networks for a given specification of a model [46].
5. Different social processes may manifest similar network structures. For example, clustering in a network might be because of structural effects, e.g., structural balance, or through node-level effects, e.g., homophily. To determine which one is the case in our observed network, a statistical model that incorporates both covariates can help. We then can assess the contribution of each covariate, and infer which social process underlies the observed network [46].
6. Localized processes might not scale to the entire network well. The com-

bination of the overall structure and the localized processes is hard to investigate without a model. (This micro-macro difference may be investigated through model simulation.)[46]

In summary, a measurement-only approach is not capable of explaining the longitudinal changes in an open source community’s network properly. It can confuse us, and it can mislead us. Our initial study, described in appendix section B.2, is an example of a measurement-only approach, which shows such limitations. This initial study guided in the proper direction; namely, trying to find a statistical model that can explain the networks’ longitudinal changes. In the following paragraphs (3.2.2), we described such a model.

### **The Statistical Model**

Longitudinal evolution of a network data is the result of many small atomic changes occurring between the consecutively observed networks. In our case, software developers are the actors in the networks, and they can form a connection with another developer, break off an existing connection, or maintain their status quo. These are the four possibilities of atomic change within our evolving networks: (1) forming a new tie; (2) breaking off an existing tie; (3) maintaining a non-connection; and (4) maintaining a connection. We assume a continuous-time network evolution, even though our observations are made at two or more discrete time points.

The state-of-the-art in studying longitudinal social networks, is the idea of *actor-oriented models* [52], based on a model of developers changing their outgoing ties as a consequence of a stochastic optimization of an *objective function*. This framework assumes that the observed networks at discrete times, are outcomes of a continuous-time Markov process. In the case of open source de-



velopers, the actor-oriented model, can be informally described as open-source-developer-oriented model, in which, it is assumed that developers are in charge of their communication and collaboration choices. They choose to have interactions with certain other developers and/or they choose to stop having interactions with another developer. In short, they have autonomy in choosing their connections.

Let the data for our statistical developer-oriented model be  $M$  repeated observations on a network with  $g$  developers. The  $M$  observed networks (at least two) are represented as directed graphs with adjacency matrices  $X(t_m) = (X_{ij}(t_m))$  for  $m = 1, \dots, M$ , where  $i$  and  $j$  range from 1 to  $g$ . The variable  $X_{ij}$  shows whether at time  $t$  there exists a tie from  $i$  to  $j$  (value 1) or not (value 0). By definition,  $\forall i, X_{ii} = 0$  (i.e. the diagonal of the adjacency matrices).

In order to model the network evolution from  $X(t_1)$  to  $X(t_2)$ , and so on, it is natural to treat the network dynamics as the result of a series of small atomic changes, and not bound to the observation moment, but rather as a more or less continuous process. In this way, the current network structure is a determinant of the likelihood of the changes that might happen next [15].

For each change, the model focuses on the developer whose tie is changing. We assume that developer  $i$  has control over the set of outgoing tie variables  $(X_{i1}, \dots, X_{ig})$  (i.e. the  $i^{th}$  row of the adjacency matrix). The network changes one tie at a time. We call such an atomic change a *ministep*. The moment at which developer  $i$  changes one of his ties, and the kind of change that he makes, can depend on attributes represented by observed covariates, and the network structure. The moment is stochastically determined by the *rate function*, and the particular change to make, is determined by the *objective function* and the *gratification function*. We cannot calculate this complex model exactly. Rather

than calculating exactly, we estimate it using a Monte Carlo Markov Chain method. The estimated model is used to test hypotheses about the forked FOSS communities.

These above three functions and their definitions taken from [51] are explained in detail the following subsections.

### Rate Function

The *rate function*  $\lambda_i(x)$  for developer  $i$  is the rate at which developer  $i$ 's outgoing connections changes occur. It models how frequently the developers make ministeps.

The rate function is formally defined [51] by

$$\lambda_i(x) = \lim_{dt \rightarrow 0} \frac{1}{dt} P(X_{ij}(t+dt) \neq X_{ij}(t) \text{ for some } j \in \{i, \dots, g\} | X(t) = x). \quad (0.7)$$

The simplest specification of the rate of change is that all developers have the same rate of change of their ties.

### Objective Function

The *objective function*  $f_i(s)$  for developer  $i$  is the value attached to the network configuration  $x$ .

The idea is that, given the opportunity to make a change in his outgoing tie variables  $(X_{i1}, \dots, X_{ig})$ , developer  $i$  selects the change that gives the greatest increase in the objective function. We assume that if there is difference between developers in their objective functions, these differences can be represented based on the model covariates [51].

Let's denote the present network by  $x = X(t)$ . The new network that would be the result of a ministep (i.e. changing a single tie variable  $x_{ij}$  into its

opposite  $1 - x_{ij}$  ) by developer  $i$  is denoted by  $x(i \mapsto j)$ . This choice is modeled in the following way:

Let  $U(j)$  denote a random variable which represents the attraction for  $i$  to  $j$ . We assume that these  $U_j$  are random variables distributed symmetrically about 0, and independently generated for each new minstep. In this way, the developer  $i$  chooses to change his tie variable with the other developer  $j$  for whom the following value [51] is the highest:

$$f_i(x(i \mapsto j)) + U(j) \quad (0.8)$$

This is a short-sighted stochastic optimization rule: short-sighted because only the situation one step after the current step is considered; and stochastic because the unexplained part is modeled by a random variable  $U(j)$ .

$U(j)$  is chosen to have Gumbel distribution with mean 0 and scale parameter 1. In this way, the probability that developer  $i$  chooses to change  $x_{ij}$  for any particular  $j$ , given that such a change happens, is [51],

$$p_{ij}(x) = \frac{\exp(f_i(i \mapsto j) - f_i(x))}{\sum_{h=1, h \neq i}^g \exp(f_i(i \mapsto h) - f_i(x))} \quad (0.9)$$

## Gratification Function

Sometimes the order in which changes occur makes a difference for the desirability of the network. Such differences cannot be represented by the objective function. So, we need another function, called gratification function.

The *gratification function*  $g_i(x, j)$  for developer  $i$  is the value attached to this developer (in addition to what follows from the objective function) to the act of changing the tie variable  $x_{ij}$  from  $i$  to  $j$ , given the current network

configuration  $x$ .

In the case that a gratification function is included, the developer  $i$  chooses to change  $x_{ij}$  for that other developer  $j$  for whom the following value [51] is the highest:

$$f_i(x(i \mapsto j)) + g_i(x, j) + U(j) \quad (0.10)$$

And we will have

$$p_{ij}(x) = \frac{\exp(f_i(i \mapsto j) + g_i(x, j) - f_i(x))}{\sum_{h=1, h \neq i}^g \exp(f_i(i \mapsto h) + g_i(x, h) - f_i(x))} \quad (0.11)$$

### Markov Chain Transition Rate Matrix

The components of the developers-oriented model, described above, define a continuous-time Markov chain on the space  $\chi$  of all directed graphs on this set of  $g$  developers. This Markov chain is used to estimate the model parameters stochastically, instead of calculating them exactly, which is not possible for us.

This Markov chain has a transition rate matrix. The transition rate matrix (also called intensity matrix), for this model is given by

$$\begin{aligned} q_{ij}(x) &= \lim_{dt \rightarrow 0} \frac{1}{dt} \mathbb{P}(X(t + dt) = X(i \mapsto j) | X(t) = x) \\ &= \lambda_i(x) p_{ij}(x) \end{aligned} \quad (0.12)$$

Expression (0.12) shows the rate at which developer  $i$  makes ministeps, multiplied by the probability that he changes the arc variable  $X_{ij}$ , if he makes a ministep.

## Markov Chain Simulation

Our Markov chain can be simulated by repeating the following steps [51]:

Start at time  $t$  with directed graph  $x$ .

1. Define  $\lambda_+ = \sum_{i=1}^g \lambda_i(x)$  and let  $\Delta t$  be a random variable with the exponential distribution with parameter  $\lambda_+(x)$ .
2. The developer  $i$  is chosen randomly with probabilities  $\frac{\lambda_i(x)}{\lambda_+(x)}$
3. Given this  $i$ , choose developer  $j$  randomly with probabilities

$$p_{ij}(x) = \frac{\exp(f_i(i \mapsto j) + g_i(x, j) - f_i(x))}{\sum_{h=1, h \neq i}^g \exp(f_i(i \mapsto h) + g_i(x, h) - f_i(x))}$$

4. Now change  $t$  to  $t + \Delta t$  and change  $x_{ij}$  to  $(1 - x_{ij})$

## Model Specification

In the previous sections, we described why we need a statistical model to describe the longitudinal trends, rather than numerical measurements. We then described a framework for such a model; namely the developer-oriented model. We described the components of this statistical model; the rate function; the objective function, and the gratification function.

In this section, we describe what network effects can be measured and tested in these functions. Specifically, we need to supply these functions with a specific model for the rate, objective and gratification functions. These functions depend on unknown parameters that need to be estimated based on the data. The estimation is explained in section 3.2.4.

In the following, we explain what model effects can be included in the objective and gratification functions, to be estimated later.

## Objective Function

The following weighted sum represents the objective function:

$$f_i(\beta, x) = \sum_{k=1}^L \beta_k s_{ik}(x) \quad (0.13)$$

Parameters  $\beta = (\beta_1, \dots, \beta_L)$  is to be estimated. Functions  $s_{ik}(x)$  can be the following [51]:

### 3.2.3 Structural Effects

For the structural effects, some of the following effects were used. The mathematical formulas for the following effects are included in the appendix section C.

1. The reciprocity effect, which reflects the tendency toward reciprocation of connections. A high value for its model parameter will indicate a high tendency of developers for reciprocated interactions.
2. The closure effects (e.g. in friendship networks, it means, friends of friends tend to become friends) For example,
  - (a) Transitive triplets effect, which models the tendency toward network closure. It reflects the preference of developers to be connected to developers with similar outgoing ties.
  - (b) Transitive ties, which reflects network embeddedness. It is similar to transitive triplets effect, however, it only count for the existence of at least one two-path  $i \rightarrow h \rightarrow j$ , rather than counting how many such two-paths exists.

- (c) Balance effect, which reflects how similar the outgoing ties of developer  $i$  are to the outgoing ties of the other developers to whom  $i$  is connected.
  - (d) Number of developers at distance two, which reflects and is an inverse measure of, network closure.
3. Three-cycles, may be interpreted as the tendency toward local hierarchy. It is similar to reciprocity defined for three developers, and is the opposite of hierarchy.
  4. Density effect which reflects tendency to have connections at all (i.e. the out-degree of a developer  $i$ ).
  5. Activity, which reflects the tendency of developers with high in-degree/out-degrees to send out more outgoing connections because of their current high in-degree/out-degree.
  6. Covariate effects: Developers' covariates may influence the formation or termination of ties. For example:
    - (a) Covariate V-related popularity, which reflects the sum of covariate V for all developers to whom developer  $i$  is connected
    - (b) Covariate V-related activity, which reflects the developer  $i$ 's out-degree multiplied by his covariate V value.
    - (c) Covariate V-related dissimilarity, which reflects the sum of differences in covariate V values' between developer  $i$  and all developers to whom developer  $i$  is connected.

We use the following developer attributes as covariates:

- Developer’s level of contribution/activity (e.g. code commits per month, or mailing list posts per month)
  - Developer’s seniority as a development community member (i.e. how long they have been in the community)
7. in-in degree assortativity, which reflects the tendency of developers with high in-degree to be connected to other developers with high in-degrees
  8. in-out degree assortativity, which reflects which reflects the tendency of developers with high in-degree to be connected to other developers with high out-degrees
  9. out-in degree assortativity, which reflects which reflects the tendency of developers with high out-degree to be connected to other developers with high in-degrees
  10. out-out degree assortativity, which reflects which reflects the tendency of developers with high out-degree to be connected to other developers with high out-degrees

### 3.2.4 Behavior-related Effects

Properties of developers can be called their behavior. For a behavior-related variable  $s_{ik}^z(x, z)$ , we can estimate the following effects.

Let’s denote a dependent behavior variable as  $Z$ .

1. Shape, which reflects the drive toward high values on the variable  $Z$ .
2. Similarity:



- (a) Average similarity, which reflects the preference of the developer  $i$  to be similar to his alters, with regard to the behavior variable  $Z$
  - (b) Total similarity, which reflects the preference of the developer  $i$  to be similar to his alters. The total influence of the alters depends on the number of the alters
3. Average alter, which reflects that the developers who have alters with higher average values of  $Z$ , have a stronger tendency to having high values of  $Z$ .

### Gratification Function

The following weighted sum represents the gratification function:

$$g_i(\gamma, x, j) = \sum_{h=1}^H \gamma_h r_{ijh}(x) \quad (0.14)$$

Some possible functions  $r_{ijh}(x)$  are the following [51]:

1. Breaking off a reciprocated tie:

$$r_{ij1}(x) = x_{ij}x_{ji}$$

2. Number of indirected links for creating a new tie:

$$r_{ij2}(x) = (1 - x_{ij}) \sum_h x_{ih}x_{hj}$$

3. Effect of dyadic covariate  $W$  on breaking off a tie:

$$r_{ij3}(x) = x_{ij}w_{ij}$$

### Markov Chain Monte Carlo (MCMC) Estimation

The described statistical model for longitudinal analysis of open source software development communities is a complex model and cannot be exactly

calculated, but it can be stochastically estimated. We can simulate the longitudinal evolution, and estimate the model based on the simulations. Then we can choose an estimated model that has a good fit to the network data. The following section described the simulation and estimation procedures.

Network evolution can be simulated using a MCMC estimation method. The method of moments (MoM) can be used in the following way [51]:

Let  $x^{obs}(t_m), m = 1, \dots, M$  be the observed networks, and the objective function be  $f_i(\beta, x) = \sum_{k=1}^L \beta_k s_{ik}(x)$ . In MoM, the goal is to determine the parameters  $\beta_k$  such that, summed over  $i$  and  $m$ , the expected values of the  $s_{ik}(X(t_{m+1}))$  are equal to the observed values. In other words, MoM fits the *observed* to the *expected*. The observed target values are:

$$s_k^{obs} = \sum_{m=1}^{M-1} \sum_{i=1}^g s_{ik}(x^{obs}(t_{m+1})) \quad (k = 1, \dots, L) \quad (0.15)$$

The simulations run as follows [51]:

1. The distance between two directed graphs  $x$  and  $y$  is defined as

$$\|x - y\| = \sum_{i,j} |x_{ij} - y_{ij}| \quad (0.16)$$

and let the observed distances for  $m = 1, \dots, M-1$  be

$$c_m = \|x^{obs}(t_{m+1}) - x^{obs}(t_m)\| \quad (0.17)$$

2. Use the  $\beta$  parameter vector and the rate of change  $\lambda_i(x) = 1$
3. Do the following for  $m = 1, \dots, M-1$

- (a) Define the time as 0, and start with  $X_m(0) = x^{obs}(t_m)$
- (b) Simulate the model, as described in 3.2.2 until the first time point, where

$$\|X_m(R_m) - x^{obs}(t_m)\| = c_m \quad (0.18)$$

4. For  $k = 1, \dots, L$ , calculate the following statistics:

$$S_k = \sum_{m=1}^{M-1} \sum_{i=1}^g s_{ik}(X_m(R_m)) \quad (0.19)$$

The simulation output will be the random variables  $(S, R) = (S_1, \dots, S_L, R_1, \dots, R_{M-1})$ . The desirable outcome for the estimation is the vector parameter  $\hat{\beta}$  for which the expected and the observed vectors are the same.

So far, we discussed a number of metrics for social networks. We will investigate their statistical significance and effects to describe the different social networks that we encounter, to see if there is any systematic differences between conflict-driven and non-conflict-driven forks.

## Hypothesis Testing

So far, we have described what we can model, how to fit a model that explains the longitudinal changes in an open source community (or any similar network). Now, we test several hypothesis about these models. For example, we test whether a particular effect's (for example, reciprocity's) statistical significance on the model, and make meaningful statements about whether the network dynamics depends on this parameter with a p-value, indicating the significance level. Such test is explained in section 3.2.4.

Alternatively, we compare two models for two open source project com-

munities, and make conclusions about the difference in several model effects, for example, reciprocity and three-cycles, using the method described in section 3.2.4.

Lastly, we compare several open source communities (or such networks) and test the statistical significance of the differences between their models, using the methods described in section 3.2.4. This allows us to compare within-group and across-group comparisons of longitudinal models of open-source projects. In this study, we do all three discussed forms of tests of significance. In the following subsections, the details of these tests are described.

### Single Parameter Test

Using the actor-oriented model, t-type test of single parameters can be done using the parameter estimates and their standard errors (S.E.). For example, for testing the null hypothesis that a component  $k$  of the parameter vector  $\beta$  is 0, the test has approximately a standard normal distribution. For this null hypothesis:

$$H_0 : \beta_k = 0; \quad (0.20)$$

the t-test can be done using the following t-statistic:

$$\frac{\hat{\beta}_k}{S.E.(\hat{\beta}_k)} \quad (0.21)$$

Given a p-value  $\leq 0.05$  would indicate a strong evidence that the network dynamics depends on the corresponding parameter's effect.

### Multi-Parameter Differences Between Groups

Given the parameter estimates  $\hat{\beta}_a$  and  $\hat{\beta}_b$ , and their standard errors  $S.E._a$  and  $S.E._b$ , the differences between two groups can be tested using the following

t-statistics. This t-statistics, under the null hypothesis of equal parameters, has an approximately standard normal distribution:

$$\frac{\hat{\beta}_a - \hat{\beta}_b}{\sqrt{S.E._a^2 + S.E._b^2}} \quad (0.22)$$

### Multivariate Analysis of Variance Between Multiple Groups

To test the statistical significance of the mean differences between groups (in our case, the different categories of forking), we may use Multivariate Analysis of Variance (MANOVA) method. In contrast to ANOVA, MANOVA method tests for the difference in two or more vectors of means.

If the data does not exhibit statistically significant differences, then our study will not reject the null hypotheses, and our assumptions cannot be verified. We argue that based on the guiding theories and the collaboration data we gathered, if the forking reasons were socially-related, they should be reflected in the data, and so, should be reflected in the longitudinal model of their collaboration network evolution. This is a conjecture, and our study will try to show that this conjecture is true, using sound statistical modeling.

#### 3.2.5 Results

The tables 16-25 shows the results of the social network analysis longitudinal modeling of the projects. The rate parameter models how frequently the developers make mini-steps. The effects and their estimates that were found to be statistically significant are marked with asterisks, with their standard deviations and t statistics.

We describe what the findings mean and how to interpret the findings in section 4.

TABLE 16: Estimates parameters in the model for Project Kmailio (Conflict-Driven). The outdegree density, transitive triplets, out-out-degree assortativity, and developers' seniority effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	1.26	(0.30)	-
Rate 2	1.99	(0.40)	-
Rate 3	2.08	(0.42)	-
Rate 4	2.51	(0.48)	-
Rate 5	2.93	(0.56)	-
Rate 6	2.69	(0.52)	-
Rate 7	2.18	(0.46)	-
Rate 8	2.45	(0.51)	-
Rate 9	2.77	(0.50)	-
outdegree (density)	-2.11***	( 0.25)	-8.38
reciprocity	1.19 <sup>†</sup>	( 0.61)	1.95
transitive triplets	2.22**	( 0.73)	3.04
3-cycles	-1.10 <sup>†</sup>	( 0.62)	-1.78
out-out degree( $\hat{1}/2$ ) assortativity	-1.00*	( 0.50)	-2.01
devSeniority alter	0.00	( 0.00)	1.39
devSeniority ego	0.00*	( 0.00)	1.99
devSeniority ego x devSeniority alter	0.00	(31.61)	0.00
devScActivity alter	-0.00	( 0.01)	-0.33
devScActivity ego	0.03	( 0.02)	1.50
devScActivity ego x devScActivity alter	-0.00	( 0.00)	-1.54
int. devScActivity ego x devSeniority ego	-0.00	(31.61)	-0.00

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 17: Estimates parameters in the model for Project ffmpeg (Conflict-Driven). The outdegree density, transitive triplets, reciprocity, 3-cycles and out-out-degree assortativity are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	5.21	(0.54)	-
Rate 2	4.52	(0.43)	-
Rate 3	4.52	(0.46)	-
Rate 4	5.62	(0.62)	-
Rate 5	4.08	(0.40)	-
Rate 6	4.69	(0.54)	-
Rate 7	4.39	(0.53)	-
Rate 8	27.77	(5.86)	-
Rate 9	3.46	(0.32)	-
outdegree (density)	-2.20***	(0.06)	-39.64
reciprocity	-2.52*	(1.27)	-1.98
transitive triplets	3.63***	(0.85)	4.27
3-cycles	-2.80**	(1.05)	-2.67
out-out degree( $\hat{1}/2$ ) assortativity	-1.42***	(0.40)	-3.56
devSeniority ego	N.A.	(N.A.)	-
devScActivity ego	N.A.	(N.A.)	-
int. devScActivity ego x devSeniority ego	N.A.	(N.A.)	-

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 18: Estimates parameters in the model for Project Amarok (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	1.29	(0.33)	-
Rate 2	2.15	(0.63)	-
Rate 3	2.12	(0.63)	-
Rate 4	5.09	(1.79)	-
Rate 5	3.96	(1.16)	-
Rate 6	6.21	(2.27)	-
Rate 7	1.57	(0.39)	-
Rate 8	0.62	(0.20)	-
Rate 9	0.68	(0.21)	-
outdegree (density)	-4.33***	( 0.25)	-17.21
reciprocity	-3.47	( 6.23)	-0.56
transitive triplets	2.91***	( 0.85)	3.42
3-cycles	-2.58	( 8.58)	-0.30
out-out degree( $\hat{1}/2$ ) assortativity	0.14	( 0.38)	0.37
devSeniority alter	-0.00	( 0.00)	-1.01
devSeniority ego	0.00	( 0.00)	1.37
devSeniority ego x devSeniority alter	0.00	(31.61)	0.00
devScActivity alter	-0.01	( 0.01)	-1.02
devScActivity ego	0.01**	( 0.01)	2.79
devScActivity ego x devScActivity alter	-0.00	( 0.00)	-0.26
int. devScActivity ego x devSeniority ego	-0.00	(31.61)	-0.00

†  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;



TABLE 19: Estimates parameters in the model for Project Apache CouchDB (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	1.31	(0.33)	-
Rate 2	2.24	(0.67)	-
Rate 3	2.21	(0.68)	-
Rate 4	5.36	(1.94)	-
Rate 5	4.11	(1.19)	-
Rate 6	6.40	(2.35)	-
Rate 7	1.62	(0.41)	-
Rate 8	0.65	(0.21)	-
Rate 9	0.70	(0.22)	-
outdegree (density)	-4.37***	( 0.26)	-17.02
reciprocity	-3.97	( 7.73)	-0.51
transitive triplets	2.98***	( 0.84)	3.53
3-cycles	-1.57	( 5.86)	-0.27
out-out degree( $\hat{1}/2$ ) assortativity	0.14	( 0.40)	0.34
devSeniority alter	-0.00	( 0.00)	-0.89
devSeniority ego	0.00	( 0.00)	1.20
devSeniority ego x devSeniority alter	-0.00	(31.61)	-0.00
devScActivity alter	-0.01	( 0.01)	-1.04
devScActivity ego	0.02**	( 0.01)	3.02
devScActivity ego x devScActivity alter	-0.00	( 0.00)	-0.36
int. devScActivity ego x devSeniority ego	-0.00	(31.61)	-0.00

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 20: Estimates parameters in the model for Project Pidgin (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	1.31	(0.34)	-
Rate 2	2.22	(0.67)	-
Rate 3	2.17	(0.67)	-
Rate 4	5.20	(1.88)	-
Rate 5	4.07	(1.22)	-
Rate 6	6.35	(2.40)	-
Rate 7	1.59	(0.40)	-
Rate 8	0.63	(0.20)	-
Rate 9	0.70	(0.22)	-
outdegree (density)	-4.38***	( 0.26)	-16.91
reciprocity	-4.93	(10.13)	-0.49
transitive triplets	2.96***	( 0.85)	3.47
3-cycles	-2.87	( 9.06)	-0.32
out-out degree( $\hat{1}/2$ ) assortativity	0.17	( 0.40)	0.43
devSeniority alter	-0.00	( 0.00)	-0.72
devSeniority ego	0.00 <sup>†</sup>	( 0.00)	1.83
devSeniority ego x devSeniority alter	0.00	(31.61)	0.00
devScActivity alter	-0.02	( 0.01)	-1.27
devScActivity ego	0.01*	( 0.01)	2.57
devScActivity ego x devScActivity alter	-0.00	( 0.00)	-0.22
int. devScActivity ego x devSeniority ego	-0.00	(31.61)	-0.00

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 21: Estimates parameters in the model for Project MPlayer (Non-Conflict-Driven). The outdegree density, transitive triplets, and developers' source code activity effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	1.28	(0.33)	-
Rate 2	2.17	(0.65)	-
Rate 3	2.11	(0.66)	-
Rate 4	5.02	(1.77)	-
Rate 5	3.99	(1.16)	-
Rate 6	6.07	(2.21)	-
Rate 7	1.57	(0.39)	-
Rate 8	0.63	(0.19)	-
Rate 9	0.69	(0.21)	-
outdegree (density)	-4.35***	( 0.25)	-17.18
reciprocity	-6.27	( 9.93)	-0.63
transitive triplets	2.79***	( 0.81)	3.44
3-cycles	-2.40	( 4.89)	-0.49
out-out degree( $\hat{1}/2$ ) assortativity	0.24	( 0.36)	0.65
devSeniority alter	-0.00	( 0.00)	-0.68
devSeniority ego	0.00	( 0.00)	1.58
devSeniority ego x devSeniority alter	0.00	(31.61)	0.00
devScActivity alter	-0.02	( 0.01)	-1.29
devScActivity ego	0.01*	( 0.00)	2.47
devScActivity ego x devScActivity alter	-0.00	( 0.00)	-0.30
int. devScActivity ego x devSeniority ego	-0.00	(31.61)	-0.00

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 22: Estimates parameters in the model for Project freeglut (Non-Conflict-Driven). The outdegree density, transitive triplets effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	$t$ stat.
Rate 1	2.32	(0.51)	-
Rate 2	3.40	(0.62)	-
Rate 3	11.19	(4.12)	-
Rate 4	2.65	(0.50)	-
Rate 5	2.36	(0.48)	-
Rate 6	2.36	(0.52)	-
Rate 7	2.07	(0.49)	-
Rate 8	1.58	(0.39)	-
Rate 9	4.06	(0.96)	-
outdegree (density)	-2.80***	(0.22)	-13.01
reciprocity	0.28	(0.28)	1.00
transitive triplets	0.43***	(0.09)	4.84
3-cycles	-0.18	(0.12)	-1.50
out-out degree( $\hat{1}/2$ ) assortativity	0.16	(0.10)	1.61

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 23: Estimates parameters in the model for Project Ceph (Not forked). The outdegree density, reciprocity, transitive triplets, 3-cycles, out-out degree assortativity and developers' alter source code activity effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	7.03	(0.83)	-
Rate 2	5.40	(0.48)	-
Rate 3	5.64	(0.56)	-
Rate 4	6.22	(0.67)	-
Rate 5	6.77	(0.64)	-
Rate 6	22.78	(3.49)	-
Rate 7	11.85	(1.02)	-
Rate 8	7.92	(0.77)	-
Rate 9	7.19	(0.64)	-
outdegree (density)	-2.85***	( 0.08)	-36.95
reciprocity	1.03***	( 0.22)	4.76
transitive triplets	1.85***	( 0.25)	7.27
3-cycles	-2.77***	( 0.42)	-6.59
out-out degree( $\hat{1}/2$ ) assortativity	-0.51***	( 0.11)	-4.59
devSeniority alter	0.00	( 0.00)	0.48
devSeniority ego	0.00	(31.61)	0.00
devSeniority ego x devSeniority alter	-0.00	(31.61)	-0.00
devScActivity alter	-0.02***	( 0.00)	-3.73
devScActivity ego	0.00	( 0.00)	1.41
devScActivity ego x devScActivity alter	0.00	( 0.00)	0.02
int. devScActivity ego x devSeniority ego	-0.00	(31.61)	-0.00

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 24: Estimates parameters in the model for Project OpenStack Neutron (Not forked). The outdegree density, reciprocity, transitive triplets, 3-cycles, out-out degree assortativity and developers' seniority effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	1.94	(0.17)	-
Rate 2	2.75	(0.24)	-
Rate 3	3.25	(0.29)	-
Rate 4	3.00	(0.25)	-
Rate 5	3.09	(0.28)	-
Rate 6	9.52	(1.15)	-
Rate 7	4.68	(0.34)	-
Rate 8	4.39	(0.33)	-
Rate 9	3.73	(0.31)	-
outdegree (density)	-3.45***	(0.06)	-57.09
reciprocity	-3.80***	(0.81)	-4.71
transitive triplets	2.65***	(0.34)	7.80
3-cycles	-2.65*	(1.10)	-2.41
out-out degree( $\hat{1}/2$ ) assortativity	-0.54***	(0.12)	-4.36
devSeniority alter	0.00*	(0.00)	2.03
devSeniority ego	0.00***	(0.00)	3.43
devSeniority ego x devSeniority alter	-0.00	(0.00)	-0.39
devScActivity alter	0.04 <sup>†</sup>	(0.02)	1.94
devScActivity ego	0.09 <sup>†</sup>	(0.05)	1.96
devScActivity ego x devScActivity alter	-0.02	(0.03)	-0.62
int. devScActivity ego x devSeniority ego	-0.00*	(0.00)	-2.04

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

TABLE 25: Estimates parameters in the model for Project GlusterFS (Not forked). The outdegree density, transitive triplets, out-out degree assortativity and developers' alter seniority effects are statistically significant. The interpretations are discussed in section 4.

Effect	par.	(s.e.)	<i>t</i> stat.
Rate 1	4.56	(0.74)	-
Rate 2	3.68	(0.65)	-
Rate 3	4.61	(1.01)	-
Rate 4	6.78	(1.50)	-
Rate 5	3.96	(0.64)	-
Rate 6	2.70	(0.40)	-
Rate 7	2.65	(0.51)	-
Rate 8	5.55	(1.16)	-
Rate 9	8.19	(1.77)	-
outdegree (density)	-2.78***	( 0.12)	-23.25
reciprocity	-0.78	( 0.68)	-1.14
transitive triplets	2.29***	( 0.37)	6.15
3-cycles	-0.88 <sup>†</sup>	( 0.53)	-1.65
out-out degree( $\hat{1}/2$ ) assortativity	-0.54**	( 0.19)	-2.88
devSeniority alter	-0.00***	( 0.00)	-3.38
devSeniority ego	-0.00	( 0.00)	-0.13
devSeniority ego x devSeniority alter	-0.00	(31.61)	-0.00
devScActivity alter	-0.09 <sup>†</sup>	( 0.05)	-1.83
devScActivity ego	0.02	( 0.04)	0.39
devScActivity ego x devScActivity alter	0.01	( 0.01)	0.45
int. devScActivity ego x devSeniority ego	0.00	( 0.00)	0.73

<sup>†</sup>  $p < 0.1$ ; \*  $p < 0.05$ ; \*\*  $p < 0.01$ ; \*\*\*  $p < 0.001$ ;

Table 26 illustrates the distinct patterns associated with Conflict-Driven vs. Non-Conflict-Driven forking categories. Negative three-cycles, positive transitive triplets, and out-out degree assortativity effects were statistically significant in the conflict-driven forks, and not statistically significant in the non-conflict-driven forks. Also, developers source code activity effect was statistically significant for non-conflict-driven forks, and non-significant for conflict-driven forks. The interpretations of these effects are discussed in the following section 4.

TABLE 26: Conflict-Driven vs. Non-Conflict-Driven parameter effects

<b>Effects</b>	<b>Conflict-Driven</b>	<b>Non-Conflict-Driven</b>
Negative three-cycles	Statistically Significant (*)	
Positive transitive triplets	Statistically Significant (*)	
Out-out degree assortativity	Statistically Significant (*)	
Developers source code activity		Statistically Significant (*)

## 4 Discussion

In this research, we investigated the pre-fork time periods for open source projects that had forked. Our aim was to investigate whether there are patterns left in the projects records that could be used to identify the pre-fork dynamics, and use them as key indicators to identify problems, and design an intervention, if needed, before the project disintegrates and the irreversible damage is done.

Several stakeholders benefit from such important research findings, as involve-



ment in open source software development, either as an individual developer, or a business/government sponsor, is an investment of time, energy and resources. The findings of our research helps decision-makers make better informed decisions both when choosing to be affiliated with or involved in certain open source projects, and when there have been involved in the projects.

We analyzed both the contents and the relationship representations of the developers in forked open source software projects, that had gone through various types of forking. Contents of the developer communication messages were analyzed using time series analysis and sentiment analysis, and the interaction patterns were analyzed using social network analysis methods for longitudinal change.

Our research goals were whether we can identify if there is a problem; whether we can do this early enough that we can actually do something about it, and whether we can do this using objectively automate-able processes without having to be intimately familiar with the projects. Our results show that monitoring the communication patterns, both with regard to message contents, and relationship representations address all three research questions.

Our findings show that the conflict-driven forks experienced a period of pre-fork negativity, manifested in the developers' communication sentiments. Our anomaly detection algorithms detected these anomalies. In contrast, the non-conflict-driven forks did not exhibit any anomalies, or showed periods of

positive spikes, which is indicative of positive communication sentiments.

Our findings also show that in conflict-driven forks, the social graph tended to have local hierarchies [52], which was manifested as statistically significant negative three-cycles effects combined with statistically significant positive transitive triplets effects. This indicated several local hierarchies within the project community, which depending on the context, can be an indicator of several factions.

Our findings also show that conflict-driven forks had communities where there existed a tendency for heavy mailing-list poster developers to be interacting with to other heavy mailing-list poster developers, which was manifested as a statistically significant out-out degree assortativity effect. In contrast, non-conflict-driven forks has communities where there existed a tendency for heavy code contributors to be interacting with other heavy code contributors, which was manifested as a statistically significant developer source code activity effect. This can be interpreted as a unhealthy vs. healthy dynamic, whereas, in the non-conflict-driven forked projects more code contributions translated into more conversations, so, the communication was backed by solid code. Whereas, in conflict-driven forked projects, the could have been much talk and not much code. One cannot help but think of Linus Torvald’s quote *“Talk is cheap. Show me the code.”*

In summary, we represented longitudinal dynamics in conflict-driven vs.

non-conflict-driven forks. We suggested indicators of unhealthy dynamics to suggest intervention. We detected pre-fork anomalies suggestive of negative communication sentiments, and we expressed underlying properties and values of community behavior as model effects with their statistical significance. Further research is needed to be able to generalize. This is a good starting point for gaining an understanding of longitudinal changes of underlying properties of an open source projects' community.

## 5 Threats to Validity

One of the threats to validity is the size of the data we had access to and could attain. The projects in this research study were selected from a pool of projects, based on a filtering criteria that included availability of their data. Given access, a larger number of projects as the sample size could result in a more robust investigation, that may give a better external validity and the power to generalize. This is an important limitation of our study, because even though, our quantitative analysis featured more projects than typical studies in this domain, still, it is not possible to infer and generalize the results and expect the predictions to be valid for all open source projects.

Secondly, there is an inherent limitation with regard to any study that uses data from online communications, as the sole source of communication. It is possible for the developers to have had interactions in person, or via other private or unattainable means of communication. The assumption that all the communication can be captured by mining repositories is intuitively imperfect, but inevitable.

Additionally, social interactions data is noisy, and our study might be affected by the noise. The trade-off between representativeness and signal-noise ratio is another limiting decision in the context of our work. Hence, our approach findings are not exact and clean signals, but rather an approximation.

Third, the statistical model we use to model the longitudinal evolution of collaboration networks is estimated stochastically, rather than being calculated exactly. The stochastic process might not always arrive at the same results. To

counter this issue, we ran the algorithm several times to double-check for such irregularities.

## **Acknowledgments**

I would like to thank my academic adviser, Prof. Carlos Jensen, and Prof. Drew Gerley, Prof. Ron Metoyer, Prof. Cindy Grimm, Prof. Alex Groce and Prof. Chris Scaffidi for their help. I would also like to thank Derric Jacobs, of the Sociology Department at Oregon State University, for his help with search for theories and lending his books to me; and again, many thanks goes to Prof. Drew Gerkey, of the Department of Anthropology at Oregon State University for his advice and pointers to theories about human behavior in social sciences. I would also like to thank the open source developers of the projects studied for making their data available, without which this study would not have been possible.

# Appendices

## A Appendix A: List of projects forked for undesirable reasons

TABLE 27: List of all projects forked because of “*personal differences among the developer team*” conflict-driven [47] in chronological order. N/A means Not Applicable

Original	Forked	Date	M.L. data accessi- ble?	Collected?
GNU Emacs	X Emacs	1991, ?	No, only after 2000	N/A
NetBSD	OpenBSD	1995, Oct	Yes, scarce, unusable	N/A
xMule	aMule	2003, Aug	No, only 2006-2007	N/A
lMule	xMule	2003, Jun	No	N/A
Sodipodi	Inkscape	2003, Nov	No	N/A
Nucleus CMS	Blog:CMS	2004, May	No, only after 09/2004	N/A
BMP	Audacious	2005, Oct	No	N/A
ntfsprogs	NTFS-3G	2006, Jul	No	N/A
OpenWRT	FreeWRT	2006, May	No, only after 10/2006	N/A
QtiPlot	SciDavis	2007, Aug	No	N/A
Kamailio	OpenSIPS	2008, Aug	Yes	Yes
Blastwave.org	OpenCSW	2008, Aug	No	N/A
jMonkeyEngine	Ardor3D	2008, Sept	Yes, scarce, unusable	N/A
Frog CMS	Wolf CMS	2009, Jul	No	N/A
Aldrin	Neil	2009, ?	No	N/A
Ffmpeg	libav	2011, Mar	Yes	Yes

TABLE 28: List of all projects forked because of the need for “*more community-driven development*” conflict-driven [47] in chronological order. N/A means Not Applicable

Original	Forked	Date	M.L. data accessible?	Collected?
Nethack	Slash’EM	1996, ?	No	N/A
GCC	EGCS	1997, ?	No	N/A
SourceForge	Savane	2001, Oct	No	N/A
PHPNuke	PostNuke	2001, Sum	No, not found	N/A
QTExtended	OPIE	2002, May	No	N/A
GraphicsMagick	Graphics	2002, Nov	No, only after 2003	N/A
freelut	OpenGLUT	2004, Mar	Yes	Yes
Mambo	Joomla!	2005, Aug	No	N/A
SER	Kamailio	2005, Jun	No, only after 2006	N/A
PHPNuke	RavenNuke	2005, Nov	No, not found	N/A
Hula	Bongo	2006, Dec	No	N/A
Compiere	A Dempiere	2006, Sept	No	N/A
Compiz	Beryl	2006, Sept	No, only after 06/2007	N/A
SQL-Ledger	LedgerSMB	2006, Sept	No	N/A
Asterisk	Callweaver	2007, Jun	Yes	Yes
CodeIgniter	KohanaPHP	2007, May	No, not found	N/A
OpenOffice.org	Go-oo.org	2007, Oct	No, only after 06/2011	N/A
Mambo	MiaCMS	2008, May	No	N/A
TORCS	Speed Dreams	2008, Nov	Yes, scarce, unusable	N/A
MySQL	MariaDB	2009, Jan	Yes, too large, unusable	N/A
Nagios	Icinga	2009, May	No	N/A
Project Dark-star	RedDwarf	2010, Feb	Yes, scarce, unusable	N/A
SysCP	Froxlor	2010, Feb	No	N/A
Dokeos	Chamilo	2010, Jan	No, not found	N/A
GNU Zebra	Quagga	2010, Jul	No	N/A
rdesktop	FreeRDP	2010, Mar	Yes	Yes
OpenOffice.org	LibreOffice	2010, Sept	No, only after 06/2011	N/A
Redmine	ChiliProject	2011, Feb	Yes, scarce, unusable	N/A

## B Appendix B: Initial study: Temporal analysis using the network-specific measurement approach

In our initial study [4][5][6], which was a network-specific study, we wanted to analyze the network-specific changes that happen to the community over a given period of time, e.g., three months before and three months after the year in which the forking event happened. For this network-specific study, we measured the betweenness centrality [12] of the most significant nodes in the graph, and the graph diameter over time. Figures 21, 22, and 23 show the betweenness centralities over the 1.5 year period for the Kamilio, Amarok and Asterisk projects respectively. To do temporal analysis, we had two options; 1) look at snapshots of the network state over time, (e.g., to look at the network snapshots in every week, the same way that a video is composed of many consecutive frames), and 2) look at a period through a time window. We preferred the second approach, and looked through a time window three months wide with 1.5 month overlaps. To create the visualizations, we used a 3 months time frame that progressed six days a frame. In this way, we would have had a relatively smooth transition.

There are many ways of looking at an individual's importance/prestige/status within a network. One is called *closeness centrality*. The *farness* of a node is defined as the sum of its distances to all other nodes. The *closeness* of a node is defined as the inverse of the farness. More informally, the more central a node is the lower its total distance to all other nodes. *Closeness centrality* can be used as a measure of how fast information will spread through the network [13]. Secondly, if we are looking for people who can serve as bridges between



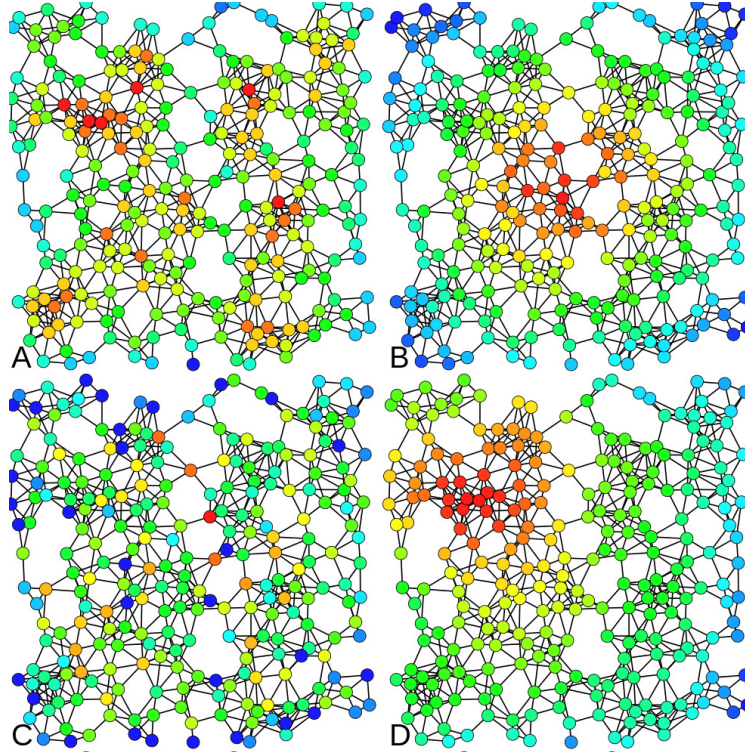


FIGURE 17: Heat-map color-coded examples are shown above. Nodes with higher centrality metric are colored with warmer color: red is the warmest color here. The same network is analyzed four times with the following centrality measures: A) Degree centrality, B) Closeness centrality, C) Betweenness centrality and D) Eigenvector centrality [48]

two distinct communities, we could measure the node's *betweenness centrality*. Betweenness centralities for mediators who act as intermediate entities between other nodes are higher [13]. Third, if cross-community collaboration is the focus, we can measure *edge betweenness centrality*. Edges connecting nodes from different communities have higher edge centrality values. In the community collaboration graph, edge betweenness or stress of an edge is the number of these shortest paths that the edge belongs to, considering all shortest paths between

all pairs of nodes in the graph. Fourth, one can claim that certain people in the community are more important than others, and whoever is close to them, is relatively more important than others. In graph terms, this is measured by *eigenvector centrality*, which is based on the assumption that connections to high-profile nodes contribute more to the importance of a node. Google’s PageRank link-analysis algorithm [43] is a variant of the eigenvector centrality measure. In short, centrality measures have been used in several studies to identify key player in a community.

In addition to the centrality measures, we planned to look into the *resilience* of the community as well. By resilience, we mean how well the network holds its structure and form when some parts of it are deleted, added, or changed. For a graph, the resilience of a graph is a measure of its robustness to node or edge failures. This could occur for instance when an influential member of the community leaves. Many real-world graphs are resilient to random failures but vulnerable to targeted attacks. Resilience can be related to the *graph diameter*: a graph whose diameter does not increase much on node or edge removal has higher resilience [13].

## B.1 Visualization

Several visualization techniques and tools are used in the field of social network analysis, for instance, Gephi [9], which is a FOSS tool for exploring and manipulating networks. It is capable of handling large networks with more than 20,000 nodes and features several SNA algorithms. We used it for dynamic network visualization. We visualized the dynamic network changes using Gephi [9].

The videos<sup>1</sup> show how the community graph is structured, using a continuous force-directed linear-linear model, in which the nodes are positioned near or far from one another proportional to the graph distance between them. This results in a graph shape between Fruchterman & Rheingold’s [22] layout and Noack’s LinLog [37].

## B.2 Initial study results and discussion

### B.2.1 Kamilio Project

Figure 20 shows four key frames from the Kamilio project’s social graph around the time of their fork (the events described here are easier to fully grasp by watching the video). A node’s size is proportional to the number of interactions the node (contributor) has had within the study period and the position and edges of the nodes change if they had interactions within the time window shown, with six day steps per frame. The 1 minute and 37 seconds video shows the life of the Kamilio project between October 2007, and March 2009. Nodes are colored based on the modularity of the network. The community starts with the GeneralList as the the biggest node, and four larger core contributors and three lesser size core contributors. The big red-colored node’s transitions are hard to miss, as this major contributor departs from the core to the periphery of the network (Video minute 1:02) and then leaves the community (Video minute 1:24) capturing either a conflict or retirement. This corresponds to the personal difference category of forking reasons.

---

<sup>1</sup>Video visualizations available at <http://eecs.oregonstate.edu/~azarbaam/OSS2014/>

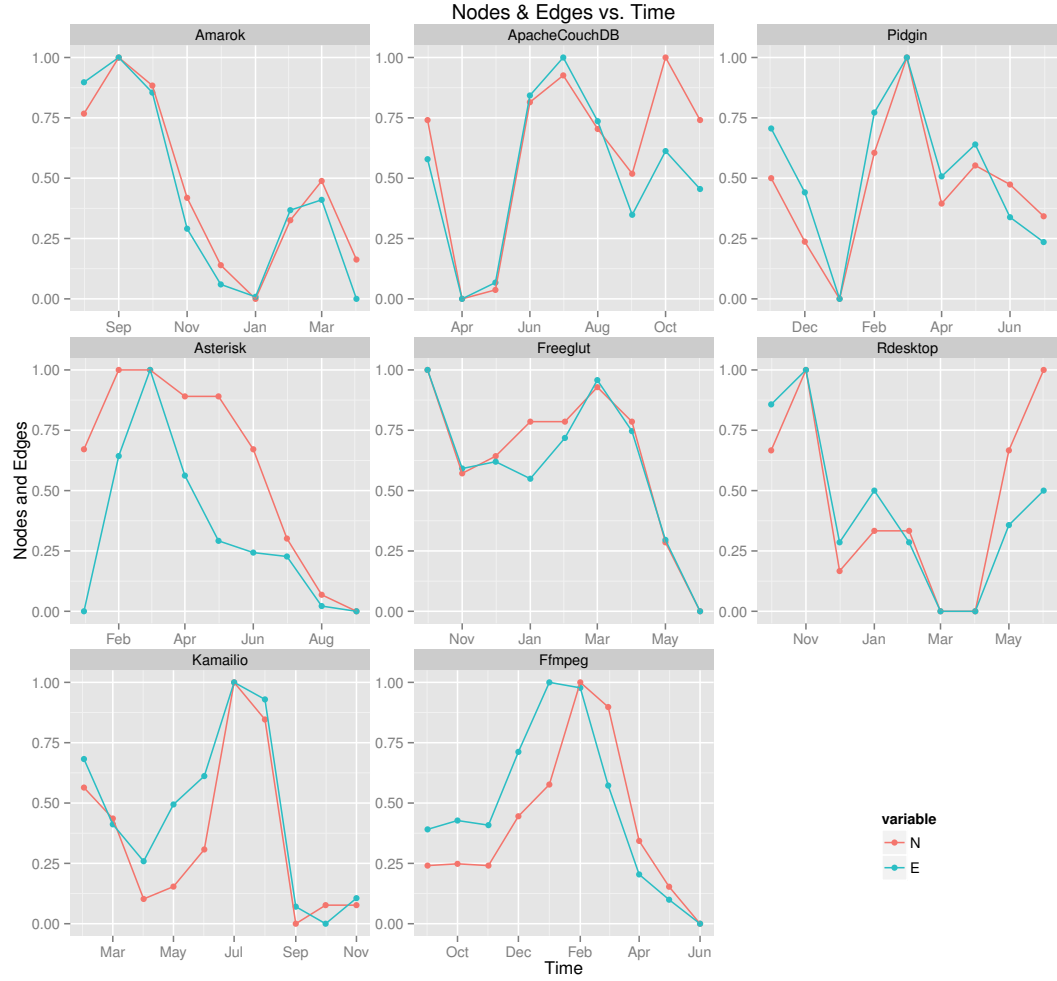


FIGURE 18: Nodes and Edges over Time. The number of nodes and number of edges are normalized to the range  $[0,1]$  to make comparison across projects meaningful, by emphasizing change in ratio, rather than the varying counts. Hence, the measurements were normalized for drawing this graph. The three projects in the first row belong to the “*technical differences*” forking reason category, the three projects in the second row belong to the “*more community-driven development*” forking reason category, and the two projects in the third row belong to the “*personal differences*” forking reason category.

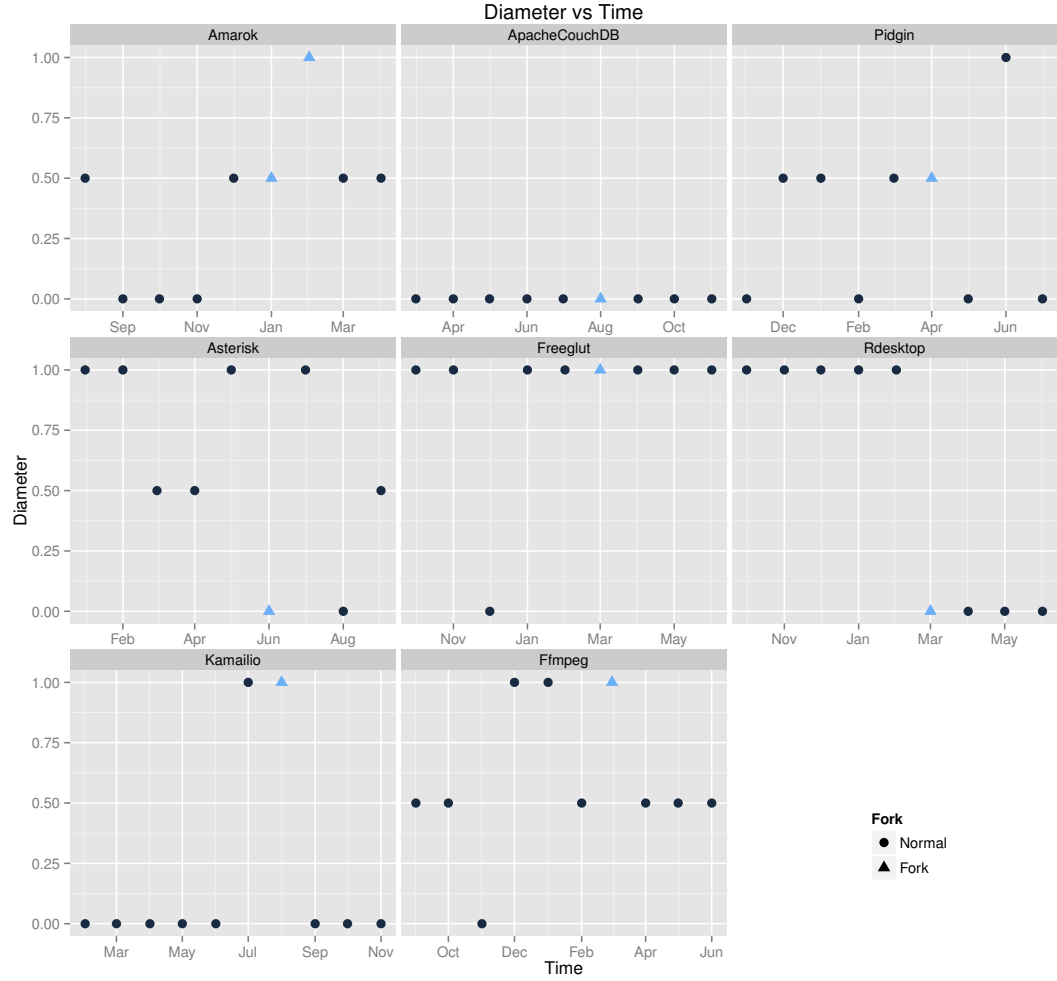


FIGURE 19: Diameter changes over Time. Note that the diameter measurements were normalized to the range  $[0,1]$ . The three projects in the first row belong to the “*technical differences*” forking reason category, the three projects in the second row belong to the “*more community-driven development*” forking reason category, and the two projects in the third row belong to the “*personal differences*” forking reason category.

Figure 21 shows the betweenness centrality of the major contributors of Kamailio project over the same time period. The horizontal axis marks the dates, (each mark represents a 3-month time window with 1.5 months overlap). The vertical

The stacked area chart displays the percentage of the top 1000 email addresses belonging to various domains over time. The x-axis represents time periods from October 2007 to March 2009. The y-axis represents the percentage of email addresses, ranging from 0% to 100%. The legend lists 30 domains, with 'gmail.com' (blue) showing a dramatic increase in share starting around June 2008, eventually accounting for nearly 50% of the top 1000 addresses. Other prominent domains include 'hotmail.com' (orange), 'yahoo.com' (yellow), and 'aol.com' (grey). Below the main chart, a bar chart shows the number of unique domains in the top 10 for each time period, with values ranging from 4 to 5.

Time Period	Number of Unique Domains in Top 10
Oct 2007 - Dec 2007	5
Nov 15, 2007 - Feb 15, 2008	5
Jan 2008 - Mar 2008	4
Feb 15, 2008 - May 15, 2008	4
Apr 2008 - Jun 2008	4
May 15, 2008 - Aug 15, 2008	5
July 2008 - Oct 2008	5
Aug 15, 2008 - Nov 15, 2008	5
Oct 2008 - Dec 2008	4
Nov 15, 2008 - Feb 15, 2009	4
Jan 2009 - Mar 2009	4

axis shows the percentage of the top betweenness centralities for each node. The saliency of the GeneralList – colored as light blue – is apparent because of its continuous and dominant presence in the stacked area chart. The chart legend

lists the contributors based on the color and in the same order of appearance on the chart starting from the bottom. Around the "Aug. 15, 2008 - Nov. 15, 2008" tick mark on the horizontal axis, several contributors' betweenness centralities shrink to almost zero and disappear. This suggests the date of fork with a month accuracy. The network diameter of the Kamilio project over the same time period is also shown in Figure 21. An increase in the network diameter during this period is noticeable; this coincides with findings of Hannemann and Klamma [24].

This technique can be used to identify the people involved in conflict and the date the fork happened with a months accuracy, even if the rival project does not emerge immediately.

### **B.2.2 Amarok Project**

The video for the Amarok project fork is available online<sup>2</sup>, and the results from our quantitative analysis of the betweenness centralities and the network diameters are shown in Figure 22. The results show that the network diameter has not increased over the period of the fork, which shows a resilient network. The video shows the dynamic changes in the network structure, again typical of a non-unhealthy network, rather than of simmering conflict. These indicators suggest that the Amarok fork in 2010 belongs to the "addition of technical functionality" rationale for forking, as there are no visible social conflict.

---

<sup>2</sup>Video visualizations available at <http://eecs.oregonstate.edu/~azarbaam/OSS2014/>

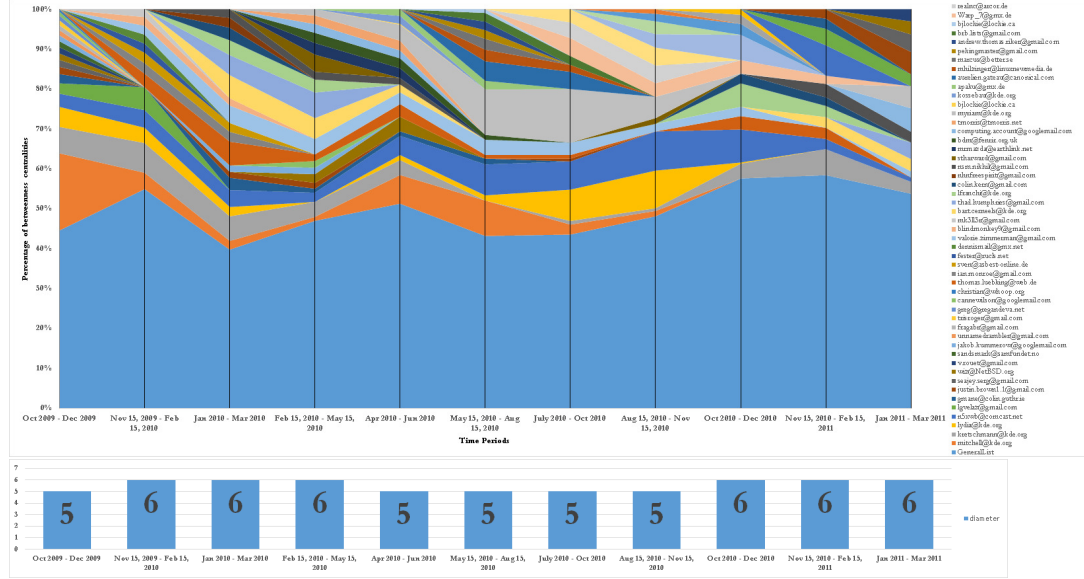


FIGURE 22: Amarak project’s top contributors’ betweenness centralities and network diameter over time between Oct. 2009 to Mar. 2011 in 3-months time windows with 1.5 months overlaps

### B.2.3 Asterisk Project

The video for the Asterisk project is also available online<sup>3</sup>, and the results from our quantitative analysis of the betweenness centralities and the network diameters are shown in Figure 23. The results show that the network diameter remained steady at 6 throughout the period. The Asterisk community was by far the most crowded project, with 932 nodes and 4282 edges. The stacked area chart shows the distribution of centralities, where we see an 80%-20% distribution (i.e., 80% or more of the activity is attributed to six major players, with

<sup>3</sup>Video visualizations available at <http://eecs.oregonstate.edu/~azarbaam/OSS2014/>



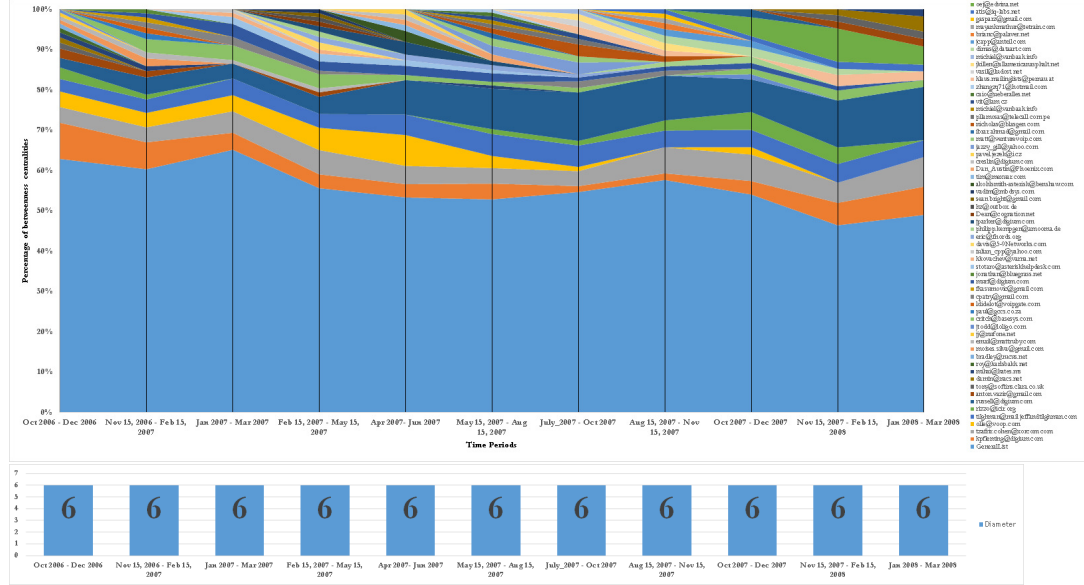


FIGURE 23: Asterisk project’s top contributors’ betweenness centralities and network diameter over time between Oct. 2009 to Mar. 2011 in 3-months time windows with 1.5 months overlaps

the rest of the community accounting for only 20%). This is evident in the video representation as well, as the top-level structure of the network holds throughout the time period. The results from the visual and quantitative analysis links the Asterisk fork to the more community-driven category of forking reasons.

#### B.2.4 Initial study conclusion

We studied the collaboration networks of three FOSS projects using a combination of temporal visualization and quantitative analysis. We based our study on two papers by Robles and Gonzalez-Barahona [47] and Hannemann and Klammer [24], and identified three projects that had forked in the recent past. We mined the collaboration data, formed dynamic collaboration graphs, and measured social network metrics over an 18-month period time window.

We also visualized the dynamic graphs (available online) and as stacked area charts over time. The visualizations and the quantitative results showed the differences among the projects in the three forking reasons of personal differences among the developer teams, technical differences (addition of new functionality) and more community-driven development. The novelty of the approach was in applying the network-specific *temporal* analysis rather than static analysis, and in the temporal visualization of community structure. We showed that this approach shed light on the structure of these projects and reveal information that cannot be seen otherwise.

More importantly, the initial study showed the limitations of a *network-specific* approach, and hence, we adopted a *population-processes* approach for our main study as explained in section 2.

## C Appendix C: Mathematical Definition of Effects in the Statistical Model

### C.1 Structural Effects for the Objective Function

1. Reciprocity

$$\sum_j x_{ij}x_{ji}$$

2. Closure effect: Transitive triplets

$$\sum_{j,h} x_{ih}x_{ij}x_{jh}$$

3. Closure effect: Transitive ties

$$\sum_j x_{ij}max_h(x_{ih}x_{hj})$$

4. Closure effect: Balance

$$\sum_{j=1}^n x_{ij} \sum_{h=1, h \neq i, j}^n |(b_0 - x_{ih} - s_{jh})|$$

$b_0$  is a constant equal to the mean of  $|x_{ih} - x_{jh}|$

5. Closure effect: The number of developers at distance two

$$\sum_j (1 - x_{ij})max_h(x_{ih}x_{hj})$$

6. Three-cycles

$$\sum_j x_{ij} \sum_h x_{jh} x_{hi}$$

7. Betweenness

$$\sum_{j,h} x_{hj} x_{ij} (1 - x_{hj})$$

8. Density (or out-degree)

$$\sum_j x_{ij}$$

9. Activity: Out-degree

$$\sum_j x_{ij} x_{j+} = \sum_j x_{ij} \sum_h x_{jh}$$

10. Covariate V-related popularity

$$\sum_j x_{ij} v_j$$

11. Covariate V-related activity

$$v_i x_{i+}$$

12. Covariate V-related dissimilarity

$$\sum_j x_{ij} |v_i - v_j|$$

13. in-in degree assortativity

$$\sum_j x_{ij} x_{+i}^{1/c} x_{+j}^{1/c}$$

with  $c = 1$  or  $2$

14. in-out degree assortativity

$$\sum_j x_{ij} x_{+i}^{1/c} x_{j+}^{1/c}$$

with  $c = 1$  or  $2$

15. out-in degree assortativity

$$\sum_j x_{ij} x_{i+}^{1/c} x_{+j}^{1/c}$$

with  $c = 1$  or  $2$

16. out-out degree assortativity

$$\sum_j x_{ij} x_{i+}^{1/c} x_{j+}^{1/c}$$

with  $c = 1$  or  $2$

## C.2 Behavior-related Effects

For a behavior variable  $s_{ik}^z(x, z)$ , the following formulas can be used.

1. Shape  $z_i$
2. Quadratic Shape  $z_i^2$
3. Total Similarity

$$\sum_j x_{ij} (sim_{ij}^z - \bar{sim}^z)$$

where

$$sim_{ij}^z = (1 - |z_i - z_j| / max_{ij} |z_i - z_j|)$$

4. Average Similarity

$$\frac{1}{x_{i+}} \sum_j x_{ij} (sim_{ij}^z - \bar{sim}^z)$$

where

$$sim_{ij}^z = (1 - |z_i - z_j| / max_{ij} |z_i - z_j|)$$

5. Average alter

$$\frac{z_i (\sum_j x_{ij} z_j)}{\sum_j x_{ij}}$$

## Bibliography

1. Asur, S., S. Parthasarathy, and D. Ucar, (2009), “*An event-based framework for characterizing the evolutionary behavior of interaction graphs*,” ACM Trans. Knowledge Discovery Data. 3, 4, Article 16, (November 2009), 36 pages. 2009.
2. Azarbakht, A. and C. Jensen, “*Drawing the Big Picture: Temporal Visualization of Dynamic Collaboration Graphs of OSS Software Forks*,” Proc. 10th Int’l. Conf. Open Source Systems, 2014.
3. Azarbakht, A. and C. Jensen, “*Temporal Visualization of Dynamic Collaboration Graphs of OSS Software Forks*,” Proc. Int’l. Network for Social Network Analysis (INSNA) Sunbelt XXXIV Conf., 2014.
4. Azarbakht, A., “*Drawing the Big Picture: Analyzing FLOSS Collaboration with Temporal Social Network Analysis*,” Proc. 9th Int’l. Symp. Open Collaboration, ACM, 2013.
5. Azarbakht, A. and C. Jensen, “*Analyzing FOSS Collaboration & Social Dynamics with Temporal Social Networks*,” Proc. 9th Int’l. Conf. Open Source Systems Doct. Cons., 2013.
6. Azarbakht, A., “*Temporal Visualization of Collaborative Software Development in FOSS Forks*,” Proc. IEEE Symp. Visual Languages and Human-Centric Computing, 2014.
7. Azarbakht, E.A. and C. Jensen, “*Longitudinal Analysis of the Run-up to a Decision to Break-up (Fork) in a Community*,” Proc. 13th IFIP International Conference on Open Source Systems. Springer, Cham, 2017.
8. Baishakhi R., C. Wiley, and M. Kim, “*REPERTOIRE: a cross-system porting analysis tool for forked software projects*,” Proc. ACM SIGSOFT 20th Int’l. Symp. Foundations of Software Engineering, ACM, 2012.
9. Bastian, M., S. Heymann, and M. Jacomy, “*Gephi: an open source software for exploring and manipulating networks*,” Int’l AAAI Conf. on Weblogs and Social Media, 2009.
10. Bezrukova, K., C. S. Spell, J. L. Perry, “*Violent Splits Or Healthy Divides? Coping With Injustice Through Faultlines*,” Personnel Psychology, Vol 63, Issue 3. 2010.

11. Bird, C., D. Pattison, R. D'Souza, V. Filkov, and P. Devanbu, "*Latent social structure in open source projects*," Proc. 16th ACM SIGSOFT Int'l. Symposium on Foundations of software engineering, ACM, 2008.
12. Brandes, U. "*A Faster Algorithm for Betweenness Centrality*", Journal of Mathematical Sociology 25(2):163-177, 2001.
13. Chakrabarti, D., and C. Faloutsos. "*Graph mining: Laws, generators, and algorithms*," ACM Computing Surveys, 38, 1, Article 2, 2006.
14. Chen, C. and Liu, Lon-Mu, "*Joint Estimation of Model Parameters and Outlier Effects in Time Series*," Journal of the American Statistical Association, 88, 284-297. 1993.
15. Coleman, J.S. "*Introduction to Mathematical Sociology*," New York etc.: The Free Press of Glencoe. 1964.
16. Crowston, K., K. Wei, J. Howison, and A. Wiggins. "*Free/Libre open-source software development: What we know and what we do not know*," ACM Computing Surveys, 44, 2, Article 7, 2012.
17. Davidson, J, R. Naik, A. Mannan, A. Azarbakht, C. Jensen, "*On older adults in free/open source software: reflections of contributors and community leaders*," Proc. IEEE Symp. Visual Languages and Human-Centric Computing, 2014.
18. Ernst, N., S. Easterbrook, and J. Mylopoulos, "*Code forking in open-source software: a requirements perspective*," arXiv preprint arXiv:1004.2889, 2010.
19. Feuerriegel S. and N. Proellocks. "*SentimentAnalysis: Dictionary-based sentiment analysis*", R package version 1.1-0. <https://github.com/sfeuerriegel/SentimentAnalysis>. 2016.
20. Ford, L. R. and D. R. Folkerson, "*A simple algorithm for finding maximal network flows and an application to the Hitchcock problem*," Canadian Journal of Mathematics, vol. 9, pp. 210-218, 1957.
21. Forrest, D., C. Jensen, N. Mohan, and J. Davidson, "*Exploring the Role of Outside Organizations in Free/ Open Source Software Projects*," Proc. 8th Int'l. Conf. Open Source Systems, 2012.



22. Fruchterman, T. M. J. and E. M. Reingold, “*Graph drawing by force-directed placement*,” *Softw: Pract. Exper.*, vol. 21, no. 11, pp. 1129-1164, 1991.
23. Guzzi, A., A. Bacchelli, M. Lanza, M. Pinzger, and A. van Deursen. “*Communication in open source software development mailing lists*,” *Proc. 10th Conf. on Mining Software Repositories*, IEEE Press, 2013.
24. Hannemann, A and , R. Klamma “*Community Dynamics in Open Source Software Projects: Aging and Social Reshaping*,” *Proc. Int. Conf. on Open Source Systems*, 2013.
25. Heider, F. *The Psychology of Interpersonal Relations*. John Wiley & Sons. 1958.
26. Howison, J. and K. Crowston. “*The perils and pitfalls of mining SourceForge*,” *Proc. Int’l. Workshop on Mining Software Repositories*, 2004.
27. Howison, J., K. Inoue, and K. Crowston, “*Social dynamics of free and open source team communications*,” *Proc. Int’l. Conf. Open Source Systems*, 2006.
28. Howison, J., M. Conklin, and K. Crowston, “*FLOSSmole: A collaborative repository for FLOSS research data and analyses*,” *Int’l. Journal of Information Technology and Web Engineering*, 1(3), 17-26. 2006.
29. Krivitsky, P. N., and M. S. Handcock. “*A separable model for dynamic networks*,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76, no. 1: 29-46. 2014.
30. Kuechler, V., C. Gilbertson, and C. Jensen, “*Gender Differences in Early Free and Open Source Software Joining Process*,” *Open Source Systems: Long-Term Sustainability*, 2012.
31. Kunegis, J., S. Sizov, F. Schwagereit, and D. Fay, “*Diversity dynamics in online networks*,” *Proc. 23rd ACM Conf. on Hypertext and Social Media*, 2012.
32. Leskovec, J., Kleinberg, J., and Faloutsos, C.: “*Graphs over time: densification laws, shrinking diameters and possible explanations*,” *Proc. SIGKDD Int’l. Conf. Knowledge Discovery and data Mining*, 2005.

33. Leskovec, J., K. J. Lang, A. Dasgupta, and M. W. Mahoney, “*Statistical properties of community structure in large social and information networks*,” Proc. 17th Int’l. Conf. World Wide Web, ACM, 2008.
34. Lopez-de-Lacalle, J. “*tsoutliers: Detection of Outliers in Time Series*”, R package version 0.6-5. <https://CRAN.R-project.org/package=tsoutliers>, 2016.
35. Nakakoji, K., Y. Yamamoto, Y. Nishinaka, K. Kishida, and Y. Ye. “*Evolution patterns of open-source software systems and communities*,” Proc. Int’l. Workshop Principles of Software Evolution, ACM, 2002.
36. Mikkonen, T., L. Nyman, “*To Fork or Not to Fork: Fork Motivations in SourceForge Projects*,” Int’l. J. Open Source Softw. Process. 3, 3. July, 2011.
37. Noack, A., “*Energy models for graph clustering*,” J. Graph Algorithms Appl., vol. 11, no. 2, pp. 453-480, 2007.
38. Nowak, M. A. “*Five rules for the evolution of cooperation*,” Science 314, No. 5805: 1560-1563. 2006.
39. Nyman, L. , “*Understanding code forking in open source software*,” Proc. 7th Int’l. Conf. Open Source Systems Doct. Cons., 2011.
40. Nyman, L., T. Mikkonen, J. Lindman, and M. Fougère, “*Forking: the invisible hand of sustainability in open source software*,” Proc. SOS 2011: Towards Sustainable Open Source, 2011.
41. Nyman, L., “*Hackers on Forking*,” Proc. Int’l. Symp. on Open Collaboration, 2014.
42. Oh, W., Jeon, S., “*Membership Dynamics and Network Stability in the Open-Source Community: The Ising Perspective*” Proc. 25th Int’l. Conf. Information Systems. 2004.
43. Page, B, B. Sergey, R. Motwani and T. Winograd, “*The PageRank Citation Ranking: Bringing Order to the Web*,” Technical Report, Stanford InfoLab, 1999.
44. Proelochs, Feuerriegel and Neumann: “*Generating Domain-Specific Dictionaries Using Bayesian Learning*”, Proceedings of the 23rd European Conference on Information Systems (ECIS 2015), Muenster, Germany, 2015.

45. R Core Team. "*R: A language and environment for statistical computing. R Foundation for Statistical Computing*", Vienna, Austria. URL <https://www.R-project.org/>. 2016.
46. Robins, G., P. Pattison, Y. Kalish, and D. Lusher. "*An introduction to exponential random graph ( $p^*$ ) models for social networks*," Social networks 29, no. 2: 173-191. 2007.
47. Robles, G. and J. M. Gonzalez-Barahona, "*A comprehensive study of software forks: Dates, reasons and outcomes*," Proc. 8th Int'l. Conf. Open Source Systems, 2012.
48. Rocchini, C. (Nov. 27 2012), Wikimedia Commons, Available: <http://en.wikipedia.org/wiki/File:Centrality.svg>, 2012.
49. Singer, L., F. Figueira Filho, B. Cleary, C. Treude, M. Storey, and K. Schneider. "*Mutual assessment in the social programmer ecosystem: an empirical investigation of developer profile aggregators*," Proc. Conf. Computer supported cooperative work, ACM, 2013.
50. Snijders, T. AB. "*Markov chain Monte Carlo estimation of exponential random graph models*," Journal of Social Structure 3, no. 2: 1-40. 2002.
51. Snijders, Tom AB. "*Models for longitudinal network data*," Models and methods in social network analysis 1: 215-247. 2005.
52. Snijders, Tom AB., GG Van de Bunt, CEG Steglich, "*Introduction to stochastic actor-based models for network dynamics*," Social networks 32 (1), 44-60. 2010.
53. Sowe, S., L. Stamelos, and L. Angelis, "*Identifying knowledge brokers that yield software engineering knowledge in OSS projects*," Information and Software Technology, vol. 48, pp. 1025-1033, Nov 2006.
54. Spence, M. "*Job market signaling*," Quarterly Journal of Economics, 87: 355-374. 1973.
55. Steglich, C., T. AB Snijders, and M. Pearson. "*Dynamic networks and behavior: Separating selection from influence*," Sociological methodology 40, no. 1: 329-393. 2010.
56. Storey, M., L. Singer, B. Cleary, F. Figueira Filho, and A. Zagalsky, "*The (R) Evolution of social media in software engineering*," Proc. Future of Software Engineering, ACM, 2014.

57. Syeed, M. M., “*Socio-Technical Dependencies in Forked OSS Projects: Evidence from the BSD Family*,” Journal of Software 9.11 (2014): 2895-2909. 2014.
58. Teixeira, J., and T. Lin, “*Collaboration in the open-source arena: the webkit case*,” Proc. 52nd ACM conf. Computers and people research (SIGSIM-CPR '14). ACM, 2014.
59. Torres, M. R. M., S. L. Toral, M. Perales, and F. Barrero, “*Analysis of the Core Team Role in Open Source Communities*,” Int. Conf. on Complex, Intelligent and Software Intensive Systems, IEEE, 2011.
60. Zachary, W., “*An information flow model for conflict and fission in small groups*,” Journal of Anthropological Research, vol. 33, no. 4, pp. 452-473, 1977.