

# Temporal Analysis of Dynamic Collaboration Graphs of Open Source Software Development: Forking

Amir Azarbakht

School of Electrical Engineering and Computer Science

Oregon State University

Corvallis, Oregon 97331

Email: azarbaam@eecs.oregonstate.edu



**Abstract**—How can we understand FOSS collaboration better? Can social issues that emerge be identified and addressed as they happen? Can the community heal itself, become more transparent and inclusive, and promote diversity? We propose a technique to address these issues by quantitative analysis and temporal visualization of social dynamics in FOSS communities. We propose using social network analysis to identify unhealthy dynamics; this will help predict formation of unhealthy dynamics and which gives the community a heads-up when they can still take action to ensure the sustainability of the project.

**Keywords.** Free/Open Source Software, Social Dynamics, Temporal Analysis, Forking, Visualization, Temporal Visualization, Social Network Analysis, FOSS, FLOSS.

## 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. Free/Open Source Software (FOSS) development efforts are prime examples of how community can be leveraged in software development, groups are formed around communities of interest, and depend on continued interest and involvement in order to stay alive [19].

Though the bulk of collaboration and communication in FOSS communities occurs online and is publicly accessible, there are many open questions about the social dynamics in FOSS communities. Projects might 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 might lead to forking of the project and dilution of the community. Forking, either as a violent split when there is a conflict or as a friendly divide when new features are experimentally added both affect the community [3].

Most recent studies of FOSS communities have tended to suffer from an important limitation. They

treat community as a static structure rather than a dynamic process. In this study, we propose to use temporal social network analysis to study the evolution and social dynamics of FOSS communities. With these techniques we aim to identify measures associated with unhealthy group dynamics, e.g. a simmering conflict, as well as early indicators of major events in the lifespan of a community. One set of dynamics we are especially interested in, are those that lead FOSS projects to fork. We used the results of a study of forked FOSS projects by Robles and Gonzalez-Barahona [23] as the starting point for our study, and tried to gain a better understanding of the evolution of these communities.

This paper is organized as follows: We present related literature on online social communities. We then present the gap in the literature, and discuss why the issue needs to be addressed. After that, in methodology, we describe how data gathering, the analysis, and the visualization of the findings is proposed to be carried out. At the end, we present preliminary results, discussion and threats to validity.

## 2 RELATED WORK

The social structures of FOSS communities have been studied extensively. Researchers have studied the social structure and dynamics of team communications [4][12][13], identifying knowledge brokers and associated activities [25], project sustainability [20], forking [19], their topology [4], their demographic diversity [15], gender differences in the process of joining them [14] and the role of the core team in their communities [26], etc. All of these studies have tended to look at community as a static structure rather than a dynamic process. This makes it hard to determine cause and effect, or the exact impact of social changes.

The study of communities has grown in popularity in part thanks to advances in social network analysis. From the earliest works by Zachary [27] to the more recent works of Leskovec et al. [16][17], 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 [27]. Zachary investigated the fission of a community, the process of communities splitting into two or more parts. He found that fission could be predicted by applying the Ford-Fulkerson min-cut algorithm [7] 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.”

Community splits in FOSS are referred to as forks, and are relatively common. Forking is defined 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.” Robles and Gonzalez-Barahona [23] 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. They classified these into six main categories. (Table 1 .) which we build on extensively. They identified a gap in the literature in case of “how the community moves when a fork occurs”.

TABLE 1: The main reasons for forking as classified by Robles and Gonzalez-Barahona [23]

Reason for forking	Example forks
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
Legal issues	X.Org & XFree

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 also 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.

The communication patterns of FOSS developers in a bug repository were examined by Howison et

Fig. 1: The main reasons for forking as classified by Robles and Gonzalez-Barahona [23]

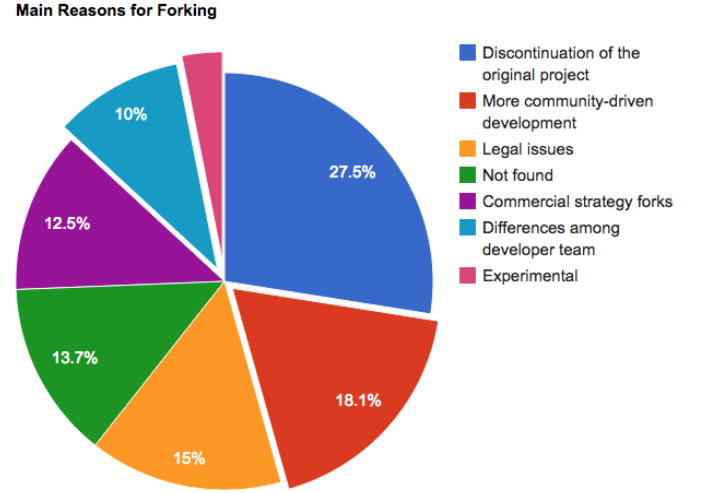


TABLE 2: The measures of diversity used by Asur et al. [1]

Metrics	Meaning
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)

al. [12]. They calculated out-degree centrality as their metric. Out-degree centrality measures the proportion of the number 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 approach is similar to theirs in how we form collaboration graphs and perform our temporal analysis. 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. [15]. They defined five network statistics 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. [15] argued that one

possible reason could be that the community structure consolidates as projects mature.

Community dynamics was the focus of a recent study by Hannemann and Klamma [10] 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 the survivors survived 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 maximal 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 approach builds on top of their findings, and the evolution of betweenness centralities and network diameters for the projects in our study are depicted in the following sections.

To date, most studies on FOSS have only been carried out on a small number of projects, and using snapshots in time. To our knowledge, no study has been done of project forking that has taken into account the temporal dimension.

### 3 MOTIVATION

To better understand and measure the evolution, social dynamics of FOSS projects, and integral components to understanding their evolution and direction, we need new and better tools. 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. We want to map the dynamics of communities to real world phenomena. Identification is the first step to rectify an undesired 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 case of FOSS, where near half the project contributors are volunteers [8]. Oh et al. [21] 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 unconstructive 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 OSS 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” [21].

Recipes for success or stagnation, sustainability or fragmentation could be identifiable, leading to a set of best practices and pitfalls.

## 4 RESEARCH GOALS AND METHODS

We propose that the social interactions data reflects the changes the community goes through, and is able to describe the context surrounding a forking event. Robles and Gonzalez-Barahona [23] classify the main reasons for forking into six classes, listed in Table 1.

Three of the six listed reasons for forking are socially related, and so should arguably be reflected in the social interaction data. As an example, if a fork occurs because of a desire for “more community-driven development”, we expect to see an interaction patterns in the collaboration data showing a strongly-connected core that is hard to penetrate for the rest of the community. In other words, in this case, the power stays in the hands of the same people throughout, as new people come and go.

We aim to analyze, quantify and visualize how the community is structured, how it evolves, and the degree to which community involvement changes over time. To this end, we adopted a methodology that we describe in the next section.

### RESEARCH OBJECTIVE 1: WHAT ARE THE SOCIAL PATTERNS ASSOCIATED WITH DIFFERENT TYPES OF FORKING?

**R.Q. 1.1** *Is there a prototypical fork for personal differences reason for forking?*

**R.Q. 1.2** *Is there a prototypical fork for more community-driven development reason for forking?*

**R.Q. 1.3** *Is a labeled project as a technical differences fork really only a technical differences fork?*

Are there patterns that exemplify these categories? What establishes an inflection point (fork)? Which metrics are indicative of inflection?

**R.Q. 1.4** *What are the determining factors?*

**R.Q. 1.5** *Where are the determining factors? In Mailing List data? In Code?*

**R.Q. 1.6** *What do the determining factor look like?*

TABLE 3: The measures of diversity used by Kunegis et al. [15]

Network property	Network is diverse when	Diversity Measures
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

## RESEARCH OBJECTIVE 2: MATCH BETWEEN OUR ANALYSIS AND THE REAL WORLD? OR, DOES OUR ANALYSIS REFLECT WHAT HAPPENED? IF YES, HOW WELL?

**R.Q. 2.1** *Does my analysis of the situation match what people in that community remember?*

**R.Q. 2.2** *If I show the analysis results to an open source developer, would they draw the same conclusion as our analysis results?*

## RESEARCH OBJECTIVE 3: CAN WE USE THIS METHOD TO KNOW *how* AND *why* PROJECTS FORK?

**R.Q. 3.1** *Current community folks, have them reflect on and future*

**R.Q. 3.2** *Can we use these to know how and why projects fork?*

## 5 METHODOLOGY

### 5.1 Phase 1: Data Collection

The study of forks by Robles and Gonzalez-Barahona [23] included information on 220 forks and their reasons. We applied three selection criteria to those projects. A project was short-listed if it was recent, i.e. the fork had happened after the year 2000; data was available; and they had a community of more than a handful of contributors. This three stage filtering process resulted in the projects listed in Table 6.

Data collection involved analyzing mailing list archives. We collected data for the year in which the fork happened, as well as for three month before and three months after that year in order to capture the social context context at the time of the fork.

### 5.2 Phase 2: Creating Communication Graphs

Many social structures can be represented as graphs. The nodes represent actors/players and the edges represent the interaction between them. Such graphs can be a snapshot of a network – a static graph – or a changing network, also called a dynamic graph. In this phase, we processed the data to form a communication

graph of the community. We were looking for how people interacted with each other. We decided to treat the general mailing list as a person, because the bulk of the communication was targeted at it, and most newcomers start by sending their questions to the general mailing list. Each communication effort was captured with a time-stamp. This allowed us to form a dynamic graph, in which the nodes would exist if and only if they had an interaction with another node during the period we were interested in.

### 5.3 Phase 3: Temporal Visualization and Temporal Evolution Analysis

In this phase, we wanted to analyze the changes that happen to the community over a given period of time, i.e. three months before and three months after the year in which the forking event happened. We measured betweenness centrality [5] of the most significant nodes in the graph, and the graph diameter over time. Figures [5][6][7] show the betweenness centralities over the 1.5 year period for the Kamailio, 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 of 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 had a relatively smooth transition.

There are many ways of looking at an individual's importance. 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 [6]. Secondly, if we are looking for people who can serve as bridges between 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 [6]. Third, if cross-community collaboration is the focus, we can measure *edge betweenness centrality*.

TABLE 4: Projects forked because of “personal differences among the developer team” [23] sorted in chronological order, and their data availability status

Original	Forked	Date	Data available?	Collected?
GNU Emacs	X Emacs	1991, ?	Only after 2000	N/A
NetBSD	OpenBSD	1995, Oct	Yes, but scarce	N/A
xMule	aMule	2003, Aug	Only 2006-2007	N/A
lMule	xMule	2003, Jun	-	N/A
Sodipodi	Inkscape	2003, Nov	Yes	Req.
Nucleus CMS	Blog:CMS	2004, May	Only after Sept 2004	N/A
BMP	Audacious	2005, Oct	Yes	Req.
ntfsprogs	NTFS-3G	2006, Jul	-	N/A
OpenWRT	FreeWRT	2006, May	Only after Oct 2006	N/A
QtiPlot	SciDavis	2007, Aug	-	N/A
<b>Kamailio</b>	OpenSIPS	2008, Aug	Yes	<b>Yes</b>
Blastwave.org	OpenCSW	2008, Aug	-	N/A
jMonkeyEngine	Ardor3D	2008, Sept	Yes, but scarce	N/A
Frog CMS	Wolf CMS	2009, Jul	-	N/A
Aldrin	Neil	2009, ?	-	N/A
<b>Ffmpeg</b>	libav	2011, Mar	Yes	<b>Yes</b>

TABLE 5: Projects forked because of the need for more community-driven development by Robles and Gonzalez-Barahona [23] sorted in chronological order

Original	Forked	Date	Data available?	Collected?
Nethack	Slash'EM	1996, ?	-	N/A
GCC	EGCS	1997, ?	-	N/A
SourceForge	Savane	2001, Oct	-	N/A
PHPNuke	PostNuke	2001, Sum	Not found	N/A
QTExtended	OPIE	2002, May	-	N/A
GraphicsMagick	Graphics	2002, Nov	Only after 2003	N/A
<b>freeglut</b>	OpenGLUT	2004, Mar	Yes	<b>Yes</b>
Mambo	Joomla!	2005, Aug	-	N/A
SER	Kamailio	2005, Jun	Only after 2006	N/A
PHPNuke	RavenNuke	2005, Nov	Not found	N/A
Hula	Bongo	2006, Dec	-	N/A
Compiere	A Dempiere	2006, Sept	No Dev. mailing list	N/A
Compiz	Beryl	2006, Sept	Only after Jun 2007	N/A
SQL-Ledger	LedgerSMB	2006, Sept	No Dev. mailing list	N/A
<b>Asterisk</b>	Callweaver	2007, Jun	Yes	<b>Yes</b>
CodeIgniter	KohanaPHP	2007, May	Not found	N/A
OpenOffice.org	Go-oo.org	2007, Oct	Only after Jun 2011	N/A
Mambo	MiaCMS	2008, May	-	N/A
TORCS	Speed Dreams	2008, Nov	Yes, but scarce	N/A
MySQL	MariaDB	2009, Jan	Yes	No
Nagios	Icinga	2009, May	Yes	Req
Project Darkstar	RedDwarf	2010, Feb	Yes, but scarce	N/A
SysCP	Froxlor	2010, Feb	-	N/A
Dokeos	Chamilo	2010, Jan	Not found	N/A
GNU Zebra	Quagga	2010, Jul	-	N/A
<b>rdesktop</b>	FreeRDP	2010, Mar	Yes	<b>Yes</b>
OpenOffice.org	LibreOffice	2010, Sept	Only after Jun 2011	N/A
Redmine	ChiliProject	2011, Feb	Yes, but scarce	N/A

TABLE 6: Forked projects for which collaboration data was collected

Projects	Reason for forking	Year
Kamailio & OpenSIPS	Differences among developer team	2008
ffmpeg & libav	Differences among developer team	2011
Asterisk & Callweaver	More community-driven development	2007
rdesktop & FreeRDP	More community-driven development	2010
freeglut & OpenGLUT	More community-driven development	2004
Amarok & Clementine Player	Technical (Addition of functionality)	2010
ApacheCouchDB & BigCouch	Technical (Addition of functionality)	2010
Pidgin & Carrier	Technical (Addition of functionality)	2008

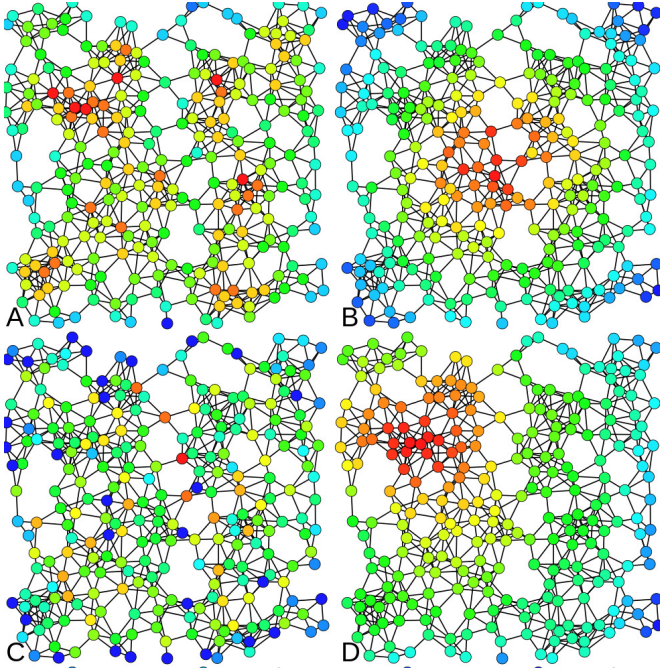


Fig. 2: Heat-map color-coded examples of nodes with high centrality metric are shown above. The same network is analysed four times with the following centrality measures: A) Degree centrality, B) Closeness centrality, C) Betweenness centrality and D) Eigenvector centrality [24]

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 [22] 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 [6].

#### 5.4 Phase 4: Temporal Visualization

Several visualization techniques and tools are used in the field of social network analysis, for instance, Gephi [2], which is a FLOSS tool for exploring and manipulating networks. It is capable of handling large networks with more than 20,000 nodes and features several SNA algorithms. It is customizable with plugins and we used it for dynamic network visualization. We visualized the dynamic network changes using Gephi [2]. The videos 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 each other proportional to the graph distance between them. This results in a graph shape between Fruchterman & Rheingold’s [9] layout and Noack’s LinLog [18].

## 6 RESULTS AND DISCUSSION

### 6.1 Kamailio Project

Figure 4 shows four key frames from the Kamailio 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 Kamailio 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.

Figure 5 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 axis shows the percentage of the top betweenness centralities for each node. The saliency of the GeneralList – colored as light blue – is apparent due to its continuous and dominant presence in the stacked area chart. The chart legend



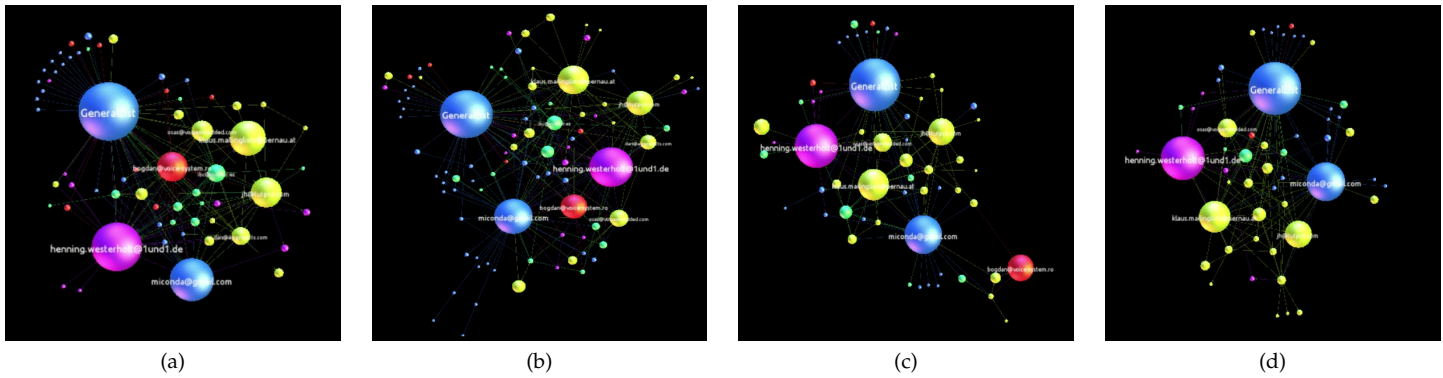


Fig. 4: Snapshots from video visualization of Kamilio's graph (Oct. 2007 - Mar. 2009) in which a core contributor (colored red) moves to the periphery and eventually departs the community.

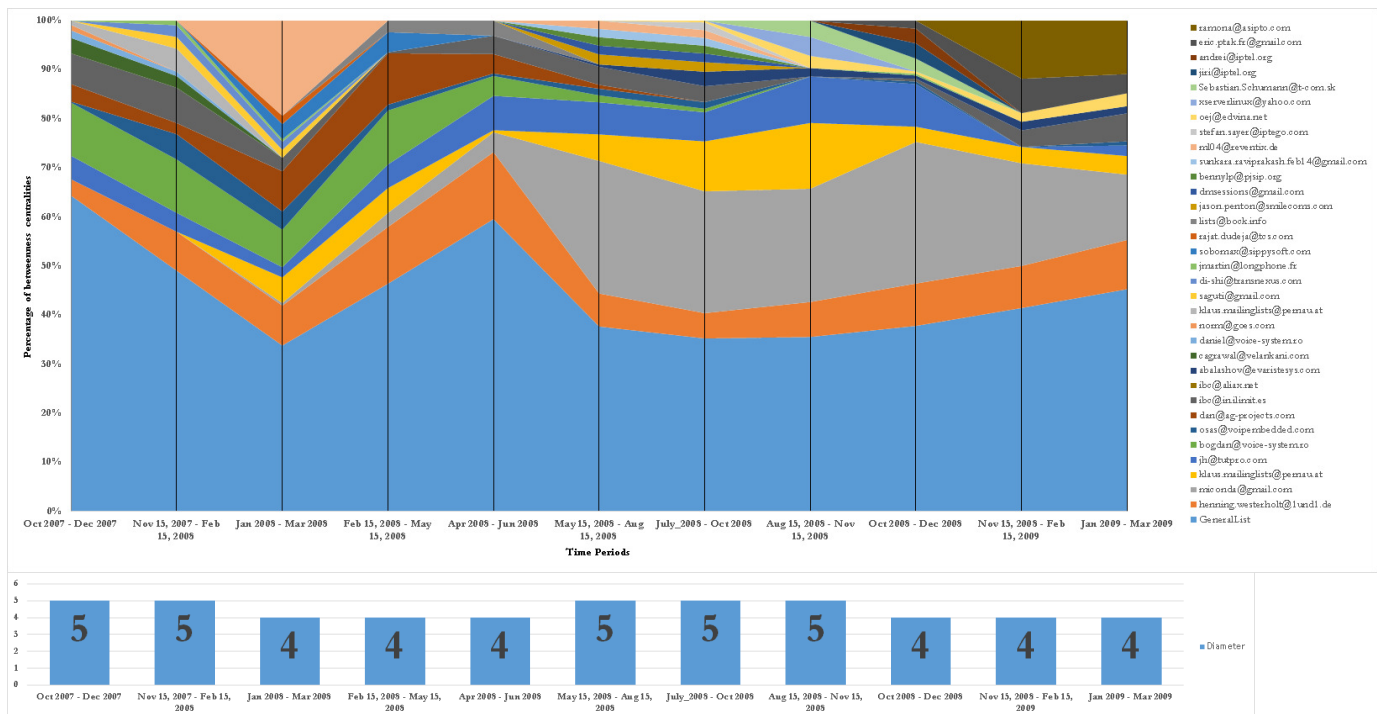


Fig. 5: Kamilio top contributors' betweenness centralities and network diameter over time (Oct. 2007 to Mar. 2009) in 3-month time windows with 1.5-month overlaps

lists the contributors based on the color and in the same order of appearance on the chart starting from the bottom. One can easily see that 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 helps identify 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 5. The increase in the network diameter during this period confirms the findings of Hannemann and Klamka [10].

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.

## 6.2 Amarok Project

The video for the Amarok project fork is available online<sup>1</sup>, and the results from our quantitative analysis of the betweenness centralities and the network diameters are shown in Figure 6. 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 healthy network, rather

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

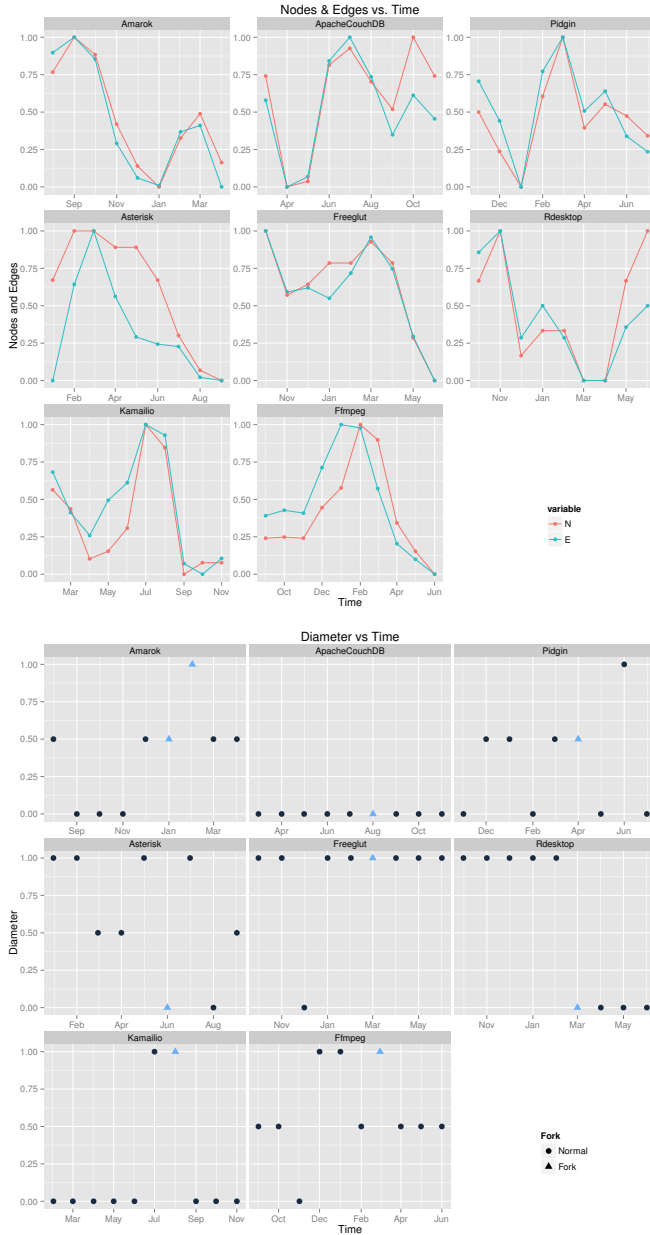


Fig. 3: Nodes and Edges Over Time

than of simmering conflict. These indicators show that Amarok fork in 2010 arguably belongs to the “addition of technical functionality” rationale for forking, as there are no visible social conflict.

### 6.3 Asterisk Project

The video for the Asterisk project is also available online, and the results from our quantitative analysis of the betweenness centralities and the network diameters are shown in Figure 7. 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 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.

## 7 TIMELINE

I want to graduate when XXX

## 8 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 [23] and Hannemann and Klamma [10], and identified three projects that had forked in the recent past. We mined the collaboration data, formed dynamic collaboration graphs, and measured social network analysis metrics over an 18-month period time window.

We also visualized the dynamic graph (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 personal differences representative project was identifiable, and so was the date it forked, with a month accuracy. The novelty of the approach was in applying the 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.

## 9 THREATS TO VALIDITY

The presented findings may not be generalized to all OSS projects. The projects studies in this paper were selected from a pool of candidate projects, partly because data about them was available. Given access, a better sampling approach has to be adopted, which could result in a more robust investigation. Furthermore, the proposed technique uses the data from online communications. The assumption that all the communication can be captured by mining repositories is intuitively imperfect, but inevitable. Hence, to minimize the effect of this assumption, we plan to complement the quantitative approach with a qualitative approach of interviewing key individuals from the community as future work.



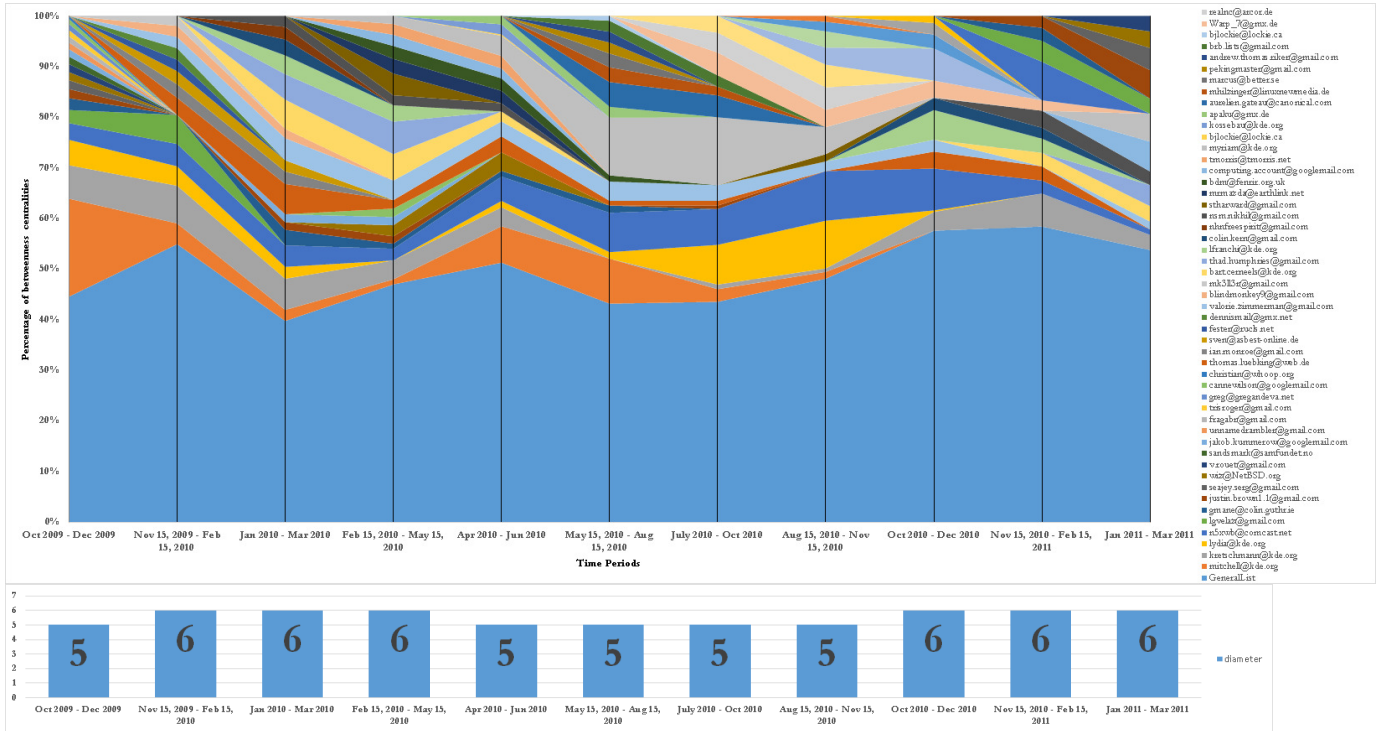


Fig. 6: Amarok project's top contributors' betweenness centralities and network diameter over time between Oct. 2009 to Mar. 2011 in 3-months time windows with 1.5 months overlaps

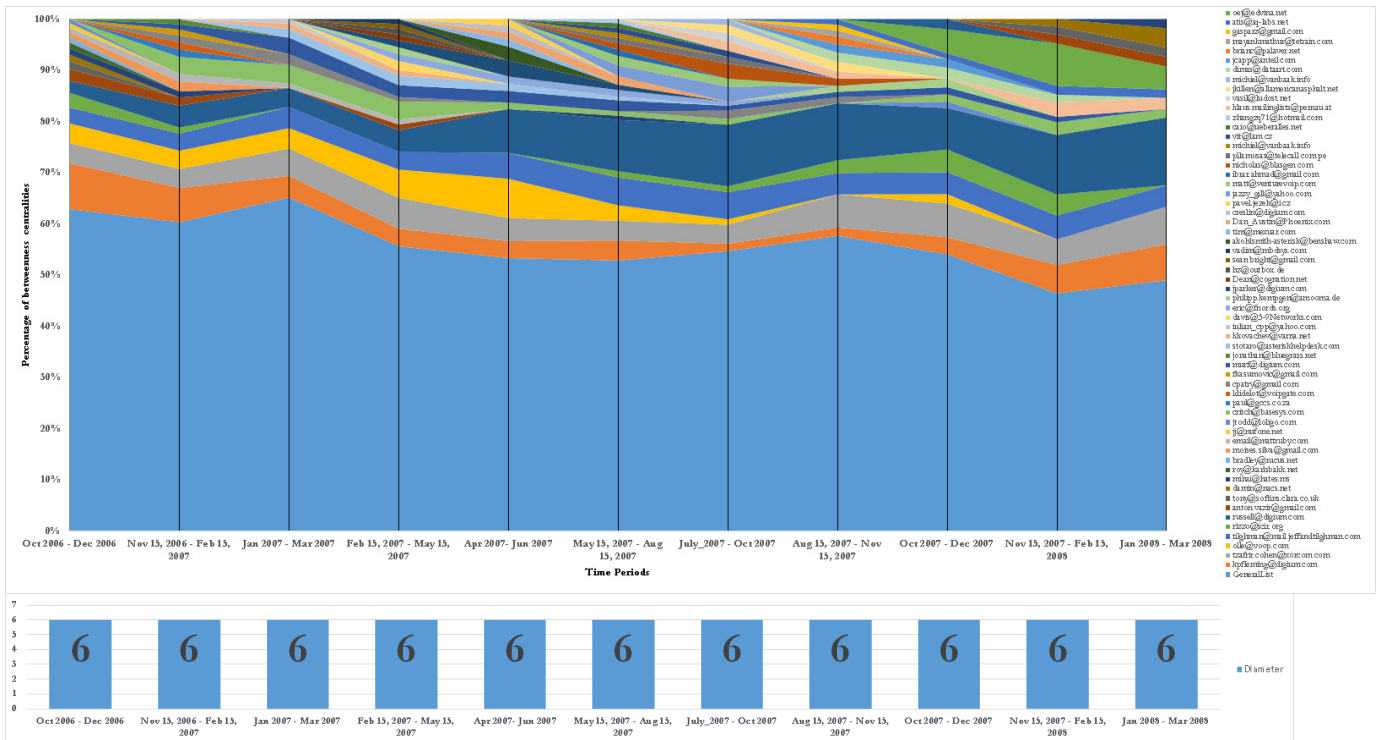


Fig. 7: Asterisk project's top contributors' betweenness centralities and network diameter over time between Oct. 2006 to Mar. 2008 in 3-months time windows with 1.5 months overlaps

## ACKNOWLEDGEMENT

The author would like to thank his academic adviser, Prof. Carlos Jensen, and his committee members, Profs. Margaret Burnett, Ronald Metoyer, and Christopher Scaffidi for their insightful guidance throughout the author's PhD program. The author 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.

## REFERENCES

- [1] Asur, S., S. Parthasarathy, and D. Ucar, (2009), "An event-based framework for characterizing the evolutionary behavior of interaction graphs," in ACM Trans. Knowledge Discovery Data. 3, 4, Article 16, (November 2009), 36 pages. 2009.
- [2] Bastian, M., S. Heymann, and M. Jacomy, "Gephi: an open source software for exploring and manipulating networks," presented at the Int. AAAI Conf. on Weblogs and Social Media, 2009.
- [3] Bezrukova, K., C. S. Spell, J. L. Perry, (2010). , "Violent Splits Or Healthy Divides? Coping With Injustice Through Faultlines," Personnel Psychology, Vol 63, Issue 3. 2010.
- [4] Bird, C., D. Pattison, R. D'Souza, V. Filkov, and P. Devanbu, "Latent social structure in open source projects," in Proc. of the 16th ACM SIGSOFT Int. Symposium on Foundations of software engineering, New York, NY, USA: ACM, pp. 24-35, 2008.
- [5] Brandes, U. (2001). "A Faster Algorithm for Betweenness Centrality", in Journal of Mathematical Sociology 25(2):163-177.
- [6] Chakrabarti, D., and C. Faloutsos. "Graph mining: Laws, generators, and algorithms," ACM Computing Surveys, 38, 1, Article 2, 2006.
- [7] 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.
- [8] Forrest, D., C. Jensen, N. Mohan, and J. Davidson, "Exploring the Role of Outside Organizations in Free/ Open Source Software Projects," in Proceedings of the 8th IFIP WG 2.13 International Conference on Open Source Systems, OSS 2012, Hammamet, Tunisia, September 10-13, 2012.
- [9] Fruchterman, T. M. J. and E. M. Reingold, (Nov. 1991). "Graph drawing by force-directed placement," Softw: Pract. Exper., vol. 21, no. 11, pp. 1129-1164.
- [10] Hannemann, A and , R. Klammer "Community Dynamics in Open Source Software Projects: Aging and Social Reshaping," in Proc. of the IFIP Second Int. Conf. on Open Source Systems, pp. 80-96, 2013.
- [11] Howison, J. and K. Crowston. "The perils and pitfalls of mining SourceForge," In Proceedings of the Int. Workshop on Mining Software Repositories (MSR 2004), pp. 7-11. 2004.
- [12] Howison, J., K. Inoue, and K. Crowston, "Social dynamics of free and open source team communications," in Proc. of the IFIP Second Int. Conf. on Open Source Systems, 319-330, 2006.
- [13] Howison, J., M. Conklin, and K. Crowston, "FLOSSmole: A collaborative repository for FLOSS research data and analyses," Int. Journal of Information Technology and Web Engineering, 1(3), 1726. 2006.
- [14] 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.
- [15] Kunegis, J., S. Sizov, F. Schwaigereit, and D. Fay, "Diversity dynamics in online networks," in Proc. of the 23rd ACM Conf. on Hypertext and Social Media, USA, 2012.
- [16] Leskovec, J., Kleinberg, J., and Faloutsos, C.: "Graphs over time: densification laws, shrinking diameters and possible explanations," in Proc. of the SIGKDD Int. Conf. on Knowledge Discovery and data Mining, 2005.
- [17] Leskovec, J., K. J. Lang, A. Dasgupta, and M. W. Mahoney, "Statistical properties of community structure in large social and information networks," in Proc. of the 17th Int. Conf. on World Wide Web (WWW '08), ACM, 2008.
- [18] Noack, A., (2007) "Energy models for graph clustering," J. Graph Algorithms Appl., vol. 11, no. 2, pp. 453-480.
- [19] Nyman, L. , "Understanding code forking in open source software," in Proc. of the 7th Int. Conf. on Open Source Systems Doctoral Consortium, Salvador, Brazil, 2011.
- [20] Nyman, L., T. Mikkonen, J. Lindman, and M. Fougère, "Forking: the invisible hand of sustainability in open source software," in Proc. of SOS 2011: Towards Sustainable Open Source, 2011.
- [21] Oh, W., Jeon, S., "Membership Dynamics and Network Stability in the Open-Source Community: The Ising Perspective" in Proceedings of the Twenty-Fifth International Conference on Information Systems. 2004.
- [22] Page, B., B. Sergey, R. Motwani and T. Winograd, "The PageRank Citation Ranking: Bringing Order to the Web," Technical Report, Stanford InfoLab, 1999.
- [23] Robles, G. and J. M. Gonzalez-Barahona, "A comprehensive study of software forks: Dates, reasons and outcomes," in Proc. of the 8th Int. Conf. on Open Source Systems, Hammamet, Tunisia, 2012.
- [24] Rocchini, C. (Nov. 27 2012), Wikimedia Commons, Available: <http://en.wikipedia.org/wiki/File:Centrality.svg>, 2012.
- [25] 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.
- [26] Torres, M. R. M., S. L. Toral, M. Perales, and F. Barrero, "Analysis of the Core Team Role in Open Source Communities," in Complex, Intelligent and Software Intensive Systems (CISIS), 2011 Int. Conf. on, pp. 109-114. IEEE, 2011.
- [27] 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.