

AN ABSTRACT OF THE THESIS OF

Rithika Kiran Naik for the degree of Master of Science in Computer Science
presented on September 15, 2015.

Title: Visualizing Contribution Patterns in Open Source

Abstract approved: _____

Dr Carlos Jensen

Open Source software gives users the freedom to copy, modify and redistribute source code without having to pay any amount of money to an organization or individual. The evolution of these softwares usually depend a lot on how the participating crowd interact and co-operate with each other. Over the past few years, open source software have become widely accepted and used hence making their study a very hot topic. To understand the working philosophy of any open source software one needs to look into a variety of factors like community interaction, strength, evolution, support and contribution patterns. Bitergia is one such organization which is highly involved in getting, providing and analyzing Software Development metrics and information for open source projects. The main idea behind this thesis work is to understand contribution patterns by the community in few open source projects by the means of various visualizations. Data visualization has gained a lot of popularity as it helps in easy analysis of huge amounts

of information quickly. We have hence taken up the task of improving the existing visualizations by Bitergia into visualizations which will help in taking informed decisions clearly. Our analysis were performed on five different datasets : Eclipse Foundation, OpenStack Foundation, Puppet Labs, Red Hat and Apache Cloud-stack.

©Copyright by Rithika Kiran Naik
September 15, 2015
All Rights Reserved

Visualizing Contribution Patterns in Open Source

by

Rithika Kiran Naik

A THESIS

submitted to

Oregon State University

in partial fulfillment of
the requirements for the
degree of

Master of Science

Presented September 15, 2015
Commencement June 2015

Master of Science thesis of Rithika Kiran Naik presented on September 15, 2015.

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

Rithika Kiran Naik, Author

ACKNOWLEDGEMENTS

Thank you mummy daddy :D

TABLE OF CONTENTS

	<u>Page</u>
1 Introduction	1
1.1 Introduction to the Introduction	1
2 Bitergia	2
2.1 Services provided	2
2.2 Task at hand	4
3 Related Work	8
3.1 Introduction to the Introduction	8
4 Methodology	9
4.1 Introduction to the Introduction	9
5 User study	10
5.1 Introduction to the Introduction	10
6 Results	11
6.1 Introduction to the Introduction	11
7 Conclusion	12
7.1 Fin	12
Bibliography	12

LIST OF FIGURES

<u>Figure</u>	<u>Page</u>
1.1 Go figure.	1
2.1 Bitergia data forums	3
2.2 Timezone origin of messages	5
2.3 Timezone origin of messages for 2002,2007,2012	5
2.4 Comparison for pyramids every two years	6
2.5 3D version of pyramids every two years	7
3.1 Go figure.	8
4.1 Go figure.	9
5.1 Go figure.	10
6.1 Go figure.	11

Chapter 1: Introduction

I have done some excellent research [?].

1.1 Introduction to the Introduction

A rectangular box with a thin black border, containing the word "Box" in a serif font.

Figure 1.1: Go figure.

Chapter 2: Bitergia

Bitergia was founded by four professors from the Universidad Rey Juan Carlos, Spain in July 2012. They were researchers who were very much involved in understanding the open source community and its growth for long. This was the main idea behind starting Bitergia, the software development analytics company. The focus here was to have data analytics techniques to mine information about project performance, track developer actions, find areas of improvement and identify risks involved in further development.

2.1 Services provided

As providing software metrics was the goal here, they started off by building tools that would help them achieve the same by extracting data from multiple forums. Few such tools developed are -

1. CVSAnalY : This tool collects and organizes data from various code management systems like CVS, Subversion and Git.

2. RepositoryHandler : This is a python library which handles code repositories using a common interface.

3. MailingListStats : This is a command line based tool which retrieves information from mailing lists. It supports files in mbox and web accessible formats.

4.Bicho : Another command line tool that retrieves information from issue tracking systems like Bugzilla, Jira, Sourceforge, Allura, Github, Google Code, Launchpad etc.

5.CMetrics : This is a library used to measure size and complexity for C files in projects.

6.Sibyl : This tool is used to extract information from websites that have a Question-and-Answer format like Askbot with further support planned for Stack-Overflow.

These bunch of tools together are known as the *MetricsGrimoire Library* which is accesible for all in Github.

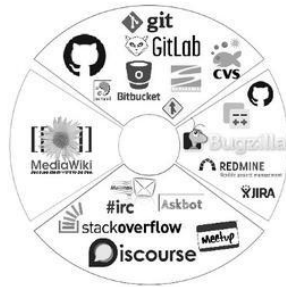


Figure 2.1: Bitergia data forums

Now that they had the tools ready to extract the data, the next step was to build a framework which would help them to analyze and then visualize the data which they would obtain from the MetricsGrimoire tools. This is known as the *VizGrimoire Library* and has two important tools in it -

1.vizGrimoireR : After extracting the needed data, this tool helps in analyzing the dataset to produce data which then can be used to answer specific type of

questions.

2.vizGrimoireJS : Having the right kind of data in hand led to the need to have a dashboard kind of environment to showcase the data in a meaningful way. This tool hence creates dashboards with reports and visualizations of the data.

2.2 Task at hand

We were presented with two cases which needed improvised visualizations. The existing visualizations were not feasible to answer different kinds of questions that one can pose. We had to work through various visualizations in order to finalize on a nearly good enough solution. We will first talk about each of the two cases presented to us and then discuss the solutions that we fixed upon in the next chapter.

CASE 1 -Where do Developers work ?

Location has always been of importance in any area of work. In the software industry, distance always was not a matter of concern because of the nature of work involved, with this being even more true for Open Source Software as its developers are assumed to be contributing according to their convenience [1]. Considering this, the group at Bitergia decided to investigate the overall spread of developers across the world for a few projects. Their main database was the data they collected from github and mailing lists for those organizations. There were certain assumptions involved like the timezones mentioned in the mail tools were correct and that they correspond to the geographical areas to some extent. They came up with these

two basic visualizations as a start to analyzing the data.

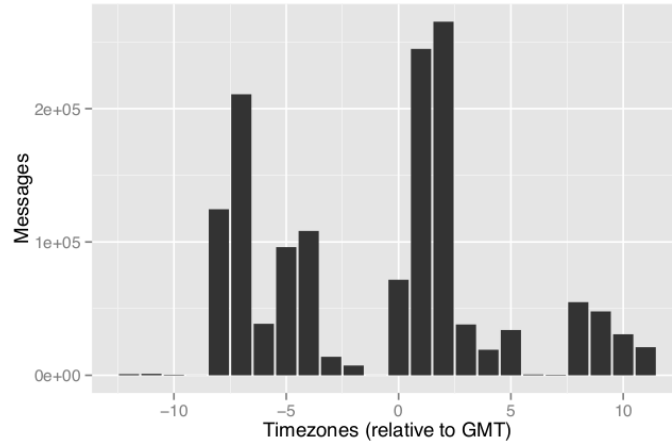


Figure 2.2: Timezone origin of messages

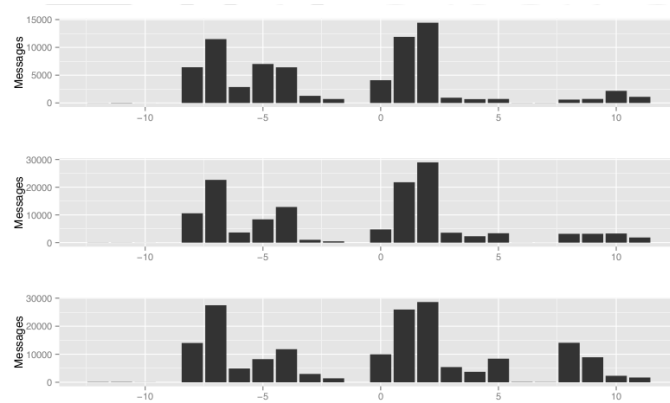


Figure 2.3: Timezone origin of messages for 2002,2007,2012

The following was the division of geography with respect to timezone -

1. America: GMT-8 to GMT-2 (US/Canada: -8 to -4)
2. Europe/Africa/Middle East: GMT to GMT+5
3. East Asia/Australia: GMT+8 to GMT+11

The visualizations being very basic didn't help solve or understand a lot about the developer spread for an organization. They were more oriented in giving an overview of messages per timezone. It was not clear about the relation between commits and authors and location tied together which needed more than a birds' eye view type of a visualization.

CASE 2 - Experience of Developers A lot has been discussed about the motivation for developers to join or stay or leave an open source project. The most sort after work being the onion model which describes the general roles that developers take and how they gravitate eventually [2, 3]. The Bitergia team focused on calculating age in project for active developers at a certain time. They again had access to datasets which gave them the data needed to calculate the same. As a result of their work the following visualizations were built -

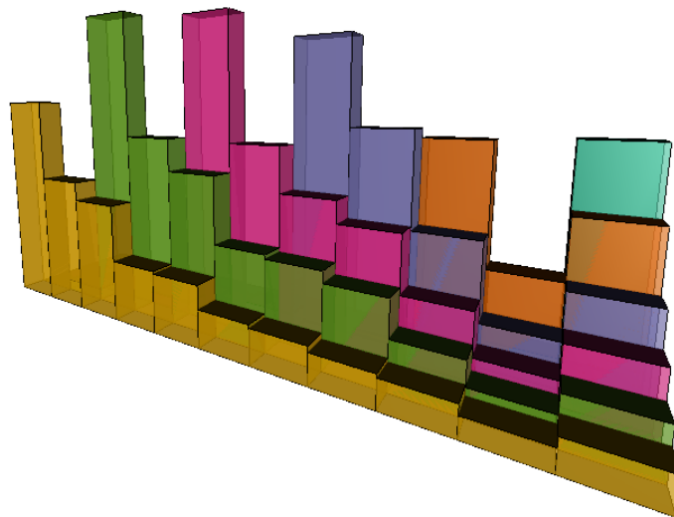


Figure 2.4: 3D version of pyramids every two years

The problem here was that, the 3D visualization, though colorful, was a little

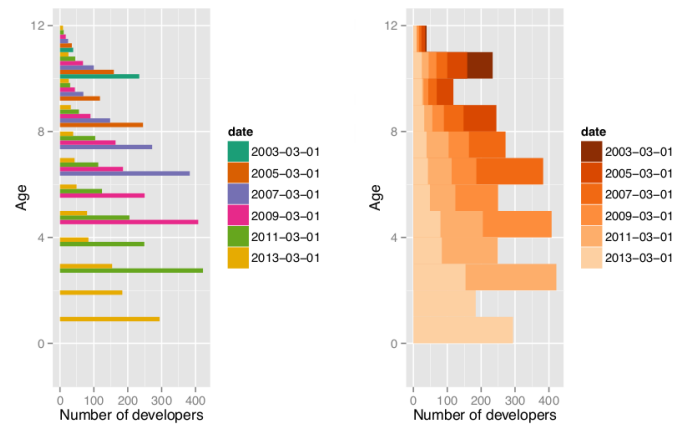


Figure 2.5: Comparison for pyramids every two years

hard to relate to understand the dynamics. WHAT ELSE IS WRONG ?????? My visuals are the same!

Chapter 3: Related Work

I have done some excellent research [?].

3.1 Introduction to the Introduction

A rectangular box with a thin black border, containing the word "Box" in a serif font.

Figure 3.1: Go figure.

Chapter 4: Methodology

I have done some excellent research [?].

4.1 Introduction to the Introduction

Box

Figure 4.1: Go figure.

Chapter 5: User study

I have done some excellent research [?].

5.1 Introduction to the Introduction



Box

Figure 5.1: Go figure.

Chapter 6: Results

I have done some excellent research [?].

6.1 Introduction to the Introduction

A rectangular box with a thin black border containing the word "Box" in a serif font.

Figure 6.1: Go figure.

Chapter 7: Conclusion

Wow, that really was excellent.

7.1 Fin

This is the end, my only friend, the end.

Bibliography

- [1] Takhteyev Y , Hilts A. Investigating the Geography of Open SourceSoftware through Github, 2010.
- [2] Crowston K , Howison J. The Social Structure of Free and Open Source Software Development. In *First Monday*, 2005.
- [3] Ye Y , Kishida K. Towards an Understanding of the motivation of Open Source Developers. In *ICSE*,2003.

