

# Treehopper – exploring version controlled software code bases using graph databases

Pradeep Gowda and Mehmet Kilicarslan

December 9, 2013

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Approach</b>	<b>2</b>
2.1	Git version control system . . . . .	3
2.2	Design considerations . . . . .	5
2.3	Technology choice . . . . .	5
2.3.1	Database . . . . .	5
2.3.2	Backend . . . . .	5
2.3.3	Front end . . . . .	5
2.4	Usage . . . . .	5
2.4.1	Installing Neo4j database . . . . .	5
2.4.2	Installing the treehopper application . . . . .	5
2.4.3	Loading repository data . . . . .	6
2.4.4	Visualising graph nodes . . . . .	6
2.4.5	Analytical dashboard . . . . .	6
2.4.6	Repository view . . . . .	6
<b>3</b>	<b>Results</b>	<b>8</b>
<b>4</b>	<b>Conclusion</b>	<b>8</b>

<b>5</b>	<b>Future work</b>	<b>8</b>
<b>6</b>	<b>Reference</b>	<b>8</b>

## 1 Introduction

A typical software project has hundreds of files, developed over months and years by numerous developers. Version control systems are an integral part of a software development practice. Version control systems not just important for maintaining the history of a project, they are also the foundation for a team to collaborate.

Version controlled code bases contain more than just the history of individual files, they are also a important artefacts in the archeology of software development.

There are many ways to visualise source code. Module dependency graphs are one of them. Dependency graphs for Object oriented programming languages have classes as the nodes and edges show the dependency between the class and where it is being used. However, there is lot of understanding captured outside the source code.

We wanted to explore the codebases interactively and answer some interesting questions like:

- Who has worked on this project for the longest time?
- What is the activity level on this project? Has there been an uptick in code commits recently?
- What is the “bus factor” on this project? That is, if one or more developers leave the team, what will be the impact?
- What is the nature of source code? Eg: what percentage is C files, HTML files etc.,

We decided on developing a software that let us find answers to these kind of questions on a “on-demand” basis and present it as a visual dashboard.

## 2 Approach

We selected `git` distributed version control system as the basis for analyzing codebases. Git was developed by Linus Torvals (the developer of Linux Operating System) as an answer to the problem of having to coordinate the work of developers worldwide that goes into to continued development of Linux Kernel. Git has become the most popular of the open source distributed version control systems among it’s peers, such as `mercurial`, `bazaar`, `darcs`, `fossil` etc., A

large number of open source projects now use git for distributed version control. Enterprises often have their own setups of git that compliments their development practices.

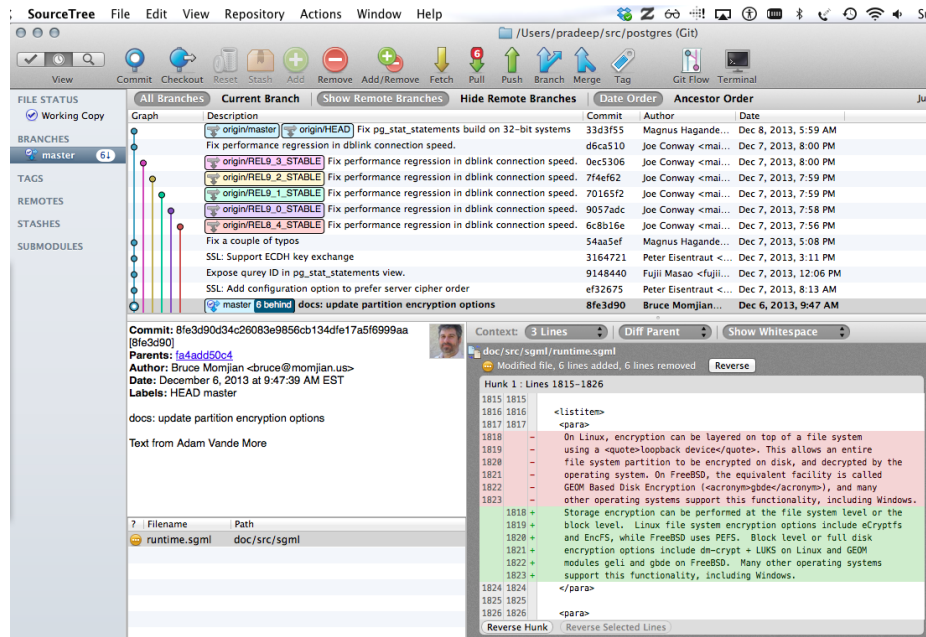


Figure 1: A GUI view of a git repository

The above image shows a snapshot of the postgresql database's source code.

The various coloured lines represent the various branches

The description corresponds to a commit, which in turn is a SHA1 hash guaranteed to be unique. Each commit has a committer. Sometimes the author of a change to the files is different than the person who commits it to the repository.

The bottom two windows show the

## 2.1 Git version control system

Every git directory maintains the complete history of changes made to the files. Git stores these changes in an internal representation called the git object storage. This storage is a directed acyclic graph.

Files in a code repository are represented by **blob** (though blobs can point to other things like symbolic links).

Directories are represented by **trees**. The trees refer to **blobs**.

A commit refers to a **tree** that represents the state of the files at the time of commit.

**refs**: References/heads/branches are bookmarks that point to a node in the DAG. They serve as reminders to the developers as to where they are working at the moment. The **HEAD** ref is a special ref that points to the currently active branch.

The following graph shows the relation between **blob**, **tree**, **refs** etc.,

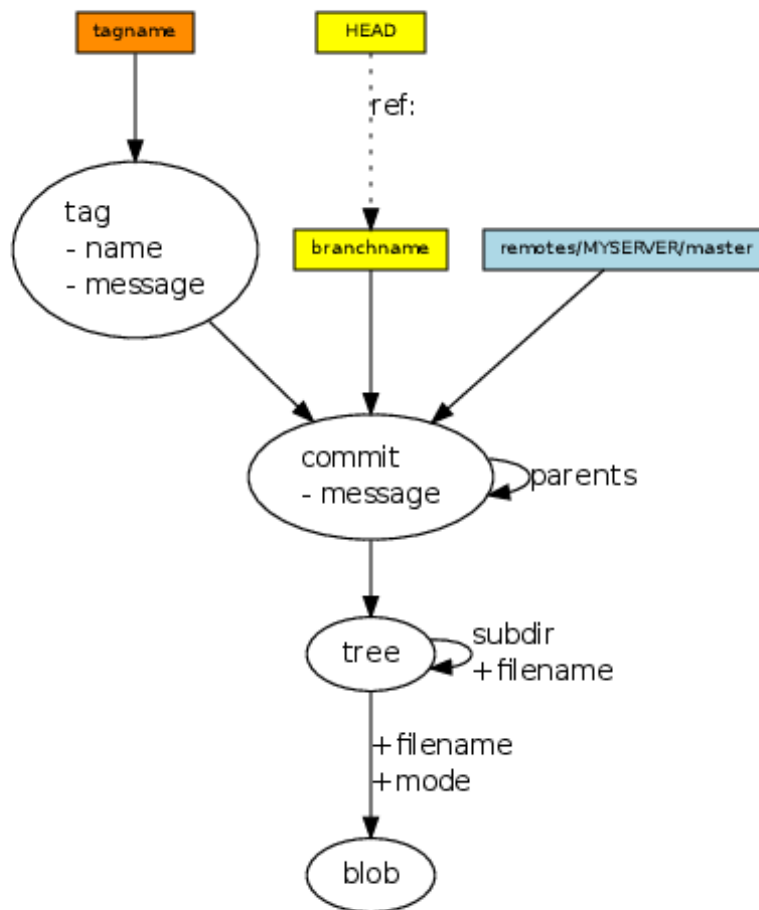


Figure 2: Git as a graph

## 2.2 Design considerations

## 2.3 Technology choice

### 2.3.1 Database

#### 2.3.1.1 Neo4j

#### 2.3.1.2 Cypher Query Language

### 2.3.2 Backend

#### 2.3.2.1 Python Git-Python

#### 2.3.2.2 Django

### 2.3.3 Front end

Zurb Foundation

## 2.4 Usage

### 2.4.1 Installing Neo4j database

Download the Neo4j database from the website – [<http://www.neo4j.org>] and unzip (into, say \$NEO4JPATH) and start the server using the command line interface

```
$ cd NEO4JPATH
$ bin/neo4j start
```

You can open <http://localhost:7474/> in the browser to see the web interface of the database server.

### 2.4.2 Installing the treehopper application

There are two parts to the application.

- Data loader – a command line interface
- Dashboard – a web interface

### 2.4.3 Loading repository data

Using the command line interface, the user can parse the git repository and upload the commit, user, tag, and file information to the graph database.

The CLI invocation is:

```
python manage.py load_git --url /Users/pradeep/src/requests --name requests
```

#### 2.4.4 Visualising graph nodes

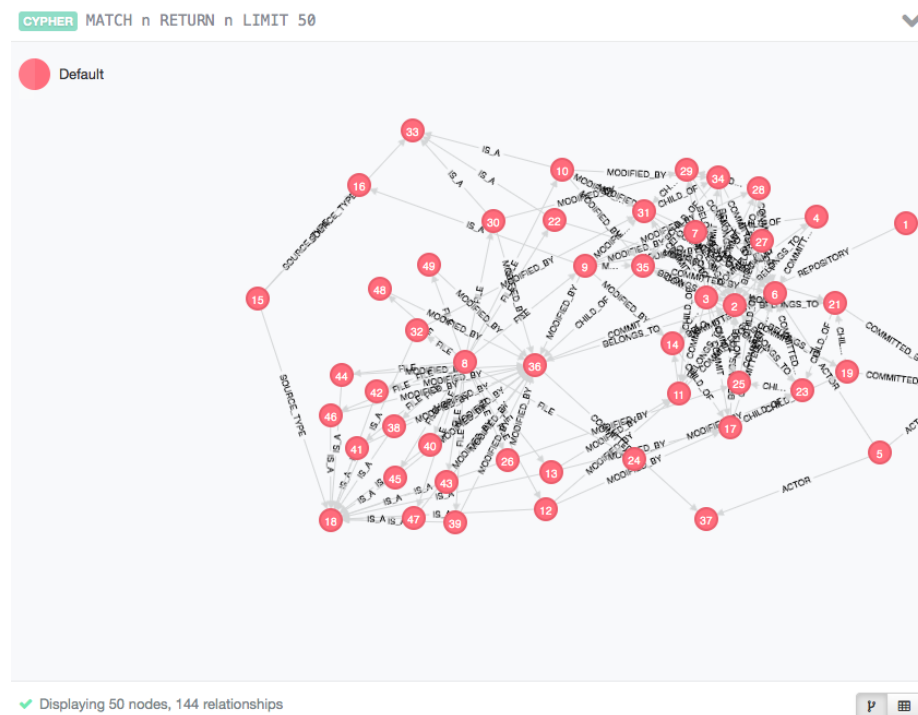


Figure 3: Query Interface

### 2.4.5 Analytical dashboard

Front page of the applications where we can see all the repositories known to the application

### 2.4.6 Repository view

Each repository known to the Application shows a dashboard like this:



# Treehopper

Code analytics dashboard for your team

## Repositories

- [postgres](#) at [/Users/pradeep/src/postgres](#)
- [requests](#) at [/Users/pradeep/src/requests](#)
- [hakyll](#) at [/Users/pradeep/src/hakyll](#)

Figure 4: Front page



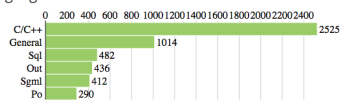
# Treehopper

Code analytics dashboard for your team

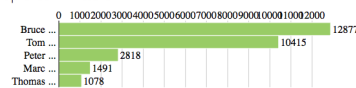
Repository: **postgres**

There are **35846** commits by **39** developers since 1996-07-09.

### Language Statistics

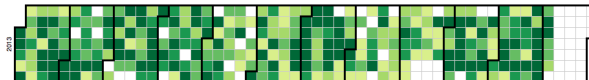


### Top 5 committers



### Calendar View

Heat map of the commits in the calendar year



### Commit History

Summary	Time	Author
docs: update partition encryption options	2013-12-06	Bruce Mornjian
docs: clarify SSL certificate authority chain docs	2013-12-06	Bruce Mornjian
Fix improper abort during update chain locking	2013-12-05	Alvaro Herrera
Clear retry flags properly in replacement OpenSSL sock_write function.	2013-12-05	Tom Lane
Avoid resetting Xmax when it's a multi with an aborted update	2013-12-05	Alvaro Herrera
build: pass EXTRA_REGRESS_OPTS to secondary regression tests	2013-12-04	Bruce Mornjian
doc: split long query into multiple lines	2013-12-04	Bruce Mornjian
Fix whitespace	2013-12-04	Peter Eisentraut
Don't include unused space in LOG_NEWPAGE records.	2013-12-03	Heikki Linnakangas
Fix full-page writes of internal GIN pages.	2013-12-03	Heikki Linnakangas

Marc G. Fournier has worked the longest on this project (since 1996-07-09).

### Release History

Release	Date	Released by
REL9_3_BETA1	2013-05-06	Tom Lane
REL9_2_BETA2	2012-05-31	Tom Lane
REL9_2_BETA1	2012-05-10	Tom Lane
REL9_1_BETA2	2011-06-09	Tom Lane
REL9_1_BETA1	2011-04-27	Tom Lane
REL9_1_ALPHA5	2011-03-28	Robert Haas
REL9_1_ALPHA4	2011-03-09	Bruce Mornjian
REL9_1_ALPHA3	2010-12-28	Peter Eisentraut
REL9_1_ALPHA2	2010-10-31	Tom Lane
REL9_1_ALPHA1	2010-09-03	Tom Lane

Figure 5: Repository view

### **3 Results**

### **4 Conclusion**

### **5 Future work**

- Support other distributed version control systems like `mercurial`.
- Handle more than one branch

### **6 Reference**

- [Git for Computer Scientists](#)