

Software Design



Git OSS-um

Gary Baker
Van Steinbrenner
Stephen White

February 5, 2019

Project Sponsor: Dr Igor Steinmacher,
Assistant professor, SICCS, Northern Arizona University

Faculty Mentor: Ana Paula Chaves Steinmacher
Version 1.0

Table of Contents

Introduction	2
Implementation Overview	3
Architectural Overview	5
Module and Interface Descriptions	7
Mining Script	7
MongoDB	10
Front-end Components with Django	12
Logging Into The Application	13
Admin Portal	13
About Us Page	13
Requesting Data	14
Data Visualization	14
Cross Comparison of Repositories	15
Python-nvd3 Framework	15
Implementation Plan	16
Conclusion	18

1. Introduction

Open source software (OSS) is a type of computer software in which source code is released under a license, where the copyright holder grants users the rights to study, change, and distribute the software to anyone and for any purpose. OSS may be developed in a collaborative public manner using a version control system named Git and hosting the project itself on a site known as GitHub¹. This version control system tracks changes made to a project and easily allows projects to have multiple contributors. It also allows for a project owner or manager to accept or decline contributions depending on their opinion of the quality of the change through what is called a pull request.

Unfortunately, a substantial amount of developers do not contribute to opensource projects. We have all been there: we do not know where to start, what files are important, or what needs to be changed. The community may be friendly and quick to provide assistance to a newcomer, or the contributor may get berated for submitting a pull request that does not adhere to the project's standards. It is so daunting that the newcomer gives up and moves on. This experience shows that there is a need in the open source community for a tool that will allow a novice passing through to quickly and efficiently determine the newcomer "friendliness" of a repository.

Our client, Dr. Igor Steinmacher, is a researcher and assistant professor at Northern Arizona University. Dr. Steinmacher and his colleagues have created a web application called FLOSScoach² to address newcomers' needs, where a user can investigate the skills needed to contribute to open source. Although FLOSScoach has been helpful, it has its shortcomings, such as:

¹ Add the website here

² FLOSScoach: <http://www.flosscoach.com>

- Data is manually fed.
- Contant manual maintenance.
- Data is not up to date.

Due to these shortcomings, we will be spearheading a solution, independent from FLOSScoach to give power back to the newcomers. Our goal is to create an autonomous process where relevant data of a requested repository is presented to the user. This will allow them to decide if the specific OSS project will support their needs and help them feel at ease when picking their first open source project.

In order to develop a successful newcomer meter for open source dynamics, we will create a web application that takes a URL leading to a GitHub project page and accesses it via the GitHub API. The collected data, such as newcomer acceptance rates and newcomer retention, among various other data points, will be presented to the newcomer as interactive graphics, which will ultimately assist users in finding the best OSS projects for them.

Within this document, we will break down our envisioned implementation of features that meet the requirements associated with our client's vision. We will be discussing our implementation and architectural overviews, modules, interfaces, and finally end with a brief recap of our design. In the following section, we will describe the implementation overview in order to give the "big picture" of our product and the technologies that we plan on using to make this product come to life.

2. Implementation Overview

As stated before, we are creating a web application address newcomers' needs to avoid being discouraged to start contributing due to their lack of knowledge on how, as well as the inherent intimidation that comes along with contributing. Our application will retrieve GitHub data

for a specified repository and display it to the user in a tabular form, accompanied by around 3-4 interactive graphics that will help the user visualize the data.

The way that this application will work is mostly simple. The user will first request information for a specific repository on our application front-end. The request is sent to the application back-end where the data is retrieved from our database and sent back to the front-end. Then, at the front-end, raw data points are formatted into a table as well as the utility Plot.ly takes that data and translates it into interactive graphics that are all displayed to the newcomer to help them form their own opinion about whether or not the repository is right for their first contribution or not.

Our application will include a few different technologies, but it mainly revolves around the web framework Django. Django is a high-level Python web framework that consists of Python, HTML, CSS, and JavaScript to create both great front-end user interfaces as well as fast and reliable back-end functionality. We have chosen for our back-end logic as well as our front-end user experience to be done using Django due to our familiarity with Python in general and the existing support and libraries that can be utilized when using Python and Django. Django is a very powerful framework and will allow us to do everything that is needed for this product while doing so with high efficiency and a streamlined development experience.

We will be storing a lot of data for our application in order to provide the newcomer with an efficient and reliable experience. In order to store the data mined from GitHub, we will be using MongoDB due to its high efficiency, ease of access, and its seamless integration with Python and Django. MongoDB is a document based database solution in which the documents consist of JSON-like format. This makes the importing of data from GitHub especially easy since the data is retrieved in JSON format from the API calls to GitHub.

There has been a lot of talk about displaying the data to the user via interactive graphics. This will be done using a utility called Plot.ly, as previously stated. Plot.ly is an open-source Python library that is used to create interactive graphics. It will allow us to create dynamic, interactive graphs of our data quickly and elegantly to help give the newcomer an easier time reading the data that is returned for their requested repository.

The last main portion of this product is hosting, this is a web application after all. In order to make this application available to the newcomers of OSS, we will be using a service called Digital Ocean. This service has a lot of features to offer our client in terms of scaling up the application as it gets bigger and also provides first class customer support to make sure that very few problems arise and that when they do, they are taken care of.

Now that we have a general idea of what technologies are at play here, it is time to take a bit more detailed look at the architecture of the product, which will be explained in the following section.

3. Architectural Overview

In the previous section, we discussed the implementation overview of Git-OSS-um. We will now be taking a high-level look at the architectural of our service, to better visualize what services will comprise the end product. The core of our software will be written using the Django python web framework, whilst making use of MongoDB to store and access files that we mine from the GitHub API. While there are many other components, Django will be the driver for all interaction between these components. This is visualized in Figure 1 below.

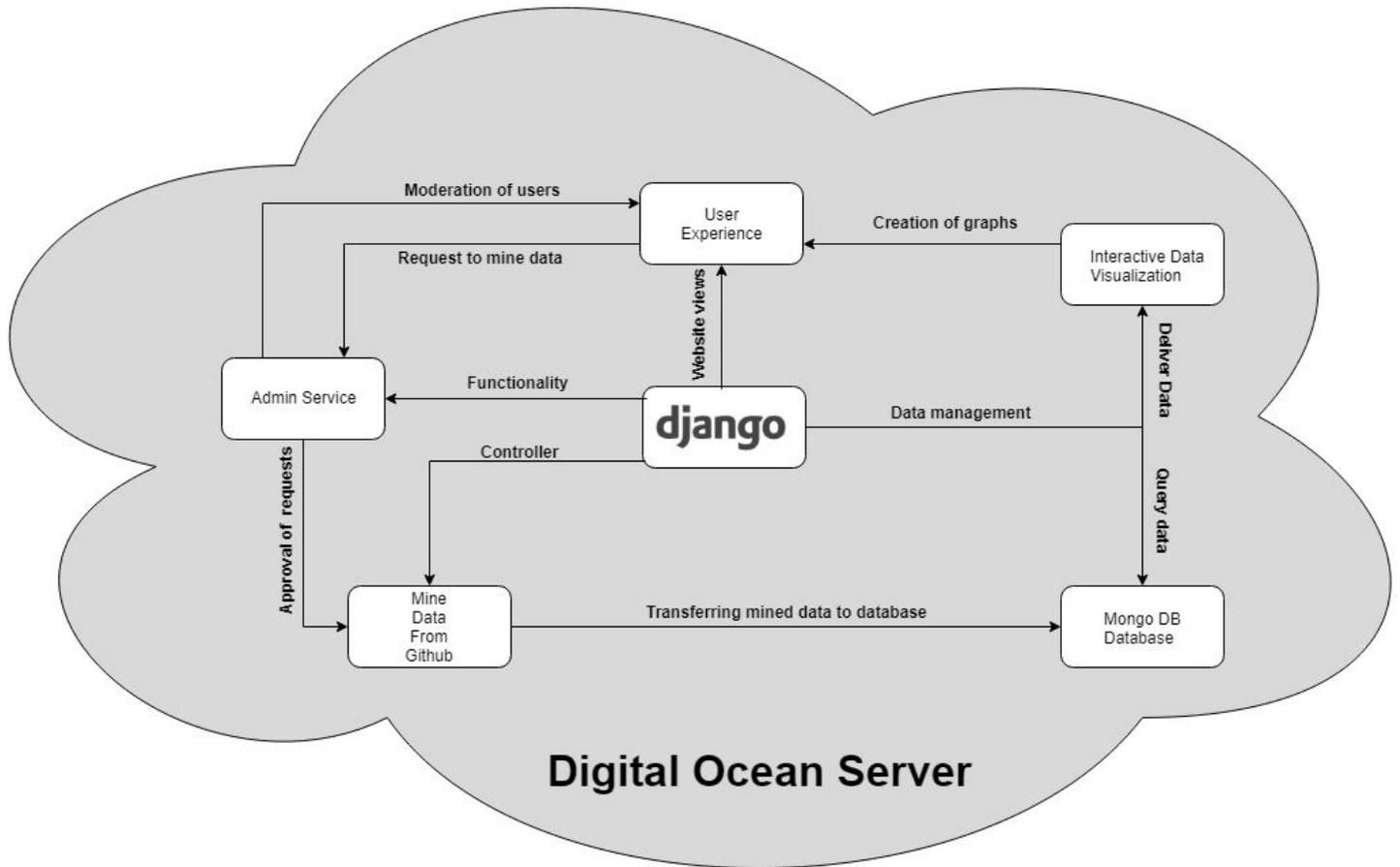


Figure 1: Architectural Overview of System Components

As explained briefly above, Django will be used to design, navigate, and operate all the necessary functions of our application. The team has decided to make use of Django 1.11 as it is a long term support (LTS) version of the framework that will be supported into 2020. It is at the center of our application, offering such necessities as admin functionality, login security, form validation, and much more.

Our team chose to use MongoDB as our database of choice to store and retrieve the data that we mine from the GitHub API. It is a document-based, NoSQL database that makes storing and locating large sets of JSON data a breeze. Using this technology, Django will be able to quickly and effectively query JSON data to build interactive data visualizations about any repository.

All extra functionality that is needed will be written in python, making it easy to plug into Django. This allows the team to quickly and efficiently write the software needed to mine data from GitHub, and store this data in a MongoDB collection. As stand alone python code, Django can handle when this code will be used, and who has the access rights to call it. Now with an understanding of the architectural overview of Git-OSS-um, we will dive further into the specific modules that comprise our software.

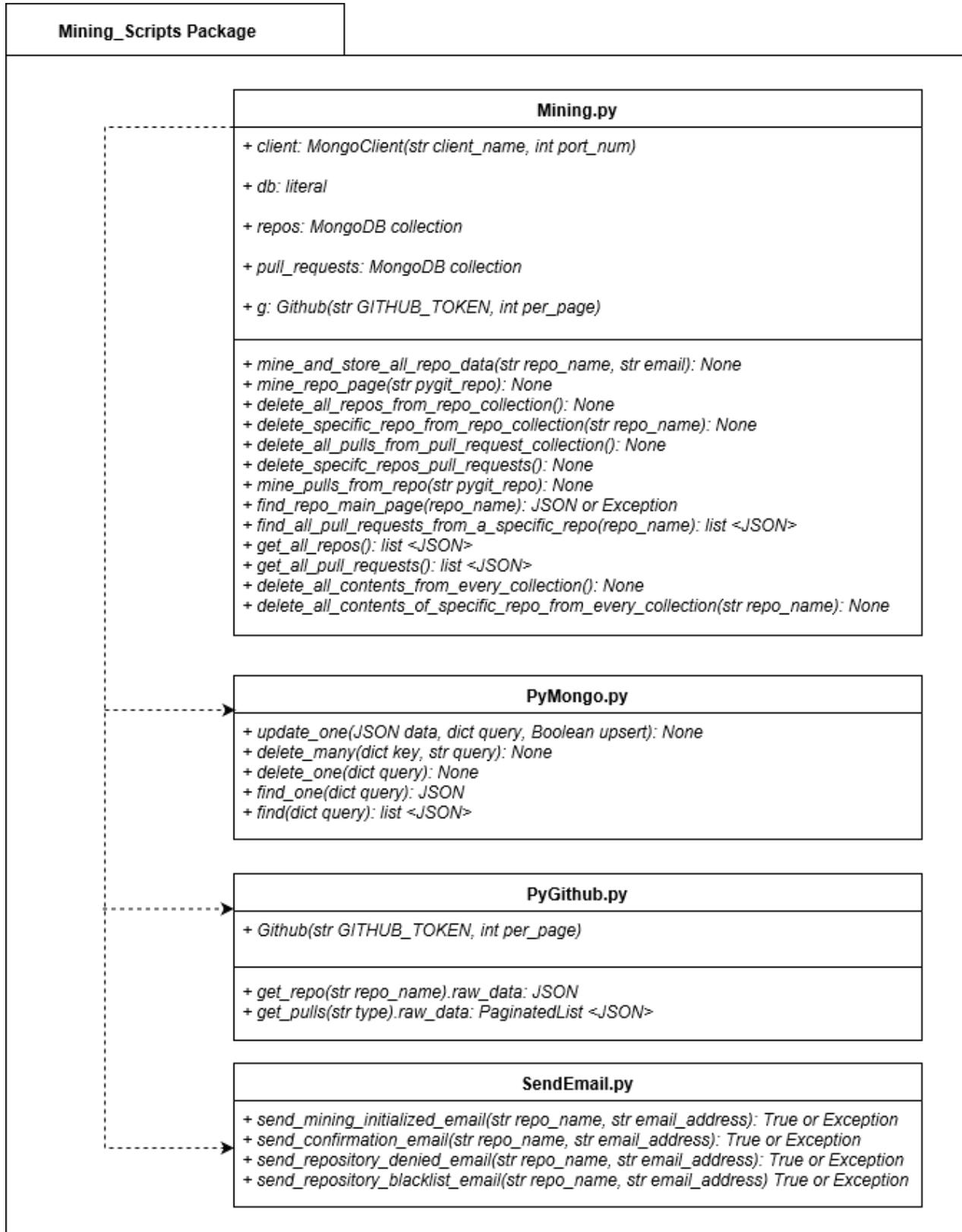
4. Module and Interface Descriptions

Within the previous section, we discussed the architectural overview of Git-OSS-um, and briefly explained its main components. We will now describe these components in greater detail in order to understand our web application in greater detail.

4.1. Mining Script

The mining.py component is accessible only by Django itself, when a request to mine GitHub data has been validated, the form has been cleaned of any potential security threats (such as HTML, or SQL injection), and the administrator has approved the request. This script contains functions written by the team that make use of several libraries in order to mine and store JSON data from GitHub, and to send an email notification to the user when this is finished.

We make use of PyGithub to quickly and efficiently mine JSON data from GitHub, Pymongo to interact with our MongoDB database collections, and SMTPLib to send email notifications to the user via the Gmail server.

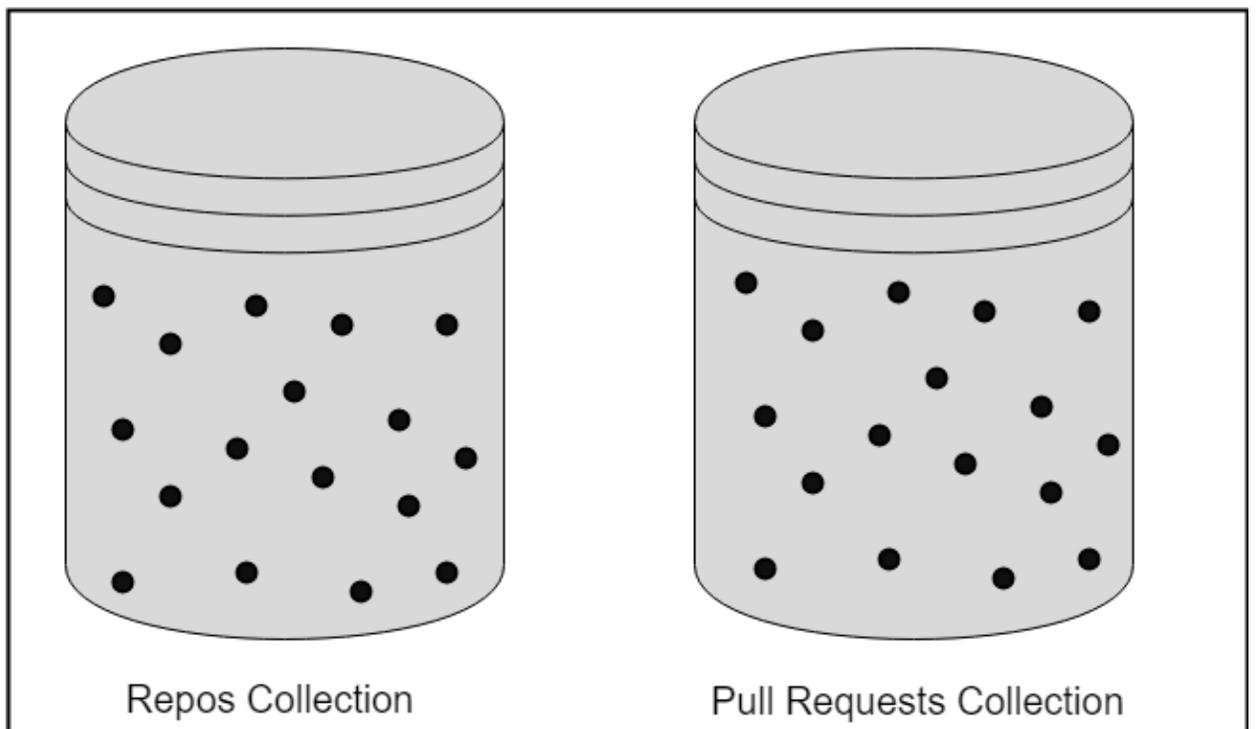


Figure_2: Mining Scripts Package Diagram

4.2. MongoDB

MongoDB will be running as a process on our DigitalOcean server on port 27017, and will never be directly accessible by the user. As mentioned in the previous section, all interaction with MongoDB will be done via the use of Pymongo within the mining.py component of our web application. We have designed our NOSQL database in way that it is simple, intuitive to learn, and accessible by our python scripts. The database is made up of two separate collections for storing JSON files from GitHub: Repos, and PullRequests. The Repos collection is responsible for containing all JSON files that represent the landing page of a GitHub repository, while the PullRequests collection will contain all pull requests that comprise any repository that we mine. Figure 3 illustrates what the database will look like.

MongoDB Database (● = JSON File)



Figure_3: MongoDB Database Overview

It is important to note that Git-OSS-um will not be managing the fields of the JSON data that have been mined from GitHub. Rather, we store the files as they are in the proper collection, and extract what data is necessary for visualization. Figures 4 and 5 display partial examples of actual data that will be stored in our MongoDB database.

```
1 // 20181110140906
2 // https://api.github.com/repos/rails/rails
3
4 {
5   "id": 8514,
6   "node_id": "MDEwOlJlcG9zaXRvcnk4NTE0",
7   "name": "rails",
8   "full_name": "rails/rails",
9   "private": false,
10  "owner": {
11    "login": "rails",
12    "id": 4223,
13    "node_id": "MDEyOk9yZ2FuaXphdGlvbWJyMjM=",
14    "avatar_url": "https://avatars1.githubusercontent.com/u/4223?v=4",
15    "gravatar_id": "",
16    "url": "https://api.github.com/users/rails",
17    "html_url": "https://github.com/rails",
18    "followers_url": "https://api.github.com/users/rails/followers",
19    "following_url": "https://api.github.com/users/rails/following{/other_user}",
20    "gists_url": "https://api.github.com/users/rails/gists{/gist_id}",
21    "starred_url": "https://api.github.com/users/rails/starred{/owner}/{repo}",
22    "subscriptions_url": "https://api.github.com/users/rails/subscriptions",
23    "organizations_url": "https://api.github.com/users/rails/orgs",
24    "repos_url": "https://api.github.com/users/rails/repos",
25    "events_url": "https://api.github.com/users/rails/events{/privacy}",
26    "received_events_url": "https://api.github.com/users/rails/received_events",
27    "type": "Organization",
28    "site_admin": false
29  },
30 }
```

Figure 4: Example JSON Representing a Repository Landing Page

```

1 // 20181110141021
2 // https://api.github.com/repos/rails/rails/pulls/9523
3
4 {
5   "url": "https://api.github.com/repos/rails/rails/pulls/9523",
6   "id": 4411988,
7   "node_id": "MDExO1B1bGxSZXF1ZXN0NDQxMTk4OA==",
8   "html_url": "https://github.com/rails/rails/pull/9523",
9   "diff_url": "https://github.com/rails/rails/pull/9523.diff",
10  "patch_url": "https://github.com/rails/rails/pull/9523.patch",
11  "issue_url": "https://api.github.com/repos/rails/rails/issues/9523",
12  "number": 9523,
13  "state": "closed",
14  "locked": false,
15  "title": "ActiveSupport::Notifications::Instrumenter#instrument should yield",
16  "user": {
17    "login": "stopdropandrew",
18    "id": 6681,
19    "node_id": "MDQ6VXN1c2Y2ODE=",
20    "avatar_url": "https://avatars1.githubusercontent.com/u/6681?v=4",
21    "gravatar_id": "",
22    "url": "https://api.github.com/users/stopdropandrew",

```

Figure 5: Example JSON Representing a Pull Request

4.3. Front-end Components with Django

We are implementing a full Django web application, which condenses our web application into one centralized package, instead of having separate front and back-end applications. Below are the components of our front-end with Django and the functionality of the components.

4.3.1. Logging Into The Application

Users will be able to sign in to our application, however, it is not necessary for browsing data. Users are encouraged to create an account so that they can post mining requests to administrators. For user sign up and login to our application, we will use Django forms. Django comes with built-in form functionality so we can build forms with ease and less code. Our login form asks users to input their username and password, we have built error handling to check our database whether a username has been registered to our product. For our sign up form for users, we ask for their name and email as a required field and their GitHub OAuth token as an optional field. If a user wants to request mining data, then they would have to register their OAuth token to mine the data.

4.3.2. Admin Portal

One of the requirements of our web application is to have an admin portal where administrators can approve or decline mining requests made by users and to manage database collections. Django offers a template interface, but can be altered and edited to accommodate our needs. The main functionality of our admin portal is a database management system where administrators can manage data in our MongoDB collections and to accept or decline mining requests made by the user. This central interface will show the administrators what requests are pending the mined JSON data that represents a repository in GitHub. This interface will also allow administrators to implement custom updates to the database, which means that an admin can see if a JSON file is up-to-date with the information that resides within it, and update that file as needed.

4.3.3. About Us Page

As a landing page, we will be creating an “About Us” page that would describe the product and team. This page will describe our web application and the core functionality and who the development team are. This page will briefly describe the main functionality such as viewing of mined data of a repository with interactive visualizations to explore the data. Alongside this, users can also cross compare up to three different repositories to further explore the data. A description of how to get a specific repository data will be explained on the page as well to clarify the process to new users. The “About Us” page will briefly introduce the development team as well to give users a background of who developed the application. The “About Us” page will briefly introduce the team’s sponsor, Igor Steinmacher, and explain how the project would be used for his research.

4.3.4. Requesting Data

In order to request data from a specific repository, a user will have to fill out a form with specific fields. Two main parts of this form would be the repository name and a user’s email address so that they can get a message saying that their request of that repository has been sent. This form will be able to be accessed in our Django admin portal, where administrators can either accept or decline the request for mining. After a decision is made, a user will get an email notification of the decision. If the request has been approved, a user’s Oauth token will be used to mine the data. After the data is mined, user’s can see the data and interact with visualizations.

4.3.5. Data Visualization

The python-nvd3 framework will be used to visualize data that is mined and existing data in our Mongo database. This framework is a Python wrapper of the nvd3 framework, which

harnesses the power of the JavaScript framework D3. The python-nvd3 framework allows us to write Python code and renders it as JavaScript and place that JavaScript in our HTML files. Because python-nvd3 also has a Django wrapper called django-nvd3, we can place python code directly into our Python files to generate our visualizations. The python-nvd3 framework also harnesses interactivity of D3. Functionality for hovering over certain parts of the visualization and constrict the visualization to certain parts of the data. Our visualizations offer an in-depth review of data pulled from a repository. For more detail, read more in the “Python-nvd3 Framework” section.

4.3.6. Cross Comparison of Repositories

The ability to compare up to three different repositories is a feature that can be used to further show what repository may be the best for the user to contribute to. The comparison will display both data and visualizations to further show users the differences of each repository. Overall, the cross comparisons of repositories and the visualizations will provide greater clarity to users.

4.4. Python-nvd3 Framework

Python-nvd3 is a python wrapper of D3.js where the library’s main goal is to create reusable charts. Python-nvd3 is a powerful tool because of its ability to keep the power that D3 initially offers and python-nvd3 renders JavaScript for us. The rendering of Javascript allows us to make a visualization in python and save the output as an HTML file and copy the contents into our web pages so that the graph can be displayed. The library also has built-in functionality for simplistic charts such as line, pie, and bar graphs, which is a perfect tool for quick visualizations and for prototyping purposes. The following sections explain what the team plans to do with the power of python-nvd3.

A key part of our front-end is interactive graphics where the user can learn new information regarding a mined repository. The python-nvd3 library lets us write python and renders JavaScript as a result. The python-nvd3 library is a wrapper that limits the overall code that we have to write. The rendering of JavaScript is displayed as an HTML file, in which we can copy and paste that JavaScript into HTML files. Because we are using Django as our front-end, the python-nvd3 library has another wrapper called django-nvd3 where we can write python directly into our web application. Writing python into a file such as view.py allows us to create template HTML files where we can reuse the graphics and interactivity.

The python-nvd3 library maintains the power of the D3.js framework, which means that we will have the power of D3 and the reusability that python-nvd3 gives us. The python-nvd3 library will be used for the reusability of charts and the interactivity of D3 so that users can explore the data of a mined repository. As users hover over parts of a visualization, they can uncover more information regarding certain data points and reveal more information about multiple repositories.

Our mining scripts obtain multitudes of information pulled from select repositories, but some of this information may not be useful to our users. It is important that we show only the relative data a user wants to know. This can be achieved by adding further interactivity to our visualizations. In select graphics, users can select what data that they would want to view on our visualizations. However, not every data option will be displayed, it is up to the team and our sponsor to select the necessary data to be displayed from mined repositories.

5. Implementation Plan

Beginning the Spring semester of 2019 in mid January, the team had a working prototype as a proof of concept in which checked off the core functionality of the product

outlined in the Requirements Specification document, signed by our team lead as well as our client.

As one of the most important tasks in the early stages of implementation, the team went about creating a checklist of all features and major milestones that needed to be met along with tentative due dates and team members responsible for each task. This checklist will be given to the client as well as a way to outline exactly what the team will be doing and what to expect by the final product deployment.

We plan to have not just the alpha build by Spring Break, but a mostly finished product in order to utilize the remaining portion of the semester for testing of the application and preparation for the final release of the product in which the newcomers can use. The following image displays our checklist featuring all of the tasks and major milestones with their due dates and responsible team member. This checklist may be subject to change but is a mostly complete timeline of the entire implementation period. The checklist is mostly detailed for before Spring Break because this is our main development phase. After Spring Break, our focus is mostly on testing, polishing the current build, and discussing any additional features with our client in which we will compromise testing time for addition of a feature.

Done	Task Description	Week Due	Who is Working on it?
x	Django About Us page and mock-ups for the rest.	1/29	Gary, Van
x	Mining.py script that will mine information from github and interact with MongoDB via python 3.6	1/29	Stephen
x	Mining_tests.py to ensure the quality/integrity of the code is high. This will be used to simulate interacting with the script, so when we move to RESTful practices, we already know the functionality works as intended	1/29	Stephen
x	Software Design Document outline	1/29	All members
x	Software Design Document Rough	2/5	All members
x	Plot.ly Bar Graph and Line Graph Prototypes (at least 1 or 2)	2/5	Van
x	Software Design Document Final	2/12	All members
x	Post requests for Mining (Front End)	2/19	Gary, Stephen
x	Admin Portal layout completed	2/19	Gary, Van
In progress	Connect Application to https://gitossum.com domain which will be actually at https://138.68.40.43:PORT_NUM	2/19	Stephen, Gary
x	Design review II Slides	2/19	All members
In progress	All front end pages layout complete with navigation	2/26	Gary
x	Admin Portal completed	2/26	Gary, Van
x	Post request data visualized on front end	2/26	All members
Coming soon	Design review II Presentation	2/28	All members
Coming soon	Tech Demo (Alpha) before Spring Break	3/12	All members
Coming soon	Software Testing Plan	4/2	All members
Coming soon	Design Review III	4/11	All members
Coming soon	Acceptance Test Demo	4/30	All members
Coming soon	Final project report, User manual and completed checkoff sheet, signed by client	5/9	All members

Figure 6: Team checklist for tasks leading up to Spring Break

6. Conclusion

Open source software is not the easiest field to be involved in, with our product, we wish to make the contribution process easier to newcomers. Our web application will be able to mine data straight from GitHub, store the data in a database for future use, and visualize that data in the form of interactive graphics which are displayed to the newcomer. In this document we have discussed the “big picture” associated with our application, implementation and architectural overviews of the product, the modules and interfaces that will be used to implement our solution, and a detailed plan of attack that maps out all tasks and milestones that must be completed by mid-March in order to deliver a mostly complete product early to leave excess time for testing as preparation for deployment to the newcomers of OSS.

By the time spring break arrives, we are confident that we will have a minimum viable product ready to present to our sponsor as both a research tool and a product to assist newcomers into contributing to open source software. As described in our document, we will create a web application that will mine data from GitHub repositories and provide the data and visualizations to our users. Our product will deliver satisfaction to the user, make the contribution process easier, and bring the power back to the newcomer.