

Github Data Analysis

Wide Project Proposal

YUAN Yingzhe
Computer Science and Engineering
HKUST
Hongkong, China

GUO Jiong
Computer Science and Engineering
HKUST
Hongkong, China

Abstract—GitHub have huge user group. Every moment there are hundreds of activities happen. To analyze the GitHub data continually, we use some opensource big data component such as Hadoop, spark, hive to construct a GitHub data processing framework, which aims at handling massive activity records and analyze data to find out patterns of users, repositories and languages. Using python crawler to get data from GitHub. By different data processing skills, the data is transport by certain channel to provide the downstream usage. Building the data warehouse base on the structure and content of the data and make three layers of the data warehouse. By spark, making data analyzation and store the result in the relational database. Finally present the graph and table on the dashboard.

I. Introduction

GitHub is an open source platform for developers to share their codes and collaborate with others. Every day there are innumerable user activities happening and these are collected by the GitHub. The data contains varieties of activities on different repos. From these data, we can extract many extra information like the recent hottest language. One important things on this mission is the data storage and data process. Large data set will cause huge pressure on the traditional relational database. Therefore, we need to seek a different way to implement the data analyze system, which can afford high load, store huge data set and have high fault tolerance. For these purposes, we choose several appropriate methods and techniques, and successfully build a system of data warehouse to deal with high load data.

II. Dataset

All the data we will use in this project are from Github. There is not an existing complete dataset to get enough well-structured data for our project, so we need to collect data and preprocess by ourselves. In this section, we will illustrate why we choose Github's data and what the data sources are.

A. Why Github

Github is the largest developer community and the largest platform for developers to share their codes and collaborate with others in the world. By analyzing the data from Github, we can get some insights about the development of open source projects, the development

of programming languages, and the development of the developer community.

B. Source I: Open API

Github provides a series of open API for developers to access most of the informations on Github. This 'Github REST API [?]' is mainly used for creating integrations, retrieve data, and automate workflows, but we can use it to query specific data needed for our project.

There are mainly two main areas of data:

- Users: Basic open informations of user, such as id, name, email, following users, repositories owned, etc.
- repositories: Basic information of public repositories, such as id, name, owner, description, stars etc.

We will provide some examples results of API query in Appendix??.

There are limits of current official REST API:

- Some need authentication, which means we need to provide a valid token to access the data.
- Some API cannot be called too frequently, which means we need to wait for a while before we can call it again.
- Most of API need an id as its query parameter, so we need to get the id first before we can query the data, which means we need to construct a massive id list.

After considering the above limitations, we decide to use the Github Archive [?] as our main data source and use the REST API to get some additional data or search some specific content.

C. Source II: Github Archive

GH Archive [?] is a project to record the public GitHub timeline, archive it, and make it easily accessible for further analysis. In fact, it also based on Github's Open API, while only use the activity API:

¹ <https://api.github.com/events>

GH Archive will crawl github's activity data in real time and sort them as json file in order of when they occurred. Activity archives are available starting from 2011-02-12, and are updated amost half day, so the total size of data is TB level.

Considering too large size of the data and the the earliest data using a different format, we will only use the latest part of them. In fact, the current amount data of data in an hour is about 700MB, so the latest data is surely enough for our project. The corresponding compression packages with specific time can be downloaded easily, and the minimal time interval of packages is one hour.

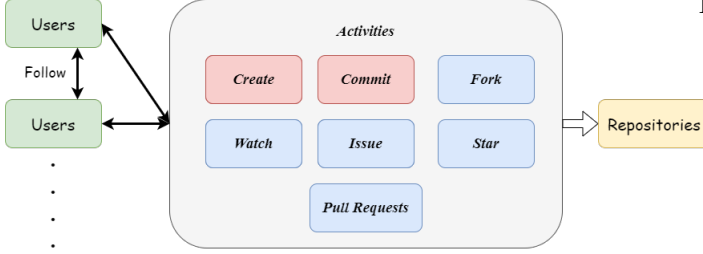


Fig. 1. This figure shows relation between user and repository and basic activities.

The archived activity data contains almost all 20+ event types provided by Github, which ranges from new commits and fork events, to adding members to a project. We will only pay attention to some key event type, such as commit, create, etc.

We will give some example of the data format in Appendix??.

D. Summary

Based on almost all activity record, Github's Open API can help us to get almost all related information. For example, from GH archive, we know user a push a commit at 12:00 in repository A, then we can get almost all details of a and A by calling Github's API.

III. Task

Github has provided some basic data aggregation for user to explore. For example, in Trending [?], Github shows the most popular repositories in the specific time range, language and spoken language. In Topics [?], Github shows repositories' topics according to tags of repositories. However, these data analysis are very basic and not wide to cover most of the information in Github.

In this project, we will use spark to construct a real-time processing framework, which aims at handling massive activity records and analyze data to give a more comprehensive analysis of Github.

There are main tasks in this project:

A. Data Preprocess

- 1) Collect data. Crawl Github Archive [?] and use Github API [?] to jointly query more detailed information about activities according to indices in Github Archive. Though Github Archive update several times a day, we will simulate it to a real-time stream to process data in real-time.

- 2) Preprocess data by Spark. For complete each json record, split it into several parts, filter useful parts and store them in different tables.
- 3) Data Persistence. Design different schema and construct key meta tables to store necessary information, statistics result or some data not easily obtained directly.

B. Data Analysis

- 1) Statistics. We will try to construct a dashboard to show some basic statistics of Github, such as the repository with the highest frequency of commits, the most frequently commits time period, etc.
- 2) Mining. There are plenty of text data from repository's file itself, commit record, issues, etc. We will try to use NLP to extract some useful information from these text data, such as more accurately topics of repositories, technical hotspots, etc.
- 3) Prediction. We will try to predict the future trend of Github, such as the number of repositories created daily in the future, trends of some programming languages, etc.
- 4) Relations. One of the main two parts of Github is Users. They connect by following each other or participating in the development of the same repository, for example, make pull requests, or give issues, etc. We will try to analyze and construct a social graph.

C. Data Visualization

- 1) Dashboard. We will try to construct a dashboard to watch the process of data process and show the results of data analysis by deploying a webpage.
- 2) Interaction. We will try to make the dashboard interactive, such as users can choose the time range of data to analyze, or choose the language to analyze, etc. Results of corresponding process not only comes from persistent database, but also from real-time stream produced by Spark.

IV. System Framework

A. Cloud Server Configuration

1) Memory and Cores: The system architecture contains many components. The metastore of hive is deployed on mysql, so that hive can be used by multiple users. Therefore, we need to choose a server with large memory to accommodate all components. Considering the system requirements, we choose a server with 4 cores and 16G memory.

2) Server Type: Our task is to build a data analysis service platform that continuously obtains log information and do data analyze jobs. We chose preemptible instances because of their relatively low price. And we have developed a set of fully automated scripts to deploy the cluster, so that if the server needs to be released unexpectedly,

the service can be easily deployed to the newly created server.

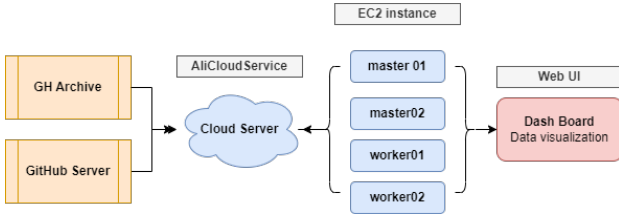


Fig. 2. A figure shows types of servers

3) Final Cluster Choice: For the final cluster, we use 4 server, which are named as master01, master02, worker01 and worker02. Different server played as different character in the cluster.

The Server details information:

- Instance type: Preemptible Instance
- Specifications: 4v CPU 16 GiB (I/O Optimized)
- Storage: 40GiB Cloud disk for each server
- Network: 100Mbps

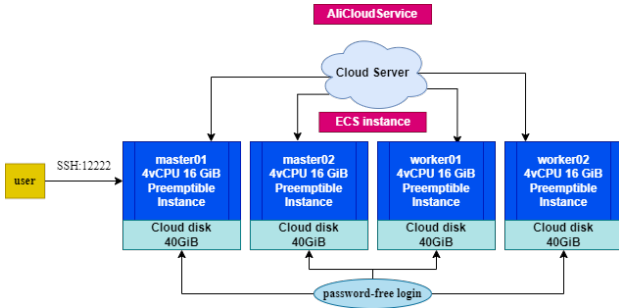


Fig. 3. A figure shows details configuration of clusters

4) Auto Configuration: Preemptive servers may be automatically released for lack of resources by Aliyun, so we need to purchase new computation instances, deploy the original big data system quickly and recover data collected when that worst situation happens. Based on this requirement, we have implemented a complete automatic configuration shell script, so that only one script need to be executed on any one single node to complete the environment configuration and deployment of the whole cluster, which greatly improves the efficiency of the cluster deployment and reduces the workload of the administrator.

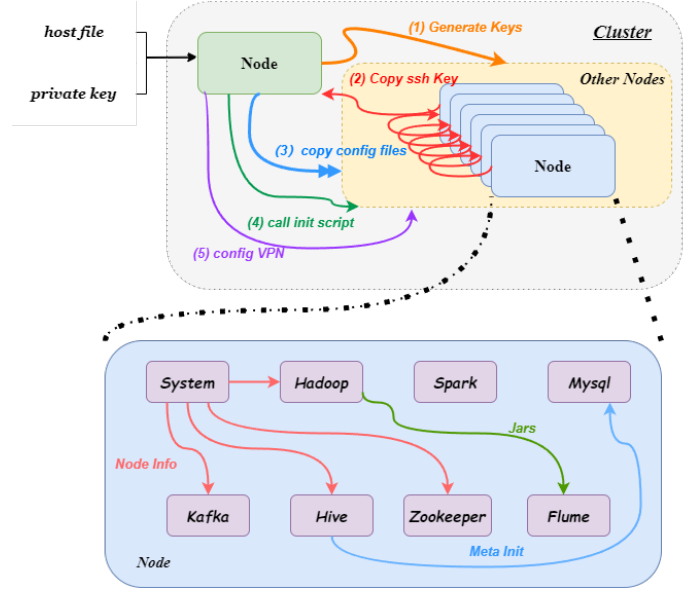


Fig. 4. A figure shows the control flow of auto configuration

Cluster: For the whole cluster:

- 1) Upload private key (for log in server initially) and host file (record the ip and name of all nodes) to any one nodes A.
- 2) Execute autoconfig shell script, which will do the following things one by one:
 - a) Require each node to generate a new ssh key pair and update sshd configuration.
 - b) Each node broadcasts its own public key to all other nodes for password-free login. Then only use the new own key to log in to other nodes after that.
 - c) Synchronize all necessary config files of components(hadoop, spark, etc) from A to all other nodes.
 - d) Execute the installation script on each node's background.
 - e) After the installation is complete, all nodes need to config a VPN proxy to access the outer Internet for downloading Github data.

All big data components will be installed when installation script is executed. The installation script will do the following things one by one on each node:

- 1) Install system software, update system environment variables and create necessary directories.
- 2) Modify configuration files of each component according to the node's role and special information, such as zookeepers'id, etc.
- 3) Download and install all components one by one, and copy the configuration files to the corresponding location.
- 4) Copy essential Jars of hadoop and spark to other components' lib directory, such as flume, which is

used for ensure consistence of the version of all components to avoid potential problems.

- 5) MySQL will be installed on worker02, modify access authority, and init database schema. The JDBC driver will be installed on all nodes.
- 6) Init hive metastore database on worker02. The hive schema will be stored in MySQL.

5) Security Control: After deployment, servers are frequently attacked from unknown sources, including DDoS, mining viruses, worms, and remote controlling, which causes the extremely high CPU usage, high disk IO, and unable to access the server. Killing processes, resetting servers, and recollecting data took up over 40% of our time

Although it has been observed that the behavior after entering the server is to inject code into the process, it is not clear by what means the server was successfully accessed and permissions obtained. We have taken the following steps to try to avoid being attacked:

- Change the default port of sshd: The default port of sshd is 22, which is the most frequently attacked port. We changed the port to 2222, which greatly reduced the number of attacks.
- Only allow strict key login: Some instruction of automated configuration steps need to be provided password, we generated a random length of 40 characters password by openssl and only use it once. After the configuration is complete, the password will be deleted, and all communication will be based on the rsa key.
- Close unnecessary ports: We only preserve port for ssh login and data visualization for external network access, and other ports of big data components are only open to the internal network.
- Install anti-virus software: We installed ClamAV on each node to scan the system periodically.

All above measures cannot completely prevent attacks, and we believe there is an unknown vulnerability. The corresponding logs are in the appendix.

B. System Architecture

Our goal is to build a distributed file system and computing analysis system based on cloud servers. After identifying the servers and clusters, we need to decide where to deploy the components and how data will flow between the different components within the architecture.

1) Technology selection: Considering our needs, we need to build a distributed file system and distributed computing system for continuously input log information, and need to store related information such as users and warehouses. So we choose the most commonly used technology in the field of big data: Hadoop architecture. Considering the data form we need, we use hive to store the data as hive tables. For faster analysis, we use spark for distributed computing and analysis. According to the

main thread, we add more components to it to maintain the continuous operation of the whole system.

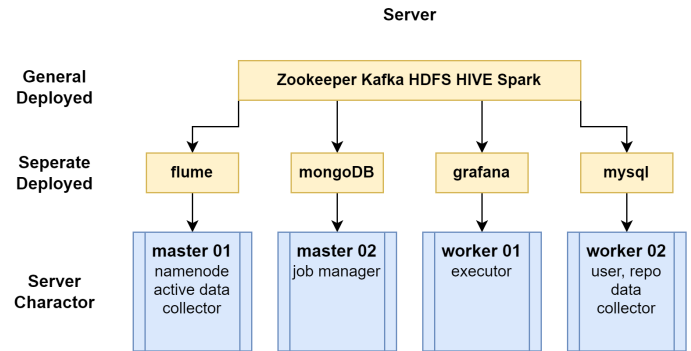


Fig. 5. A figure show the deployment of the system

2) Deployment: Balance the server workload so that the cluster can perform at its best. Therefore, it is important to plan the placement of components. Below are the locations of the different components on the server and their basic functionality.

For the entire analysis system we use the following components. Here we briefly introduce these components:

- 1) Zookeeper: ZooKeeper is a centralized service for maintaining configuration information, naming, providing distributed synchronization, and providing group services.[1] <https://zookeeper.apache.org/>
- 2) Kafka: message queue, which temporarily stores data and sends it downstream. It is the data source and key part of subsequent data stream processing.
- 3) Hadoop: Hadoop includes HDFS, Yarn and MapReduce. HDFS is the main storage place where we store data. HDFS is deployed on each server. Use Yarn as a resource manager. The resource manager is deployed on master02.
- 4) Hive: Hive data is stored in HDFS. The hive engine is very inefficient, so only use hive as storage and data management on HDFS. Hive is deployed on the entire cluster.
- 5) Spark: Spark is deployed on the entire cluster as the main component for data processing and data analysis.
- 6) MySQL: The data of users and warehouses are stored in MySQL, and MySQL is deployed on worker02. At the same time, the hive meta store is also on MySQL.
- 7) Flume: Track the data in the data capture folder and transfer it to Kafka. It will also transfer data from Kafka into a temporary folder in HDFS. The data capture process is deployed on worker02.
- 8) MongoDB: MongoDB is deployed on master02 for expose the API.
- 9) Grafana: Grafana is deployed on worker01. It's the interface of data visualization.

C. Data Flow

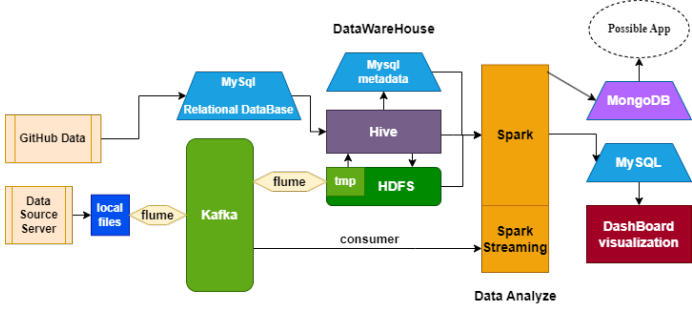


Fig. 6. A figure shows the data flow of the system

The server structure section introduces the general data flow direction, and here is a detailed description of the data flow in the overall analysis system.

1. Data source: The data sources are divided into two parts. The first part is GitHub user behavior log information. The source of this part of the information is the GitHub Archive website, from which one hour of data can be obtained every hour. We temporarily store the acquired data in a folder of worker02 through the Python program for subsequent data capture. The second part is GitHub user personal information and repo information. The source of this part of information is the official server API provided by GitHub. We use Python crawlers to obtain user and repo data through these APIs and write them into the MySQL database, and the database keeps updating these data.

2. Data transmission channel: Data transmission channels mainly include flume and kafka. The obtained log information tracks the file information through Flume, and transmits the continuously arriving data content to Kafka. Kafka internal data can be used by spark streaming through consumer. At the same time, Kafka and HDFS are connected through flume. The data accumulated in Kafka is consumed in batches and stored in the temporary folder of HDFS.

3. Data load: MySQL memory has continuously updated user and repository information. The temporary folder of HDFS has continuously added log information. We incrementally synchronize MySQL data to Hive through the spark program. Build it into a Hive table. The data of the Hive table is stored in HDFS. At the same time, the log information in the temporary folder is imported through Hive, and the original data is transferred to Hive for subsequent use by using the JSON serde parsing package.

4. Data Warehouse Construction: Build a data warehouse based on the acquired data. The details will be explained later.

5. Data analysis: Spark directly extracts the data in the Hive data warehouse to analyze various indicators through the built data warehouse, and writes the analyzed

results into multiple different databases for different purposes: Write to MySQL for subsequent data display. Write to MongoDB to expose Api externally. Write to Hive for historical data storage. Spark streaming directly obtains data through kafka for real-time data processing, so as to obtain updated and fast real-time data display.

D. Data Warehouse Construction

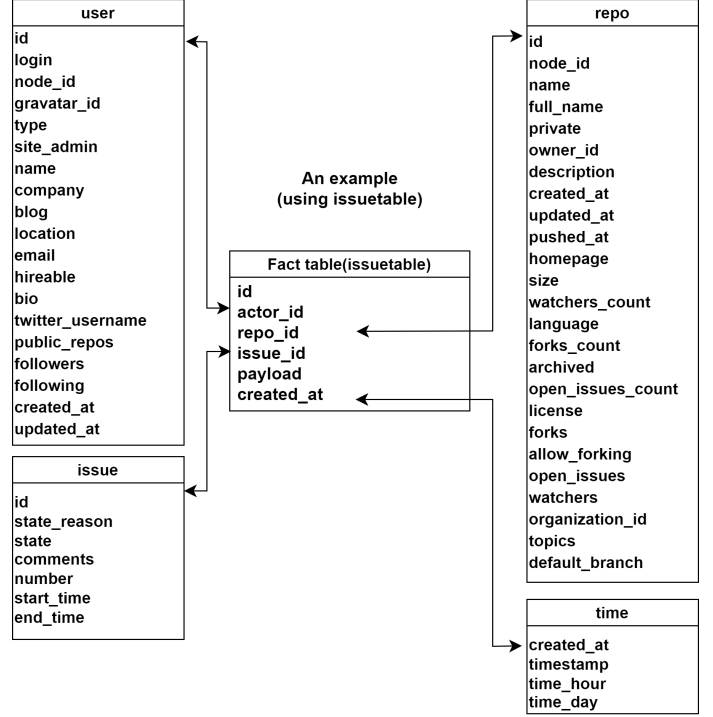


Fig. 7. A figure shows the schema of the data warehouse

- 1) Modeling: The basic structure of data warehouse modeling adopts constellation model. There is a star schema for each transaction table, and different transaction tables share the same dimension table to form the entire constellation model.

Data granularity selection: The log data itself comes from every operation of GitHub users, and there is a corresponding time record for each operation. Since the data is acquired in one batch at the same time during the data acquisition process, the time granularity is selected to be as fine as 1 hour.

Determination of dimension table and transaction table: the dimension table adopts four dimensions of user, repo, issue and time, and the transaction table is divided into 8 types of transaction tables, such as create, delete, pull, and push, according to the different behavior of users.


```

10     ...
11   }
12   ...
13 }

```

C. Get All Commits of a Repository

1 <https://api.github.com/repos/apache/spark/commits>

```

1  [
2    {
3      "sha": "0fde146e8676ab9a4aeafebb1684eb7a44660524",
4
5      "commit": {
6        "author": {
7          "name": "Juliusz Sompolski",
8          "email": "julek@databricks.com",
9          "date": "2023-03-24T01:56:26Z"
10        },
11        "committer": {
12          "name": "Hyukjin Kwon",
13          "email": "gurwls223@apache.org",
14          "date": "2023-03-24T01:56:26Z"
15        },
16        "message": "...",
17        "tree": {
18          "sha": "308faa10a134402eb4e0796c54b4e4dd000fc7dc"
19        },
20        ...
21      },
22      "author": {
23        "login": "juliuszsompolski",
24        "id": 25019163,
25        "node_id": "MDQ6VXNlcjIIMDE5MTYz",
26        ...
27      },
28      "committer": {
29        "login": "HyukjinKwon",
30        "id": 6477701,
31        ...
32      },
33      "parents": [{}]
34    },
35    ...
36  ]

```

```

19    "size": 1,
20    "distinct_size": 1,
21    "ref": "refs/heads/master",
22    "head": "6b089eb4a43f728f0a594388092f480f2ecacfd",
23    "before":
24      ↪ "437c03652caa0bc4a7554b18d5c0a394c2f3d326",
25    "commits": [
26      {
27        "sha": "6b089eb4a43f728f0a594388092f480f2ecacfd",
28        "author": {
29          "email":
30            ↪ "5c682c2d1ec4073e277f9ba9f4bdf07e5794dabe@rspt.ch",
31          "name": "rspt"
32        },
33        "message": "Fix main header height on mobile",
34        "distinct": true
35      }
36    ],
37    "public": true,
38    "created_at": "2015-01-01T15:00:01Z"
39  },
40  ...
]

```

Appendix B GH Archive

One single json record from the GH Archive.

```

1  [
2    {
3      "id": "2489651051",
4      "type": "PushEvent",
5      "actor": {
6        "id": 3854017,
7        "login": "rspt",
8        "gravatar_id": "",
9        "url": "https://api.github.com/users/rspt",
10       "avatar_url":
11         ↪ "https://avatars.githubusercontent.com/u/3854017?"
12       },
13       "repo": {
14         "id": 28671719,
15         "name": "rspt/rspt-theme",
16         "url": "https://api.github.com/repos/rspt/rspt-theme"
17       },
18       "payload": {
19         "push_id": 536863970,

```