

Freie Universität Berlin

Institute for Computer Science of Freie Universität Berlin

QIFTool - Documentation

ver1.0 - 21.10.2020

Robert Selack

robfu@zedat.fu-berlin.de

QIFTool or Query Issue Finder-Tool is a project created as a bachelor's thesis. It aims to help in the quality research field of technical debt by providing relevant discussions regarding these debts. The discussion are presented in form of issues from github. The tool uses keywords and additional metrics to find potentially interesting issues. Although it is meant for the field of technical debt the tool can be also be used to return all different kind of issues' topics.

Contents

- 1 Workflow 3**
- 2 How to run the program 3**
- 3 Configuration File - config.ini 4**
- 4 Interactive Mode 6**
 - 4.1 Functions 7
- 5 SQLite - Database 9**
 - 5.1 Repositories 10
 - 5.2 Issues 12
- 6 Expanding the tool 15**
 - 6.1 Metrics 15
 - 6.2 SQLite commands 16

1 Workflow

QIFTool is written in Python(3.8) and uses [Google's Custom Search JSON API](#) in conjunction with [Google's Custom Search Engine](#) to filter issues directly on Github. Keywords will be read out of a configuration file to determine which issues should be prefiltered. All prefiltered issues found by the engine will be inserted for research and caching purposes into a [SQLite](#) database. Afterwards it uses the official Github API ([PyGithub](#)) to look through the available pieces of metainformation inside each found issue and compares them to the other metrics set in the configuration file to only show the issues that match the requirements.

2 How to run the program

1. Download the 'qiftool.py' and 'requirements.txt' files from the repository
2. Place both files in the desired location
3. Open the terminal and navigate to the files' location
4. Install all dependencies by running 'pip3 install -r /path/to/requirements.txt' or just 'pip install -r /path/to/requirements.txt' depending on your python version
5. Run the program by using 'python3 qiftool.py'
6. By running it for the very first time the tool should have created a 'config.ini' file inside the tool's folder. Fill out the necessary parameters following the instructions in [3](#)
7. With the 'config.file' filled out run the program again just like in step 5
8. The tool should now operate properly and an interactive mode will be seen. Follow [4](#) for further instructions

3 Configuration File - config.ini

This file is created by running the program for the very first time. It is used to give the user a space to use their own parameters used by the tool. The file contains three sections for the user to fill out.

```
[DEFAULT]

path_of_database = current
path_of_download = current


[credentials]

github_api_key = randomnumberandlettersinlowercase
google_api_key = randomlettersinuppderandlowercaseandsymbols
google_cse_id = randomleterandnumbersinlowercase


[metrics]

keywords = technical debt    refactor    rewrite
issue_comments = 5
repo_contributors = 50
```

1. [DEFAULT]

- this section contains the path for the database and downloaded repositories to be stored in. The user is able to create their own path with the location of the 'qiftool.py's as a pivot. These can be changed by providiung a valid path on your machine.

2. [credentials]

- this section contains the corresponding credentials necessary to run the used APIs

(a) github_api_key

- i. register on github
- ii. [use this link](#) and click on 'generate new token' to create a new key

iii. paste the key as a parameter

(b) `google_api_key`

i. register on google

ii. [use this link](#) and click on 'Get a Key' to create a new key

iii. either choose a project or create a new one

iv. follow the instructions and paste the key as a parameter

(c) `google_cse_id`

i. login to the google account created in the prior step

ii. [use this link](#) and click on the project you used to create the google key with

iii. look for the 'Search engine ID' and paste the ID as a parameter

3.

4. `[metrics]`

These contain the metrics used for the google search. For further details for understanding each metric please refer to the tables in [5](#).

(a) `keywords`

- string of characters with each element separated by a tabulator. Note that the keywords will be used to find patterns that exactly match. So 'refactor' will find 'refactoring' but not vice versa. In addition the keywords are connected with a logical and.

(b) `issue_comments`

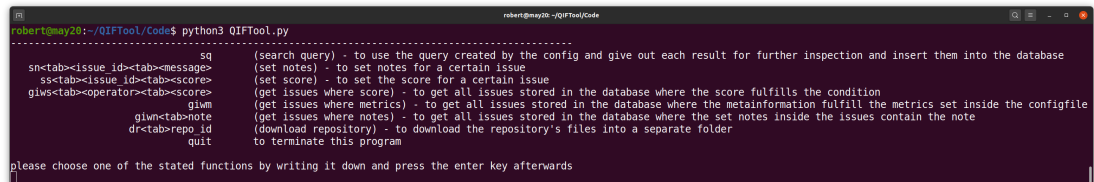
- an integer over 0. It will show issues that have at least the amount of comments set in this metric. So 5 will yield issues with 5 or more comments.

(c) `repo_contributors`

- an integer over 0. It will show issues that have at least the amount of contributors working on the corresponding repository. So 5 will yield issues with more 5 or more contributors working on its repository.

4 Interactive Mode

Once you successfully configured the configuration file in 3 an interactive mode will be seen on the console after running it. In this mode the program will wait for the user



```
robert@may20:~/QIFTool/Code$ python3 QIFTool.py
-----
sq          (search query) - to use the query created by the config and give out each result for further inspection and insert them into the database
sn<tab><issue_id><tab><message> (set notes) - to set notes for a certain issue
ss<tab><issue_id><tab><score> (set score) - to set the score for a certain issue
glvs<tab><operator><tab><score> (get issues where score) - to get all issues stored in the database where the score fulfills the condition
glwm        (get issues where metrics) - to get all issues stored in the database where the metainformation fulfill the metrics set inside the configfile
glwn<tab><note> (get issues where notes) - to get all issues stored in the database where the set notes inside the issues contain the note
dr<tab><repo_id> (download repository) - to download the repository's files into a separate folder
quit        to terminate this program
-----
please choose one of the stated functions by writing it down and press the enter key afterwards
```

Figure 1: default interactive mode

to simply type a desired function into the console and confirming it by pressing enter. After being done with a function the program goes back to displaying the interactive mode as it loops itself around it.

4.1 Functions

function	description
sq	(search query) - start the google search. The metrics set in the configuration file will be used to determine what results will be found and shown.
sn<tab><issue_id><tab><message>	(set notes) - sets a note for a certain issue inside the database issue_id - a string of numbers. Found within the <i>issue_id</i> field in either the output or database of the issue. message - a string of characters that will be inserted into the <i>notes</i> field inside the database.
ss<tab><issue_id><tab><score>	sets a score for a certain issue inside the database. issue_id - a string of numbers. Found within the <i>issue_id</i> field in either the output or database of the issue. score - a number chosen by the user to represent its relevance.
giws<tab><operator><tab><score>	(get issues where score) - displays all issues stored in the database where the score fulfills the condition set by the user. operator - all comparison operators allowed by the SQL. <, >, =, <=, >= score - a number chosen by the user to represent its relevance and compare the issues inside the database to.

function	description
giwm	(get issues where metrics) - displays all issues stored in the database where the pieces of metainformation fulfill the metrics set inside the configfile. This function yields the same functionality as the 'sn'-function but with the database being the source.
giwn<tab>note	(get issues where notes) - displays all issues stored in the database where their notes contain the note set by the user with this very function note - a string of characters. This can be used in conjunction with SQL-syntax like providing a " or % around the note.
dr<tab>repo_id	(download repository) - downloads the repository's files into a separate folder. This folder's location is set by the configuration file. The structure of the downloaded files also is identical to that of its repository.
quit	terminates this program.

5 SQLite - Database

This tool uses the SQLite version 3.33.0 (2020-8-14) library. The database created with this tool consists of two tables with one table referring to the other in a 1:n relation.

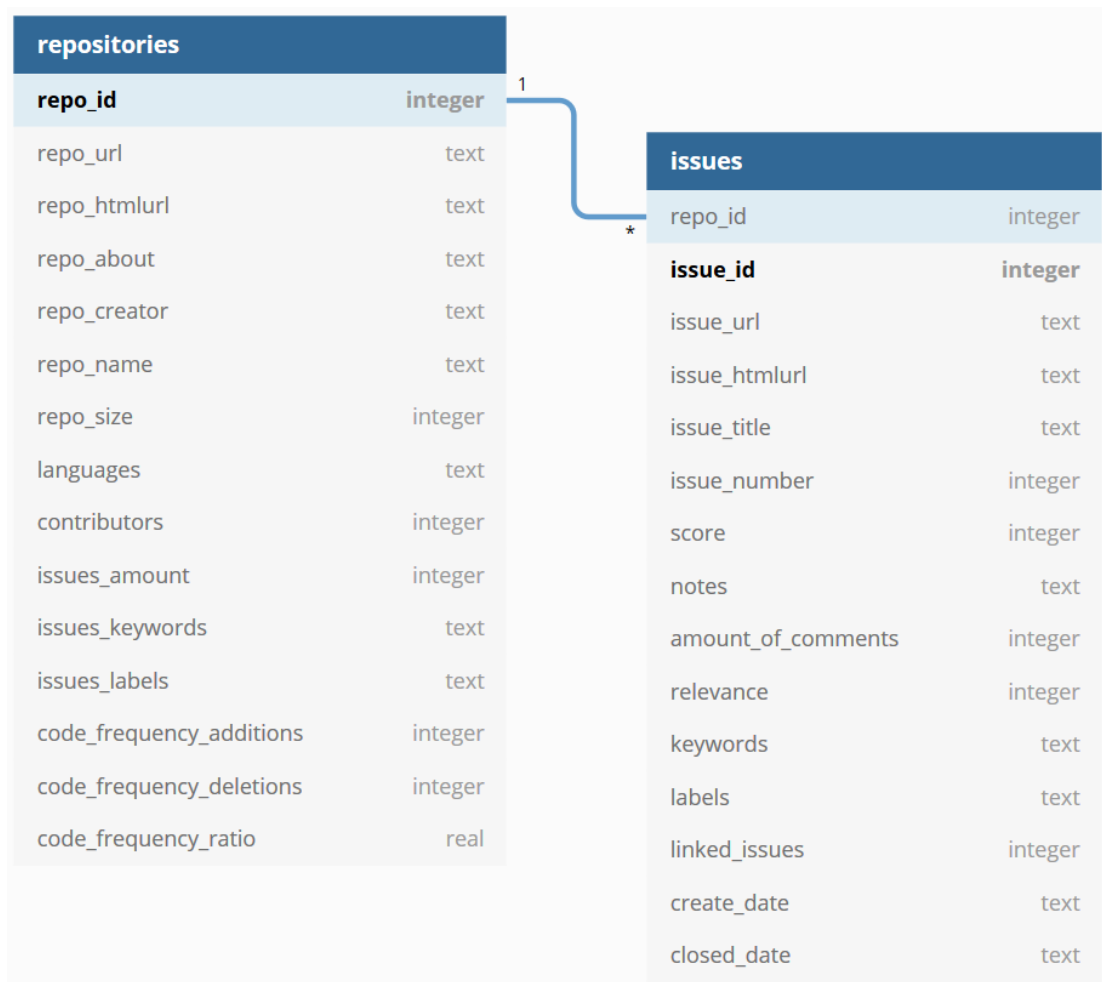


Figure 2: Database schema

5.1 Repositories

attribute	datatype	description
repo_id (primary key)	integer	identifier for a repository
repo_url	text	url for the JSON file of this repository
repo_htmlurl	text	url that refers to the web based github repository
repo_about	text	contains the 'about' of the repository found on github. It's a brief description of the repository
repo_creator	text	name of the creator of this repository or fork
repo_name	text	name of the repository
repo_size	integer	size of this repository in MB
languages	text	list of programming languages used in this repository
contributors	integer	amount of contributors of this repository
issues_amount	integer	amount of issues of this repository

attribute	datatype	description
issues_keywords	text	a list of all keywords that where used to find issues related to this repository
issues_labels	text	a list of labels that were used for the issues found
code_frequency_additions	integer	overall amount of lines of code added to this repository
code_frequency_deletions	integer	overall amount of lines of code deleted from this repository
code_frequency_ratio	real	quotient of the lines of code added and deleted. A value between 0 and 1 with 1 meaning all code that was added got deleted and 0 meaning all code that was added none got deleted.

5.2 Issues

attribute	datatype	description
repo_id (foreign key)	integer	identifier for a repository
issue_id	integer	identifier for an issue
issue_url	text	url for the JSON file of this issue
issue_htmlurl	text	url that refers to the web based github issue
issue_title	text	title of the issue
issue_number	integer	relative number of this issue created within its repository
score	integer	value to set by the user. Used for the user's own usage of a rank system. Makes it possible to rank found issues relative to each other in order to find more valuable issues easier later on
notes	text	string to set by the user. Used for the user's own organisation. Makes it possible to note interesting attributes about a certain discussion or topic sorting. It's possible to look for certain patterns inside the set notes
amount_of_comments	integer	amount of comments that the issue has

attribute	datatype	description
relevance	integer	value corresponding to the relative order of issues found with a query. The higher up (earlier) the issue has been found the higher its relevance. Third hit on page two equals a relevance of 23. These relevances change to show an all time best relevance every time the issue has been found.
keywords	text	list of all keywords that has been used to find this issue over all query iterations
labels	text	list of all labels that are used with this issue
linked_issues	text	list of all issues that are linked to within this issue. <i>This attribute has yet to be implemented</i>

attribute	datatype	description
create_date	text	date of the creation of this issue. Although the datatype is a 'text', SQLite still recognizes the string as a date due its formatting
closed_at	text	date of when the issue was closed. Although the datatype is a 'text', SQLite still recognizes the string as a date due its formatting. If the issue has not been closed yet this field is set to 'NA'

6 Expanding the tool

This tool can offer various fields of expansion. Most likely it will be an addition of other metrics. This section will explain how to proceed in order to expand the tool. Depending on the type of expansion the tool needs less steps to accomplish it.

6.1 Metrics

1. Inside `'create_config()'`: Add the desired metric as an attribute inside the `'[metrics]'` section and give it a default value.
2. Inside the `'Config'` class: Add the new metric to the `'__init__'` function
3. Inside `'read_config()'`: let the new metric be read out and stored inside the created instance of the `'Config'` class.
4. Inside `'create_database(path)'`: Add the new metric to the respective table and add the desired constraints as well as edit the `'UNIQUE'` modifier. `new_metric integer,`
5. Inside `'insert(conn, table, values)'`: Add the new metric inside the respective insert statement as well as the required `'?'`
6. Depending on what kind of metric and how accessible it is, it is necessary to create a function and preferably a class for the metric to be extracted and stored in. In case it needs to be created and takes API requests, do not forget to add the `reset_sleep(auth)'` condition. see function `'stats_code_frequency(repo, auth)'` for an example
7. Inside either `'RepoObj'` or `'IssueObj'`: add the new metric to the `'__init__'` function
8. Inside `'metric_check(conn, config_issue_comments, config_repo_contributors, issue_id, repo_id)'`: Add the new metric as a new condition. The conditions have a 2^n complexity for with n being the amount of metrics. This semantic should be rewritten in order to allow the addition of more metrics.
9. If desired to also print out the new metric, the class `'IssuePrint'` and function `'issue_print'` need to be modified by adding the new metric.

6. Expanding the tool

10. Inside `'page_iterator(auth, keywords, issue_comments, repo_contributors, google_api_key, google_cse_id, path_db)'`: the new metric needs to be called and then be stored inside the repo or issue object
11. Inside `'input_handler(init)'`: To access or make proper use of the new metric, add a new function that is callable via the interactive mode.

6.2 SQLite commands

These commands will be callable via the interactive mode and access the database and its entries directly.

1. Inside `'input_handler(init)'`: Add a condition to to access the new function. Note that the function will be determined via the use of tabulators as seperators.
2. Create a new function similar to `'get_issues_where_metrics(conn, keywords, issue_comments, repo_contributors)'` and modify the SQL statement to perform the desired action