This zip file contains The Technical Debt Dataset [1] as a set of CSV files.
In the rest of the document, there is a brief explanation of each of the tables that make up the dataset.

## PROJECTS

This table contains the basic information about the projects. It contains the project ID, the project name, as well as the links to the GitHub repository, the Jira issue tracker, and the SonarQube project key.

## SONAR_ANALYSIS

This table contains the analysis keys to join with the SONAR_MEASURES and SONAR_ISSUES tables as well as the analysis date and the revision (i.e. commit hash) that it belongs to. This table allows the integration of the different SonarQube tables with the Git information offered by tables like GIT_COMMITS.

## SONAR_MEASURES

SonarQube is one of the most common open-source static code analysis tools for static quality analysis. This table contains the different measures SonarQube analyses from the commits. As an example, in Figure 1 we can see some of the measures analysed which correspond to:

1. **vulnerabilities**
   Number of vulnerability issues.
2. **security_rating**
   A = 0 Vulnerabilities
   B = at least 1 Minor Vulnerability
   C = at least 1 Major Vulnerability
   D = at least 1 Critical Vulnerability
   E = at least 1 Blocker Vulnerability
3. **security_remediation_effort**
   Effort to fix all vulnerability issues. The measure is stored in minutes in the DB. An 8-hour day is assumed when values are shown in days.
4. **code_smells**
   Total count of Code Smell issues.
5. **sqale_index**
   Effort to fix all Code Smells. The measure is stored in minutes in the database. An 8-hour day is assumed when values are shown in days.
6. **sqale_debt_ratio**
   Ratio between the cost to develop the software and the cost to fix it. The Technical Debt Ratio formula is:
   `Remediation cost / Development cost`
7. **sqale_rating**
   Rating given to your project related to the value of your Technical Debt Ratio.
8. **duplicated_lines_density**
   = duplicated_lines / lines * 100
9. **duplicated_lines**
   Number of lines involved in duplications.
10. **duplicated_blocks**
    Number of duplicated blocks of lines.

11. **duplicated_files**
    Number of files involved in duplications.
12. **Complexity**
    It is the Cyclomatic Complexity calculated based on the number of paths through the code. Whenever the control flow of a function splits, the complexity counter gets incremented by one. Each function has a minimum complexity of 1. This calculation varies slightly by language because keywords and functionalities do.
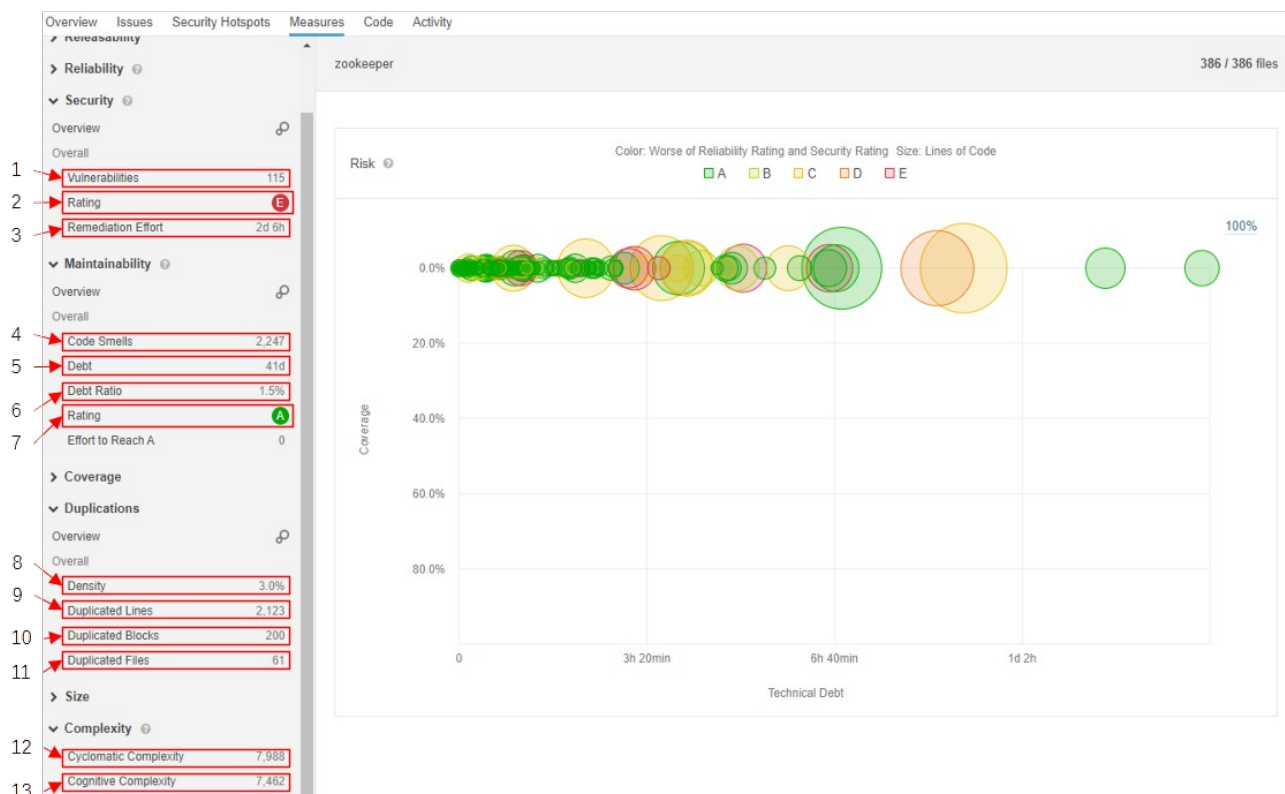13. **cognitive_complexity**
    How hard it is to understand the code's control flow.



*Figure 1: Example of some SonarQube measures.*

For more information on the SonarQube measures visit: https://docs.sonarqube.org/latest/user-guide/metric-definitions/

# SONAR_ISSUES

This table lists all of the SonarQube issues, as well as the anti-patterns and code smells detected in the analysed projects. In Figure 2 we can see some of the attributes of a SonarQube issue, which correspond to:

1. **COMPONENT**
   File containing the issue.
2. **RULE**
   SonarQube rule that has been triggered.
3. **TYPE**
   Type of the issue, which can be Bug, Vulnerability or Code Smell.
4. **SEVERITY**
   Severity of the issue, which can be Blocker, Critical, Major, Minor or Info.
5. **STATUS**
   After creation, issues flow through a lifecycle, taking one of the following statuses: Open, Confirmed, Resolved, Reopened or Closed.
6. **EFFORT**
   Estimated effort required to fix the issue.
7. **CREATION_DATE**
   Creation date of the issue.
8. **START_LINE**
   Start line of the issue in the file.
9. **TAGS**
   Custom tags of the issue.



*Figure 2: Example of SonarQube issues.*

For more information on the SonarQube issues visit:
https://docs.sonarqube.org/latest/user-guide/issues/

## SONAR_RULES

This table lists the rules monitored by SonarQube. In Figure 3 we can see some of the most relevant features of this table, which correspond to:

1. **NAME**
   Summary of the rule that serves as its name.
2. **DEF_REMEDIATION_FUNCTION**
   Remediation function used to estimate the effort to solve the issue.
3. **DEF_REMEDIATION_BASE_EFFORT**
   Estimated time to solve the issue.
4. **PLUGIN_NAME**
   Family of rules where the rule belongs to.
5. **PLUGIN_RULE_KEY**
   Rule identifier.
6. **TYPE**
   Type of the issue, which can be Bug, Vulnerability or Code Smell.
7. **SEVERITY**
   Severity of the issue, which can be Blocker, Critical, Major, Minor or Info.
8. **SYSTEM_TAGS**
   Tags used by SonarQube to better group the rules.
9. **DESCRIPTION**
   Description of the rule.

*Figure 3: Example of a SonarQube rule.*

# GIT_COMMITS

This table reports the commit information retrieved from the git log. As shown in Figure 4, the main attributes of the GIT_COMMITS table are:

1. **COMMIT_MESSAGE**
   Message summarizing the changes introduced with the commit.
2. **BRANCHES**
   Branches where this commit is present.
3. **AUTHOR**
   Developer who created the commit.
4. **COMMITTER**
   Developer who committed the commit.
5. **PARENTS**
   Commit hash of the previous commit.
6. **COMMIT_HASH**
   Commit hash used to identify the commit.



*Figure 4: Example of the main attributes of a git commit.*

# GIT_COMMIT_CHANGES

This table contains the changes performed in each commit. In Figure 5 we can see the main attributes, which are:

1. **NOTE**
   The description of what has been done in the commit.
2. **COMMITTER_ID**
   The ID of the developer who committed the commit.
3. **LINES_ADDED**
   Number of new lines added.
4. **LINES_REMOVED**
   Number of lines removed.
5. **FILE**
   The full path to the modified file.
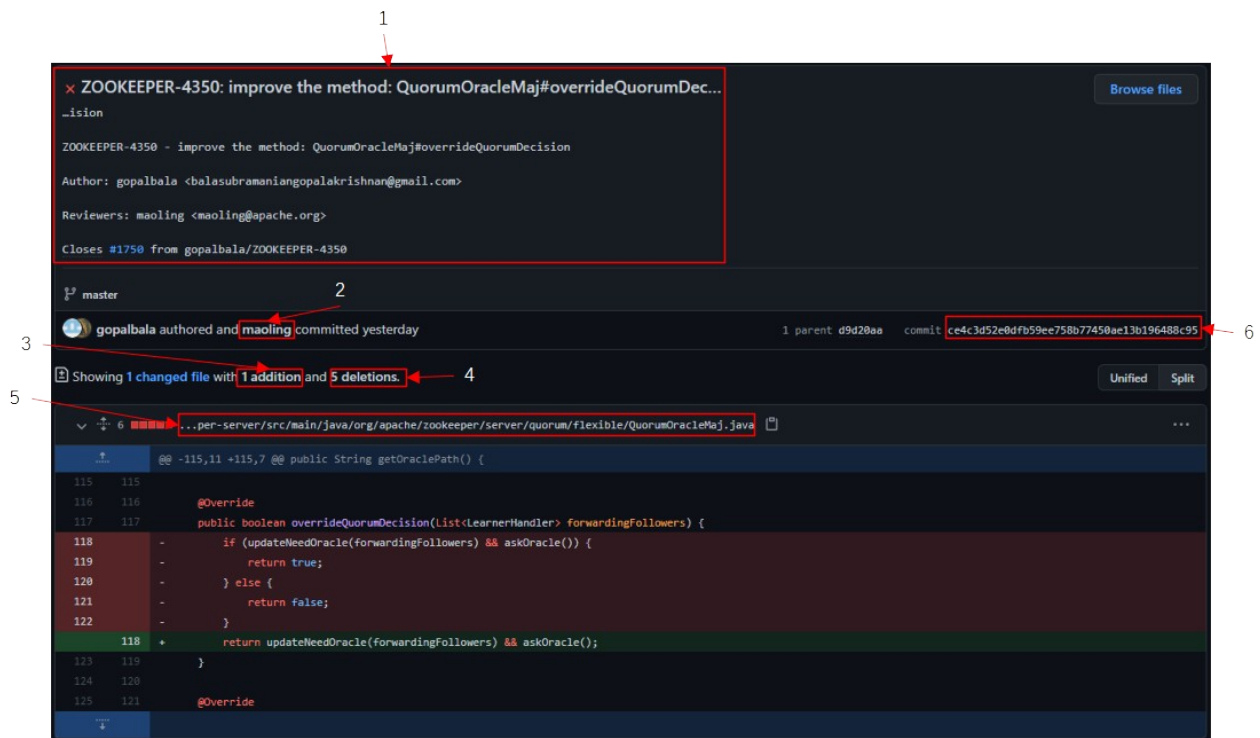6. **COMMIT_HASH**
   Hash used to identify the commit.



*Figure 5: Example of the main attributes of a git commit.*

# REFACTORING_MINER

Refactoring Miner is an open-source tool that classifies the different refactorings in the history of Java projects. This table reports the list of refactoring activities applied in the studied repositories. Figure 6 shows a small extract of the table describing its attributes.



*Figure 6: Extract of the table REFACTORING_MINER.*

# SZZ_FAULT_INDUCING_COMMITS

The SZZ algorithm tries to identify the fault-inducing commits from a project's version history. The algorithm is based on Git's blame/annotate feature and assumes that the fault-inducing commit of a fault is known. Usually, this is done by combining data from an issue tracker and form Git's log command. This table reports the results from the execution of the SZZ algorithm, which labels the fault-inducing and -fixing commits. Figure 7 shows a small extract of the table describing its attributes.



*Figure 7: Extract of the table SZZ_FAULT_INDUCING_COMMITS.*

## JIRA_ISSUES

Different organizations use Jira to track different kinds of issues, which can represent anything from a software bug to a project task, or a leave request form. This table contains the Jira issues for the analysed projects. Figure 8 shows a small extract of the table describing its attributes.



*Figure 8: Extract of the table JIRA_ISSUES.*

## Dataset schema

Figure 9 presents the Entity Relationship schema of the dataset. Regarding the fields in the tables, the fields needed to distinguish a row (i.e. primary key) have a key icon, the fields referencing a primary key of another table (i.e. foreign key) have a small arrow above their icon, and "..." means that the table has more fields than presented in the figure.
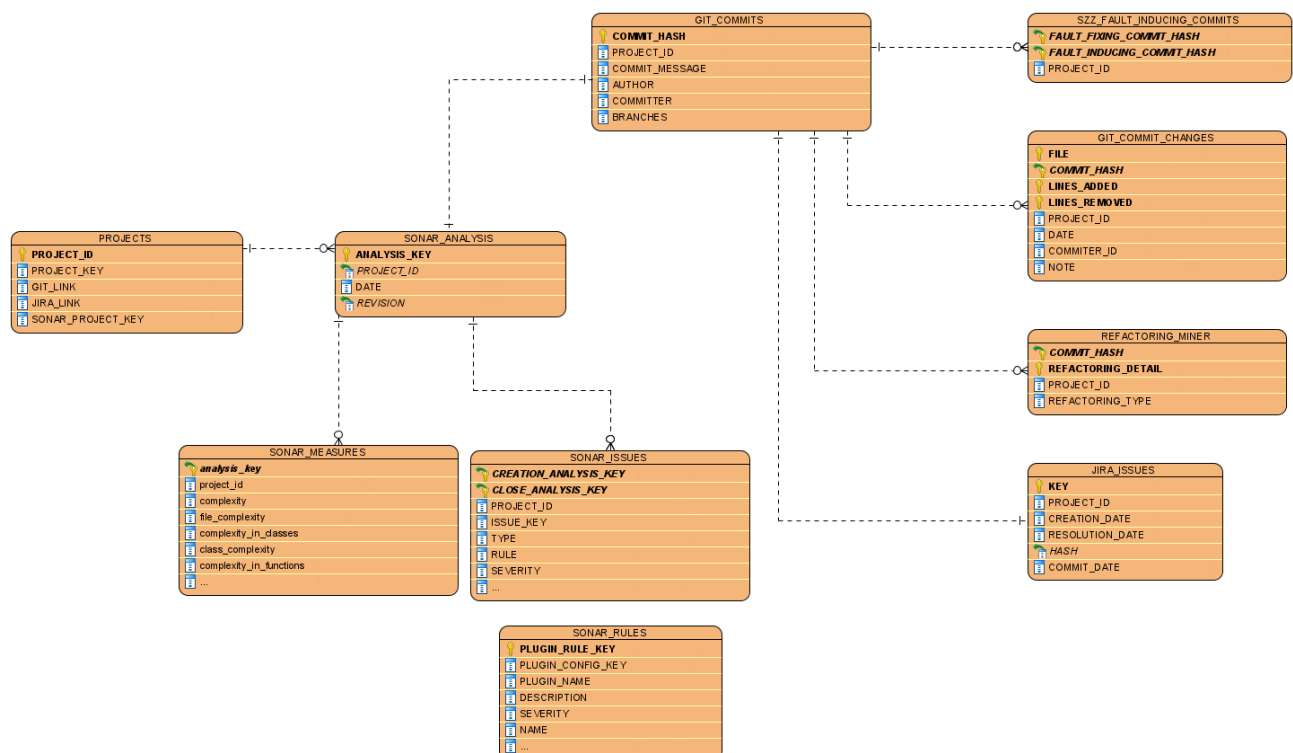


*Figure 9: Entity Relationship Diagram of the dataset.*

[1] Valentina Lenarduzzi, Nyyti Saarimäki, Davide Taibi. The Technical Debt Dataset. Proceedings for the 15th Conference on Predictive Models and Data Analytics in Software Engineering. Brazil. 2019. Download the paper