

Measuring Software Ticket Quality using Quantitative Data Analysis

Andrei-Mihai Nicolae (21473292)

April 19, 2018

ABSTRACT

Software tickets are of valuable importance to the whole computing science field - they guide engineers towards better planning, management and tracking of their progress throughout complex projects. However, there are few studies that investigate what makes for a high quality, valuable ticket. This can lead to multiple issues in a company, such as increased communication friction between developers and end users filing bug reports, as well as increased overall costs due to waste of development effort. In this research paper, we present our findings after investigating a large number of variables surrounding software tickets, such as whether the presence of stack traces influence the time to close for the ticket. Our results show that the presence and type of attachments, comments complexity (i.e. number of comments per ticket and total number of words), summary and description complexity, grammar correctness scores as well as the sentiment drawn from the comments can influence the quality of the ticket. We bring a couple of novel aspects to the research community including one of the largest dataset statistically analysed in the field, as well as state-of-the-art sentiment and grammar correctness analysis.

1. INTRODUCTION

In the past decade, technology has drastically increased its influence on virtually every aspect of our society. Therefore, software projects have inherently become more complex and require increasing number of developers in the team. Due to this, software engineers have created issue tracking systems, a means of planning, managing and tracking work products in the form of *software tickets* or *issues*.

There are multiple platforms for providing such issue tracking systems, among which the most popular are Jira [1] and Bugzilla [16]. For both platforms, the tickets are split into two main categories: feature requests (i.e. feature to be implemented into the system) and bug reports (i.e. issue encountered by an end user or discovered by a developer in the codebase). Regardless of the type of ticket, they possess various information that can be filled in by the reporter (i.e. person who created the ticket; can be both an end user or a developer in the team), providing the developers a detailed view of what is requested or what went wrong.

Even though tickets provide such comprehensive data regarding a specific task, studies have shown that fixing bugs is one of the most common tasks performed by developers [14]. One of the reasons for this is the communication friction between developers and end users [12] as developers might need clarification regarding what information the users have

provided (e.g. cannot reproduce the bug, screenshot is unclear). Another main reason for this waste of effort on solving tickets, according to Just et. al [11] and Zimmermann et. al [21], is the generally poor design of issue tracking systems. This can lead to various issues, including increased costs for the company, wasted development effort, decreased customer satisfaction and overall poor performance.

Therefore, there is a need in the community to find the answer to what makes for a high quality, valuable software ticket that would improve the overall performance of the development team and, inherently, the company. As there are many fields in a ticket, there are numerous unanswered questions, such as whether stack traces have an influence on the quality of a ticket?

In this research paper, we present and discuss our findings after running a quantitative data analysis on over 3200,000 tickets taken from more than 15 open source projects. We have implemented a Go application with multiple commands (i.e. store, analyze, plot, statistics) that can automatically fetch any number of tickets from a Jira instance, analyze them, generate plots and run statistical tests.

During the analysis part, we investigate several variables in correlation with *Time-To-Close*, which we define as the *metric of quality*. *Time-To-Close* represents the period of time between the creation and the closing of a ticket; more specifically, the creation of the ticket is marked when the status of the ticket is set to *Open* and the closing of a ticket is considered when the ticket status is set to *Closed* (or some similar status, such as *Fixed*, *Resolved*). *Time-To-Close* is our dependent variable, and as independent variables (i.e. variables controlled in order to test the effects on the dependent variable) we have set a number of ticket fields.

In order to provide answers to what factors influence quality in a software ticket, we answer the following seven research questions in this paper (Section 5):

- does the presence of attachments and their type (e.g. code snippet, screenshot) influence the *Time-To-Close* for a ticket?
- does the presence of stack traces improve the ticket quality?
- does the presence of steps to reproduce reduce the *Time-To-Close*?
- is there a relationship between the number of comments influence *Time-To-Close*?
- is there a relationship between

This study brings several contributions to the research community:

- an innovative tool was built for the purpose of this project, its strengths lying in its simplicity, speed and extensibility;
- it is one of the few studies in the field that performs a quantitative analysis rather than a qualitative one;
- it is one of the very few research projects that investigates such a large number of tickets (over 300,000) extracted from 38 different projects;
- it is, to our knowledge, the first study to conduct sentiment and grammar correctness analyses on software tickets.

In Section 2 we iterate over state-of-the art studies in the field and then continue with Section 3 where we discuss how the ticket data set was collected and analysed. We continue with describing the data set (Section 4) where we provide insights into various aspects of the data (e.g. size of database, number of tickets with attachments) and then discuss about correlations between the variables (Section 5). Finally, we provide future research directions in Section 6 and present our conclusions in Section 7.

2. RELATED WORK

The study of Bettenburg et. al [2] showed what makes for a valuable bug report through qualitative analysis. After conducting interviews with over 450 developers, one of the main factors behind a quality ticket is grammar correctness in summary and description. Another aspect which was flagged as helpful by the interviewees was the presence of stack traces and steps to reproduce in tickets. A further contribution brought by the authors to the community is the creation of a tool called Cuezilla which is able to automatically predict the quality of a bug report with an accuracy of around 40%.

Another research paper that strengthens the argument that readability is a quality factor in software tickets is the one presented by Hooimeijer et. al [9]. They conducted their analysis on over 25,000 bug reports from the Mozilla project, investigating readability, daily load, submitter reputation, the whole changelog histories and severity. They conclude that not only readability in the textual fields of a ticket influence the *Time-To-Close*, but also that the presence of attachments and the number of comments have a clear effect on the duration of triaging. On the other hand, the patch count and other similar fields did not provide significant value to the quality of tickets.

However, there are other types on information typically included in software tickets which might prove beneficial for developers and subsequently improve the ticket quality. One such type is stack traces and Schroter et. al [18] analyse how quickly tickets are closed when they either have stack traces included in their fields (e.g. description) or not. They collected their data from the Eclipse project using the InfoZilla tool proposed by Bettenburg et. al [4]. Then, they linked the stack traces to changes in source code by mining the Eclipse version controlled repository. The results showed that around 60% of the bugs that contained stack traces in their reports were fixed in one of the methods in the frame.

Moreover, more than 40% of the tickets having stack traces got fixed in the first stack frame.

Software tickets can have duplicates, usually meaning that the most important fields, such as summary or description, describe the same feature request or bug report. Even though one might believe that they cannot bring any value to a software project, the work of Bettenburg et. al [3] shows the contrary. The authors collected large amounts of data from the Eclipse open source project and ran various kind of textual and statistical analysis on the data to find answers. After the results were computed, they concluded that usually bug duplicates contain information that is not present in the master report (i.e. the original report that was filed). Moreover, developers also specified that they have often found value in these duplicate bug reports and that they can even aid automated triaging techniques.

Bettenburg et. al [5] present in their work an application called infoZilla. This tool can parse bug reports and correctly extract stack traces, patches and source code. When evaluated on over 150,000 bug reports, it proved to have a very high rate of over 97% accuracy. We applied some of the techniques shown in the study and successfully managed to retrieve stack traces as well with a great accuracy.

One other type of information that could be extracted and used for valuable purposes is the summary of a bug report, as shown in the work of Rastkar et. al [17]. The authors wanted to determine if bug reports could be summarized effectively and automatically so that developers would need only analyze summaries instead of full tickets. Firstly, they asked university students to volunteer to annotate the bug reports collected from various open sources projects (i.e. Eclipse, Mozilla, Gnome, KDE) by writing a summary of maximum 250 words in their own sentences. These human-produced annotations were then used by algorithms to learn how to effectively summarize a bug report. Afterwards, the authors asked the end users of these bug reports, the software developers, to rate the summaries against the original bug reports. They eventually learned that existing conversation-based extractive summary generators trained on bug reports produce the best results.

We tried to apply the technique presented in the work of Fischer et. al [6], where they present a way to track specific version controlled branches from information in software tickets. They managed to successfully visualize the tracked features and illustrate their non apparent dependencies which can prove very useful in projects with a large number of components.

The work of Just et. al [11] examines a rather different aspect of quality in tickets - they are investigating the quality of the underlying issue tracking systems instead. The authors ran a survey on 175 developers from Eclipse, Apache and Mozilla. Then, they applied a card sort in order to organize the comments into hierarchies to identify common patterns. Their findings can be summarized in the following top seven suggestions:

- create differentiation between bug report difficulty levels between novices and experienced reporters;
- encourage users to input as much information as possible;
- add ability to merge tickets when needed;
- recognise and reward valuable bug reporters;

- integrate reputation into user profiles to mark experienced reporters;
- to be able to easily and expressively search through tickets.

However, bug report assignment is not the only part of bug triaging that needs investigation. Lamkanfi et. al [13] show a method through which one can predict the severity of a reported bug (i.e. assign story points to the ticket) automatically. The authors conducted their research on the Mozilla, Eclipse and GNOME projects and divided their method into four steps:

- extract and organize bug reports;
- pre-process bug reports - tokenization, stop word removal and stemming;
- train the classifier on multiple datasets of bug reports - 70% training and 30% evaluation;
- apply the trained classifier on the evaluation set.

The conclusions they drew from running the evaluation process were rather interesting:

- terms such as deadlock, hand, crash, memory or seg-fault usually indicate a severe bug; however, there are other terms that can indicate the opposite, such as typos;
- when analyzing the textual description they found that using simple one line summaries resulted in less false positives and increased precision levels;
- the classifier needs a large dataset in order to predict correctly;
- depending on how well a software project is componentized, one can use a predictor per component rather than an universal one.

3. BUILDING THE DATA SET

The first step towards building the data set was to create a tool that was able to execute all the commands that we required: fetching the tickets from any Jira instance, storing them into some form of database, analyze the variables of interest, automatically plot the correlations between them and run statistical analysis on the data. After careful consideration, we decided that Go was the best way to go for various reasons:

- it is designed with simplicity in mind, thus making it easier for others who might join the project to read and understand the codebase;
- it is compiled, statically typed, which implies that it is a much faster candidate than other interpreted languages such as Ruby or Python;
- it is designed with concurrency in mind, thus it helps reduce times of execution and computing power considerably;

Throughout the entire application, we tried to apply the UNIX philosophy of creating small applications that do one thing, but do it well. Therefore, we implemented clean and simple packages where we tried to use the standard library as much as possible so that potential contributors coming to the project would find it easy to start working directly on the code. Also, we designed the tool with extensibility in mind, so that other database providers (e.g. CockroachDB) or issue tracking systems (e.g. Bugzilla) could be easily implemented - this was achieved using the elegant Go interfaces.

We split our tool into four main commands: *store*, *analyse*, *plot* and *stats*. They all follow the flow shown in figure 1 and complete the whole application cycle, in the end producing a database of analysed tickets, as well as plots and statistics for investigating correlations, saved on the filesystem.

In the following three subsections we describe how we designed and implemented the commands mentioned above.

3.1 Fetch and Store

This is the first command implemented in our tool - it is designed to fetch, from any valid Jira URL, any number of tickets and then store everything inside an instance of BoltDB [10], which is a very simple yet powerful key-value pair database. The whole process is parallelized with the possibility of scaling even up to hundreds of goroutines (i.e. lightweight threads) running concurrently.

The decision of choosing Bolt in favor of well tested, more traditional databases such as Postgres or MySQL came naturally when we noticed that we did not need to query the database often, but rather get the whole array of tickets and manipulate them. Moreover, a MongoDB or Postgres DB would have required a running server. Even though the server could have been running on the local machine, Bolt instances are actually files with the extension *db* saved on disk, thus it is trivial to create backups.

We have interacted with the Jira Server REST API which has very good documentation written by the Atlassian team. After getting familiar with all the features the REST API offers, we have opted for getting paginated issues from the server, which means slicing the whole set of tickets into multiple smaller arrays to improve performance and reduce the load on the server. We also specified to Jira exactly what set of issue fields to return and, eventually, while responses were retrieved from the instance, we began storing them into our Bolt DB instance.

The way we chose to store them was to set the issue key (i.e. unique identifier that differentiates a ticket from all others across all projects inside a Jira instance) as the key of the pair and the value is the JSON representation of the ticket. This helped us reduce the size of the database once we reached hundreds of thousands of tickets.

3.2 Analyse

Once the Bolt DB instance is completely populated with all the tickets we were interested in, the *analyse* command first fetches all the issues in the database. Then, it runs in parallel the *seven types of analysis*:

- *attachments* - checks whether attachments influence Time-To-Close and also look at what types of attachments (e.g. screenshots, archive, code snippets) influence quality the most;

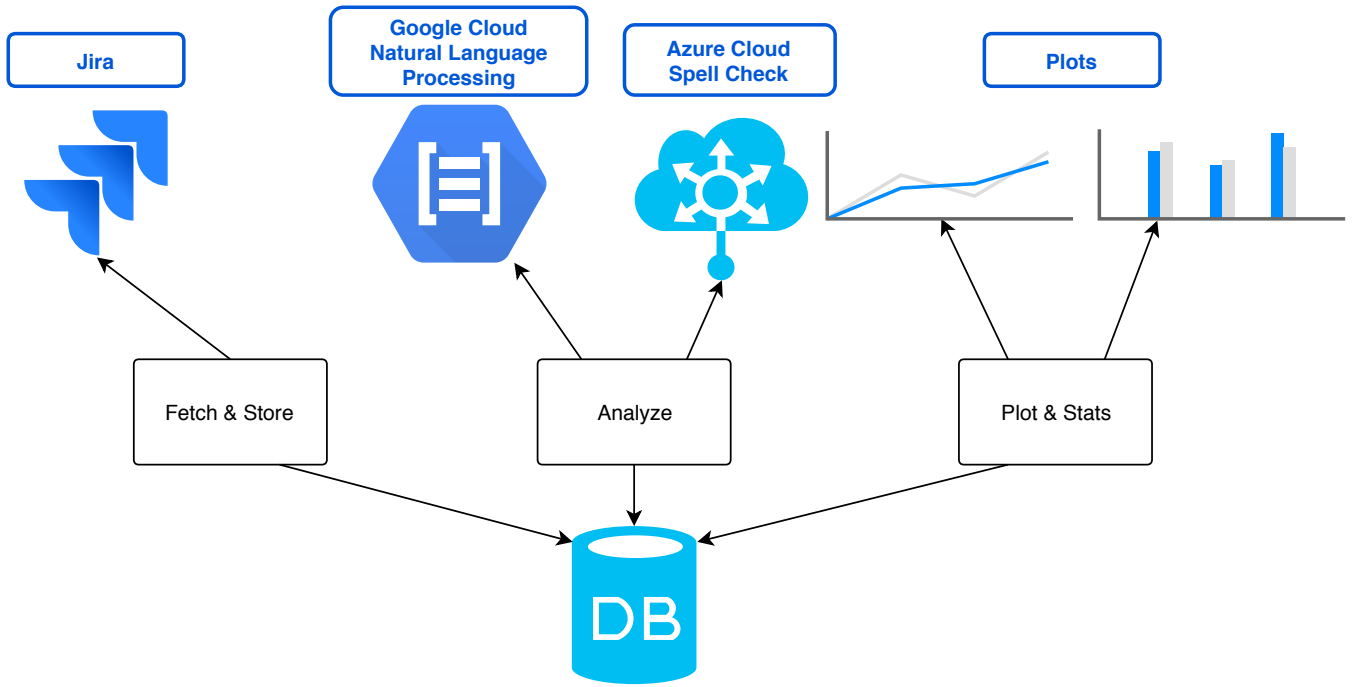


Figure 1: Application flow of the Go tool.

- *steps to reproduce* - performs a complex regex to detect steps to reproduce and then checks whether their presence influence the quality of a ticket;
- *stack traces* - runs complex regex for detecting stack traces in either summary, description or comments and verifies whether their presence influence Time-To-Close for the ticket;
- *comments complexity* - loops through all comments and counts all words; then, the tool verifies whether having many words or comments decreases or increases Time-To-Close for the ticket;
- *fields complexity* - same analysis as comments complexity, but only for summary and description;
- *grammar correctness* - it uses Azure Cloud Bing Spell Check API [15] to perform analysis on summary, description and comments; after concatenating everything and making it compatible with the API allowed formats, the tool receives back in the JSON payload not only the number of flagged tokens (i.e. grammar errors), but also their types (e.g. unknown token, misspell);
- *sentiment* - uses Google Cloud Platform’s Natural Language Processing API [8] to retrieve the sentiment score for summary, description and comments; it first concatenates everything and makes it conform to Google’s API and then sends the request, receiving a score from -1 to 1 inclusive (-1 is completely negative, 1 is completely positive).

Implementation wise, the application first checks Bolt for tickets and gets either all of them in memory (easier to parse,

but heavy on resource utilization) or slices them into smaller arrays to get processed afterwards (harder to parse as it eventually requires re-creation of the whole array of tickets before inserting back into the database, but it is much lighter on computing power). After having the tickets available, we start looping through them and perform all analysis types mentioned above. Once some sort of value is computed (e.g. grammar correctness score), it gets set in its corresponding field in the ticket struct and subsequently is stored back in the database. We chose this approach because we do not want to run the analysis every time we want to plot correlations between variables, but rather have the data already available there. Moreover, for external analysis such as hitting Google Cloud Platform APIs, there are costs for running the analysis, thus in this way we drastically reduce overall costs for the project.

3.3 Plot and Statistical Tests

Afterwards, we have created two extra commands for automatically generating plots (i.e. scatter plots and bar charts) and running statistical tests on the collected data.

For plotting the data, we first connect again to the Bolt database instance and then filter out only the issues that have the variable *Time-To-Close* set (i.e. tickets marked as closed when they were fetched from the Jira instance). Afterwards, we either save scatter plots or bar charts to disk.

Also, in order to validate our findings, the statistics command fetches all the data from the database and runs two types of tests: Welch’s T Test [20] for analysing categorical data (e.g. has/does not have attachments) and Spearman’s rank correlation coefficient [19] for investigating continuous data (e.g. how do different grammar correctness scores change *Time-To-Close*).

4. CHARACTERISING THE DATA SET

The data that we collected is stored, as previously mentioned, in a Bolt database instance which is around 6 GB in size. All tickets are stored as key-value pairs, where the key is the ticket's unique ID (e.g. KAFKA-100) and the value is the JSON encoded representation of the ticket. The tickets stored have the following fields available for investigation:

- attachments - files or images attached to tickets;
- summary - short description of the ticket;
- description - more detailed information regarding what is requested or what bug was encountered;
- time estimate - estimated number of hours to complete the task (set by triagers or developers);
- time spent - time period between opening the ticket and a specific point in time;
- created timestamp;
- issue status - jira specific statuses, including Open, Closed, Awaiting Review, Patch Submitted;
- due date - optional deadline for when the ticket should get closed;
- comments - discussion around the ticket conducted by developers, end users, triagers;
- priority - it can range from low priority (minor) to critical/blocker;
- issue type - specific Jira field that specifies whether the ticket is either a bug, a feature request or a general task.

Another component of issues that we stored and proved to be crucial in our analysis is changelog histories together with their corresponding items. They represent the whole history of a ticket and can show status transitions, story points modifications, summary/description editing etc. What made them useful for our study was that we needed to see when tickets were marked closed and, as Jira does not provide this in a separate field, we looped through all changelog history items and saved when the transition to Closed or a similar status was made.

There are other fields that can be configured inside Jira, including custom fields, but we did not collect them as they would not have been helpful in conducting the analysis. However, in addition to the fields we saved, we also stored grammar correctness scores, sentiment scores, whether they have steps to reproduce, if stack traces are present and number of words in comments, summary and description. Even though all of them, apart from grammar and sentiment scores, can be computed locally, we store them because of performance issues due to the very large size of the database.

In total, we have collected 303,138 tickets spanned across 38 projects from the Apache Software Foundation: Impala, Eagle, Groovy, Lucene, Hadoop, Kafka, Apache Infrastructure project, Tika, Solr, ActiveMQ, ZooKeeper, Velocity, Tez, Storm, Stratos, CouchDB, Cassandra, Beam, Aurora, Bigtop, Camel, CarbonData, Cloudstack, Flex, Flink, Ignite, HBase, Mesos, Ambari, Cordova, Avalon, Atlas, Cactus, Flume, Felix, Geode, Ivy and Phoenix. These projects

are using a large varieties of programming languages, ranging from Java, C, Go to Python and Ruby [7]. Moreover, in terms of contributors, the projects we selected range from small teams of people such as Tika to large numbers of developers spread across the globe, such as the people working on Kafka. We came to this final number and range of tickets and projects because of a couple of reasons:

- the study needed as much diversity as possible in order to correctly analyse and validate the data;
- these are the most important, up-to-date and contributed to Apache projects by the open source world;
- Apache is one of the few companies/foundations that use Jira exclusively for their projects and it is public (they do provide a Bugzilla alternative, but it is rarely updated compared to Jira).

More specifically, we computed the following numbers for the tickets we collected:

- out of the 303,138 tickets, 236,383 tickets have been closed (i.e. marked Closed, Resolved, Done or Completed) by the time we fetched them from the Jira instance;
- 287,120 tickets are of high priority (i.e. marked as Blocker, Critical, Major or High);
- 103,397 tickets have attachments which we split in the following categories:
 - 4,082 have code attachments (e.g. Go, Java, Python);
 - 7,171 have image attachments (e.g. png, jpg/jpeg, gif);
 - 164 have video attachments (e.g. mkv, avi, mp4);
 - 13,990 have text attachments (e.g. txt, md, doc(x));
 - 918 have config attachments (e.g. json, xml, yaml);
 - 2,491 have archive attachments (e.g. zip, tar.gz, rar, 7z);
 - 280 have spreadsheet attachments (e.g. csv, xlsx);
 - 80,015 have other attachments (i.e. any file extension not pre-defined by us in one of the above categories).
- 270,907 tickets have comments;
- 39,988 tickets have steps to reproduce (i.e. sequence of steps specified by the bug reporter that would reproduce a certain bug);
- 1,942 tickets have Java stack traces - we have applied the technique described by Bettenburg et. al [4] for extracting structured data from bug reports and even though the method proved to be efficient, the total number of stack traces is so small due to the fact that many projects are actually written in languages other than Java;
- 133,689 tickets have grammar correctness scores (i.e. number of grammar mistakes inside summary, description and comments) - they have been collected using Microsoft Azure Bing Spell Check API which is one of the best grammar checking tool in the industry;

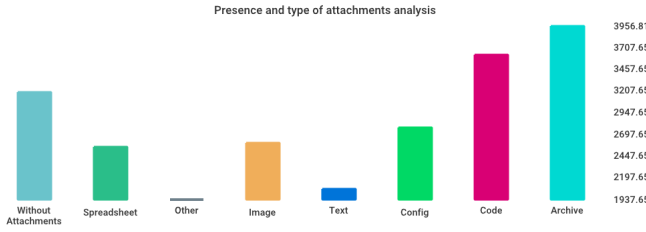


Figure 2: Application flow of the Go tool.

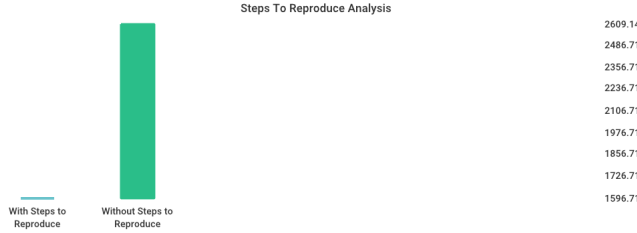


Figure 3: Application flow of the Go tool.

- 157,047 tickets have sentiment scores (i.e. positive/negative sentiment on a scale from -1 to 1 drawn from summary, description and comments) - they have been collected using Google Cloud Platform Natural Language Processing APIs.

All these numbers are also stored inside the database in a different bucket for easier access in the future. Moreover, the statistical tests we ran on the tickets are saved as well in the database.

5. CORRELATIONS

5.1 Attachments

A Welch T Test was computed to test the difference between tickets with attachments and those without attachments. The analysis revealed significant difference with $p < 0.01$ with an average of $X <change this>$ for tickets with attachments and an average of $Y <change this>$ for tickets without attachments.

5.2 Steps to Reproduce

A Welch T Test was computed to test the difference between tickets with attachments and those without attachments. The analysis revealed significant difference with $p < 0.01$ with an average of $X <change this>$ for tickets with attachments and an average of $Y <change this>$ for tickets without attachments.

5.3 Stack Traces

A Welch T Test was computed to test the difference between tickets with attachments and those without attachments. The analysis revealed significant difference with $p < 0.01$ with an average of $X <change this>$ for tickets with attachments and an average of $Y <change this>$ for tickets without attachments.



Figure 4: Application flow of the Go tool.

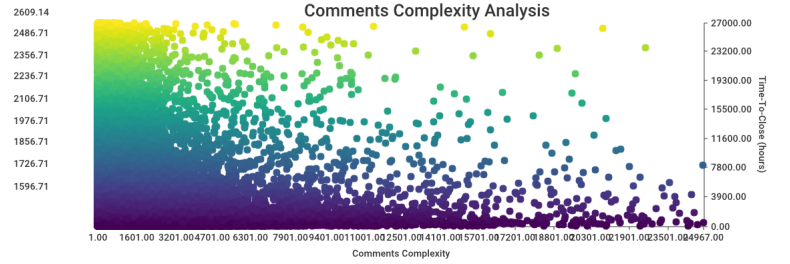


Figure 5: Application flow of the Go tool.

5.4 Comments

A correlational analysis revealed a weak positive/negative (one of them, depending on context) relationship with $r = \dots$, $p = \dots$.

5.5 Summary and Description

A correlational analysis revealed a weak positive/negative (one of them, depending on context) relationship with $r = \dots$, $p = \dots$.

5.6 Grammar Correctness

A correlational analysis revealed a weak positive/negative (one of them, depending on context) relationship with $r = \dots$, $p = \dots$.

5.7 Sentiment Scores

A correlational analysis revealed a weak positive/negative (one of them, depending on context) relationship with $r = \dots$, $p = \dots$.

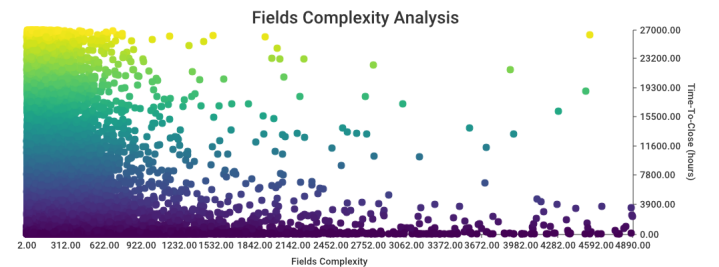


Figure 6: Application flow of the Go tool.

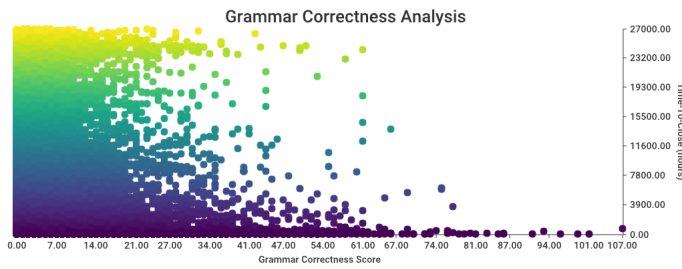


Figure 7: Application flow of the Go tool.

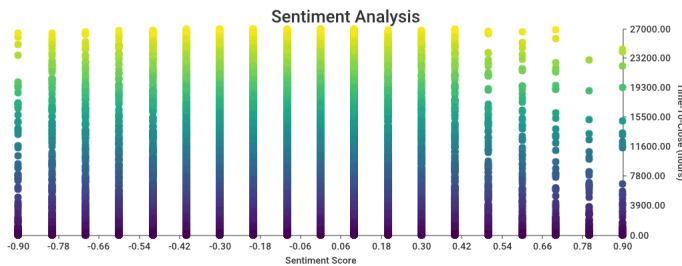


Figure 8: Application flow of the Go tool.

6. FUTURE WORK

7. CONCLUSIONS

Acknowledgments. Tim, my parents, Corina.

8. REFERENCES

- [1] Atlassian. Jira. <https://www.atlassian.com/software/jira>. [Online; accessed 02-April-2018].
- [2] N. Bettenburg, S. Just, A. Schröter, C. Weiss, R. Premraj, and T. Zimmermann. What makes a good bug report? pages 308–318, 2008.
- [3] N. Bettenburg, R. Premraj, T. Zimmermann, and S. Kim. Duplicate bug reports considered harmful? really? pages 337–345, 2008.
- [4] N. Bettenburg, R. Premraj, T. Zimmermann, and S. Kim. Extracting structural information from bug reports. pages 27–30, 2008.
- [5] N. Bettenburg, S. W. Thomas, and A. E. Hassan. Using fuzzy code search to link code fragments in discussions to source code. pages 319–328, 2012.
- [6] M. Fischer, M. Pinzger, and H. Gall. Analyzing and relating bug report data for feature tracking. 3:90, 2003.
- [7] A. S. Foundation. Apache Projects. <https://projects.apache.org/projects.html?language>. [Online; accessed 12-March-2018].
- [8] Google. GCP Natural Language Processing API. <https://cloud.google.com/natural-language>. [Online; accessed 05-March-2018].
- [9] P. Hooimeijer and W. Weimer. Modeling bug report quality. pages 34–43, 2007.
- [10] B. Johnson. Bolt. <https://github.com/boltdb/bolt>. [Online; accessed 19-February-2018].
- [11] S. Just, R. Premraj, and T. Zimmermann. Towards the next generation of bug tracking systems. pages 82–85, 2008.
- [12] M. Korkala and F. Maurer. Waste identification as the means for improving communication in globally distributed agile software development. *The Journal of Systems and Software*, 95:122–140, 2014.
- [13] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals. Predicting the severity of a reported bug. pages 1–10, 2010.
- [14] T. D. LaToza, G. Venolia, and R. DeLine. Maintaining mental models: a study of developer work habits. In *Proceedings of the 28th international conference on Software engineering*, pages 492–501. ACM, 2006.
- [15] Microsoft. Bing Spell Check API. <https://azure.microsoft.com/services/cognitive-services/spell-check>. [Online; accessed 05-March-2018].
- [16] Mozilla. Bugzilla. <https://www.bugzilla.org/>. [Online; accessed 02-April-2018].
- [17] S. Rastkar, G. C. Murphy, and G. Murray. Summarizing software artifacts: a case study of bug reports. pages 505–514, 2010.
- [18] A. Schroter, A. Schröter, N. Bettenburg, and R. Premraj. Do stack traces help developers fix bugs? pages 118–121, 2010.
- [19] C. Spearman. The proof and measurement of association between two things. *The American journal of psychology*, 15(1):72–101, 1904.
- [20] B. L. Welch. The generalization of student’s problem when several different population variances are involved. *Biometrika*, 34(1/2):28–35, 1947.
- [21] T. Zimmermann, R. Premraj, J. Sillito, and S. Breu. Improving bug tracking systems. In *Software Engineering-Companion Volume, 2009. ICSE-Companion 2009. 31st International Conference on*, pages 247–250. IEEE, 2009.