



# Study the correlation between the readme file of GitHub projects and their popularity<sup>☆</sup>

Tianlei Wang<sup>a</sup>, Shaowei Wang<sup>a,\*</sup>, Tse-Hsun (Peter) Chen<sup>b</sup>

<sup>a</sup> Department of Computer Science, University of Manitoba, Canada

<sup>b</sup> Department of Computer Science and Software Engineering, Concordia University, Canada

## ARTICLE INFO

### Article history:

Received 9 November 2022

Received in revised form 18 April 2023

Accepted 18 July 2023

Available online 22 July 2023

### Keywords:

Documentation

GitHub

Readme file

Random forest

## ABSTRACT

A readme file plays an important role in a GitHub repository to provide a starting point for developers to reuse and make contributions. A good readme could provide sufficient information for users to learn and start a GitHub repository and might be correlated to the popularity of a repository. Given the importance of the role that a readme file plays, we aim to study to understand the correlation between the readme file of GitHub repositories and their popularity. We analyze readme files of 5,000 GitHub repositories across more than 20 languages. We study the relationship between readme file related factors and the popularity of GitHub repositories. We observe that: (1) Most of the studied readme file related factors (e.g., the number of lists, the number and frequency of updates on the readme file) are statistically significantly different between popular and non-popular repositories with non-negligible effect size. (2) After controlling repository-specific factors (e.g., repository topics and license information), the number of lists and the frequency of updates are the most significantly important factors that discriminate between popular and non-popular repositories. (3) The most of updates were made to update references in popular repositories, while in non-popular repositories most updates are for the content of how to use the repository.

*Editor's note: Open Science material was validated by the Journal of Systems and Software Open Science Board.*

© 2023 Elsevier Inc. All rights reserved.

## 1. Introduction

GitHub, the largest open-source projects hosting site in the world, provides a platform for developers and organizations to share and host their open-source projects. GitHub hosts more than 200 million repositories, as of March, 2022 (GitHub, 2022c). Millions of developers and organizations contribute to these repositories by adding new features and resolving bugs (Bao et al., 2019; McDonald and Goggins, 2013; Jiang et al., 2017). Developers also learn from the hosted projects on GitHub and utilize them in their projects (Gharehyazie et al., 2019, 2017).

A readme file plays an important role in a repository to provide a starting point for developers to understand and learn the repository. As the official documentation supported by GitHub, a readme file typically provides the basic information related to the repository, such as what the project does, how users can get started with the project, etc (GitHub, 2022a). The information provided in the readme file largely affects the contribution

and the usage of a repository. A well-organized and maintained readme could provide sufficient information for users and reduce the learning curve of getting started on the project for users (Steinmacher et al., 2014), which probably correlates to the popularity of the project. For instance, a repository with a well-maintained readme file, in which all information updated to date, may be easy for a user to get started on it, compared with a repository that has a significant amount of out-of-dated information in the readme file, as a result, attracts more developers to use it and contribute to it. Therefore, we aim to understand the correlation between various factors related to readme file of a GitHub repository and the popularity of the repository, e.g., whether the organization and update activities of a readme file correlate to the popularity of a repository. Notice that our aim is not to uncover causal relationships between various factors and the popularity of GitHub repositories. Instead, we compare the differences between popular and non-popular GitHub repositories, aiming to help highlight good practices in improving readme files.

In this study, we performed an analysis of the readme files of 5,000 GitHub repositories across different domains and languages. We study the relationship between readme file related factors and the popularity of GitHub repositories to understand

<sup>☆</sup> Editor: Gabriele Bavota.

\* Corresponding author.

E-mail addresses: [wangt316@myumanitoba.ca](mailto:wangt316@myumanitoba.ca) (T. Wang), [shaowei.wang@umanitoba.ca](mailto:shaowei.wang@umanitoba.ca) (S. Wang), [peterc@encs.concordia.ca](mailto:peterc@encs.concordia.ca) (T.-H. Chen).

whether the factors related to readme files are different between popular and non-popular GitHub repositories, and which factors are most important for distinguishing them. For this purpose, we structure our inquiry using 18 factors along three dimensions: the content presented and organized in the readme file (e.g., the number of lists, images, and links), updates on the readme file (e.g., the number and frequency of updates), and repository meta information (e.g., topics, age, license, and size). More specifically, we structure our study along the following three research questions:

- **RQ1: Is there a relationship between the studied readme file related factors and the popularity of GitHub repositories?**

In RQ1, we compare the readme file related factors between popular and non-popular repositories and examine if the differences between those two groups are statistically significant. We observe that all studied readme file related factors (e.g., number of links, images, and updates frequency) except heading content are statistically significantly different between popular and non-popular repositories with non-negligible effect size. Popular and non-popular repositories share similar headings (e.g., usage and installation).

- **RQ2: What are the most important readme file related factors that differentiate popular and unpopular GitHub repositories?**

In RQ2, we build a classification model to investigate which readme file related factors are the most important for discriminating between popular and non-popular repositories by controlling factors related to the repository. We observe that after controlling factors related to the repository (e.g., repository topics and license information), the number of lists, links, and the frequency of updates are statistically significantly important factors that discriminate between popular and non-popular repositories. The frequency/number of readme updates and the number of lists and links positively correlate with the likelihood of a repository being popular.

- **RQ3: What content is updated in readme files between popular and non-popular GitHub repositories?**

For RQ3, we perform qualitative analysis to understand what content is updated in readme file of popular and non-popular repositories. We observe that the largest portion of the updates (47%) were performed to update references in popular repositories, while the largest portion of updates were performed to update the content of “how” in non-popular ones (34%). A remarkable portion of updates on readme files were performed to improve the presentation and update badges for both popular and non-popular repositories.

In summary, the quality of the readme file (e.g., proper organization of content and ensuring content is up to date) is positively correlated with the popularity of repositories after controlling for repository-specific factors. Developers should pay attention to improving the readme file, e.g., improving the organization of content and ensuring the content is up to date. Future research is encouraged to develop approach to help maintain readme automatically (e.g., detecting and updating broken/obsolete references) and generate readme file templates automatically. To enrich future research on this direct, we make our replication publicly available [https://github.com/TianleiWang0414/WIP2022-GitHub\\_Readme-code/](https://github.com/TianleiWang0414/WIP2022-GitHub_Readme-code/).

**Paper Organization.** Section 2 introduces the background information about the GitHub readme file and related work. Section 3 describes our data collection process and studied factors.

Section 4 presents the result of our research questions. Section 5 describes the discussion on update activities of readme files, implications of our findings, and the threats to the validity of our observation. Finally, Section 6 concludes the paper.

## 2. Background & related work

### 2.1. Readme file on GitHub

Readme files are text files, typically written in markdown, that provide an overview of information to allow an individual to understand how to run, when to use, and what tools to apply for a given project. Usually, readme files are the first documentation/file of a project that is visited. Readme files usually follow a strict naming convention, as the name suggests, they are often written as “readme.md”. Suggested by GitHub official documentation, readme files typically contain the following information: (1) What the project does; (2) Why the project is useful; (3) How users can get started with the project; (4) Where users can get help with your project; (5) Who maintains and contributes to the project (GitHub, 2022a). Readme files are part of the GitHub repositories and they can be updated accordingly as the projects evolve. Fig. 1 shows an example readme file of a GitHub repository.<sup>1</sup> Various pieces of information are labeled in the example, such as heading, image, badge, and link.

### 2.2. Related work

#### 2.2.1. Studying and leveraging the content presented in readme files

Readme files are considered an important source for understanding repositories not only by developers but also by researchers. Therefore, there has been a number of studies that leverage readme files to conduct software engineering tasks (e.g., project recommendation and project categorization). Zhang et al. (2017) proposed an approach that leverages the content of readme files together with their sources (e.g., source code), to detect similar repositories on GitHub. They assume that similar repositories have similar readme files in terms of content. Similarly, Koskela et al. (2018) leveraged the content presented in readme files together with the tags of repositories to recommend GitHub repositories for users based on their profiles. Sharma et al. (2017) categorized GitHub projects based on the content in readme files. Coelho et al. (2018) leverage the information presented in readme files to identify unmaintained GitHub projects. Trockman et al. (2018) conducted an empirical study to understand the repository badges in the npm Ecosystem shown on the readme files. Ikeda et al. (2019) investigate the content of readme file for JavaScript packages and its relationship with the type of a package. Different from prior studies that focus on leveraging readme file as a source to perform other tasks, our study focuses on studying the readme file itself and investigating its association with the popularity of GitHub repositories. Hassan and Wang (2017) proposed technique to automatically extract software build commands from software readme files and Wiki pages for software building.

Several studies have been done to investigate the content of readme files of GitHub repositories. Liu et al. (2022) study the patterns of readme files and determine the degree to which readme files are aligned with the official guidelines. They found that the majority of readme files do not align with the GitHub guidelines and repositories whose readme files follow the GitHub guidelines tend to receive more stars. Prana et al. (2019) studied the content of 393 GitHub readme files. Their study shows that the information about “what the project does” and “how users

<sup>1</sup> <https://github.com/alibaba/Sentinel>

## Sentinel: The Sentinel of Your Microservices

Sentinel CI passing codecov 45% Maven Central v1.8.4 license Apache 2 chat on gitter — badge

### Introduction heading

As distributed systems become increasingly popular, the reliability between services is becoming more important than ever before. Sentinel takes "flow" as breakthrough point, and works on multiple fields including **flow control**, **traffic shaping**, **circuit breaking** and **system adaptive protection**, to guarantee reliability and resilience for microservices.

Sentinel has the following features:

- **Rich applicable scenarios:** Sentinel has been wildly used in Alibaba, and has covered almost all the core-scenarios in Double-11 (11.11) Shopping Festivals in the past 10 years, such as "Second Kill" which needs to limit burst flow traffic to meet the system capacity, message peak clipping and valley fills, circuit breaking for unreliable downstream services, cluster flow control, etc.
- **Real-time monitoring:** Sentinel also provides real-time monitoring ability. You can see the runtime information of a single machine in real-time, and the aggregated runtime info of a cluster with less than 500 nodes.

### Documentation

See the [Sentinel](#) for the document website link

See the [中文文档](#) for document in Chinese.

### 2. Define Resource

Wrap your code snippet via Sentinel API: `SphU.entry(resourceName)`. In below example, it is `System.out.println("hello world");`; code block

```
try (Entry entry = SphU.entry("HelloWorld")) {
    // Your business logic here.
    System.out.println("hello world");
} catch (BlockException e) {
    // Handle rejected request.
    e.printStackTrace();
}
// try-with-resources auto exit
```

### Who is using

These are only part of the companies using Sentinel, for reference only. If you are using Sentinel, please [add your company here](#) to tell us your scenario to make Sentinel better :) image



Fig. 1. Example readme file of a GitHub repository.

can get started with the project" of a repository is common while information on the purpose and status is rare. They developed a classifier to predict categories of sections in the readme files with an F1 score of 0.746. Ikeda et al. (2019) studied the readme file of JavaScript packages, and they show that information related to usage, installation, and license are commonly contained by such readme files. Treude et al. (2020) recruited technical editors to assess various software documentation, including readme files for R GitHub Projects and Stack Overflow threads in terms of ten dimensions such as readability, cohesion, and structure. They shed light on challenges for software documentation and directions for future studies. Different from these studies that focus on studying the content of readme files, our work aims to investigate the relationship between readme file related features of GitHub projects and their popularity.

#### 2.2.2. Studying popularity of GitHub repositories

Many studies have analyzed software code repositories to understand what makes a code repository popular. Weber and Luo (2014) proposed a classifier using 38 features in various dimensions (e.g., code, project) to classify popular and non-popular Python GitHub projects. Zhu et al. (2014) studied the patterns of folders used by 140k Github projects and their relationship with project popularity. They find that the standard folders, such as documents, testing, and examples, are not only among the most frequently used, but their presence in a project is associated with increased chances that a project's code will be forked. Aggarwal et al. (2014) investigated the correlation between the documentation change of 90 GitHub projects and their popularity. They observed that the consistent popularity of projects attracts consistent work on documentation. Different from their study,

we not only investigate the features related to the revision of readme files, but also other features in terms of structure and content of the readme files. [Borges et al. \(2016\)](#) identified four patterns of popularity growth of projects on GitHub in terms of stars. They observed that the main factors that impact the number of stars of GitHub projects are programming language and the application domain. We consider programming language and application domain (we consider tags of a repository) as controlling factors in this study. Different from previous studies, we focus on examining the relationship between features related to readme files and Github projects' popularity and controlling several repository related factors, such as programming language, license, topics, size, and age.

The most related work is a study by [Fan et al. \(2021\)](#), which studied the correlation between features from various dimensions (i.e., code, reproducibility, and documentation) of AI repositories and their popularity. Their results show that popular and non-popular AI repositories are significantly different in various features, such as the number of links, images in the README file, and license. Different from Fan et al.'s study which focuses on investigating academic AI repositories, our study focus on more general GitHub repositories.

### 3. Data preparation

#### 3.1. Repository data collection

To have enough readme files to answer our RQs, we randomly sampled 5,000 GitHub repositories from Google big query as of January, 2022 ([Google, 2020](#)). We retrieved the owner information as well as the repository URL from the Google big query. We then use GitHub API ([GitHub, 2022b](#)) to retrieve the README.md files in the main/master branch for the sampled 5,000 repositories. The studied repositories host projects that span more than 20 programming languages including popular ones, such as Python, Java, JavaScript, Go, and C++.

To ensure the quality of our studied repositories, we future filter repositories based on certain criteria. We first removed the repositories that do not have a readme file (i.e., "README.MD"). We filtered out 132 such repositories. We removed the repositories that have default content in the readme file. GitHub initializes the content of a readme file with the repository name if the creator does not modify the readme file. The reason why we removed such repositories is that they do not contain any useful information for our RQs (only contain the name of the repositories). We removed the repositories whose readme file is not fully written in English. We removed such repositories since it requires translation of the whole file to extract useful information, since translation may introduce bias. We used a library called *langdetect*<sup>2</sup> to detect the language written in the readme files. We further removed 614 repositories during the above process. GitHub hosts repositories (academic/tutorial repositories) other than software repositories that could introduce bias ([Munaiah et al., 2017](#); [Liu et al., 2022](#)). We removed the repositories that have less than 50% of the files that are not written in programming languages and have only one contributor. Last, we removed the repositories that have less than one year from their creation date to our collected date. We did so to ensure the repositories have enough time to attract attention from the community and obtain starts. We ended up with 1,566 repositories.

**Table 1**

Statistics of popular and non-popular repositories. The value in each cell is ordered as mean/median/min/max.

Factor	Popular	Non-popular
Star	13,398.2/4374/677/339,873	0/0/0/0
Fork	3029.7/939/25/86,266	0.25/0/0/6
Watch	516.8/202/7/8,472	2.4/2/1/57
Size (in KB)	137,704.4/14764.0/0/22,309,053	18,339.4/603/0/1,222,458
Age (in day)	2280.8/2137/770/4,808	2446.8/2420/77/3832

#### 3.2. Data labeling

To conduct our analysis, we need to categorize repositories as popular or unpopular. For this purpose, we collected repositories' metadata including stars, watches, number of pull requests, number of contributors, and forks via Github API to find what can be represented for popularity. With this information, we applied Spearman's rank correlation to them. [Fig. 3](#) shows that there is a high correlation between forks, watches, and stars. Thus, we believe it is appropriate to select stars as the proxy for popularity following prior studies ([Han et al., 2019](#); [Borges et al., 2016](#)). We find that the number of stars is highly skewed and only a small portion of repositories received many stars (see [Fig. 2](#)). We label the repositories that have a star of 0 as non-popular and the ones that have a star larger than 100 as popular. Note that we use a clear-cut since we would like to mitigate the bias from the repositories with a very similar number of starts while belonging to two different categories and provide a clear boundary between popular and non-popular repositories. For instance, if we considered a threshold of 50 as the boundary to categorize the repositories with at least 50 stars as popular repositories and the ones with less than 50 as non-popular repositories, then a repository with 49 stars and a repository with 51 would be considered as different classes although their stars are very close. Nevertheless, such a clear-cut threshold probably would introduce bias to our study and cause a threat to the validity of our study. We would discuss it in more detail in Section 5.3. We collected 810 non-popular repositories and 756 popular ones. [Table 1](#) shows the information of two groups of repositories.

#### 3.3. Factors collection

In RQ1, we compare the factors related to the readme file between popular and non-popular repositories (see Section 4.1 for details). In RQ2, we build a classification model to understand the most important readme file related factors discriminating between popular and non-popular repositories by controlling factors that are related to their corresponding repositories.

In this subsection, we discuss the studied factors and how we collect them. The quality of a readme file probably affect the popularity of the corresponding repository. A well-maintained and well-organized could provide sufficient information for users and reduce the learning curve of getting started on the project for users. We collect the factors related to readme file's quality from two dimensions: readme content and readme update. Readme content and update consider the quality of the readme file from static and dynamic perspectives, respectively. The content dimension reflects if a readme file is well-organized, and the update dimension reflects if a readme file is well-maintained. In terms of content, we consider the factors that are captured at the moment of data collection. For content dimension, we first consider the factors that characterize the presentation of a readme file (i.e., *R\_NumCodeBlock*, *R\_NumImage*, *R\_NumIndent*, and *R\_NumList*). We expect that a readme file with a good presentation provides a good impression to visitors and correlates with the popularity of the repository. GitHub users usually use

<sup>2</sup> <https://pypi.org/project/langdetect/>



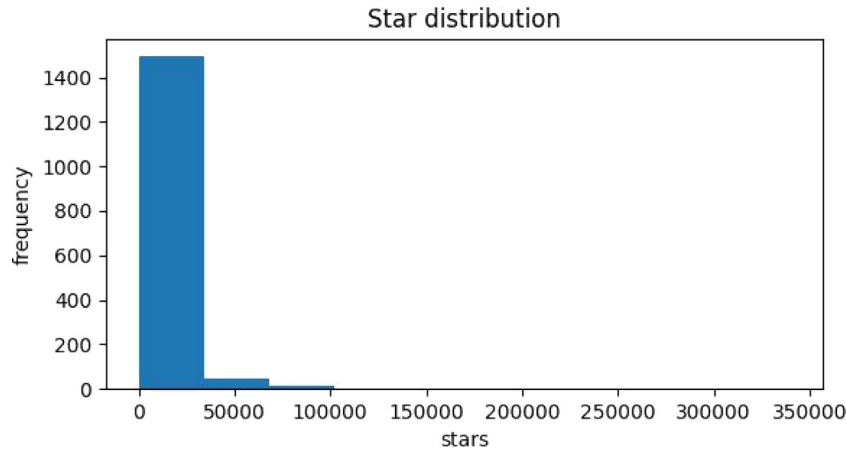


Fig. 2. Distribution of repositories with different numbers of stars.

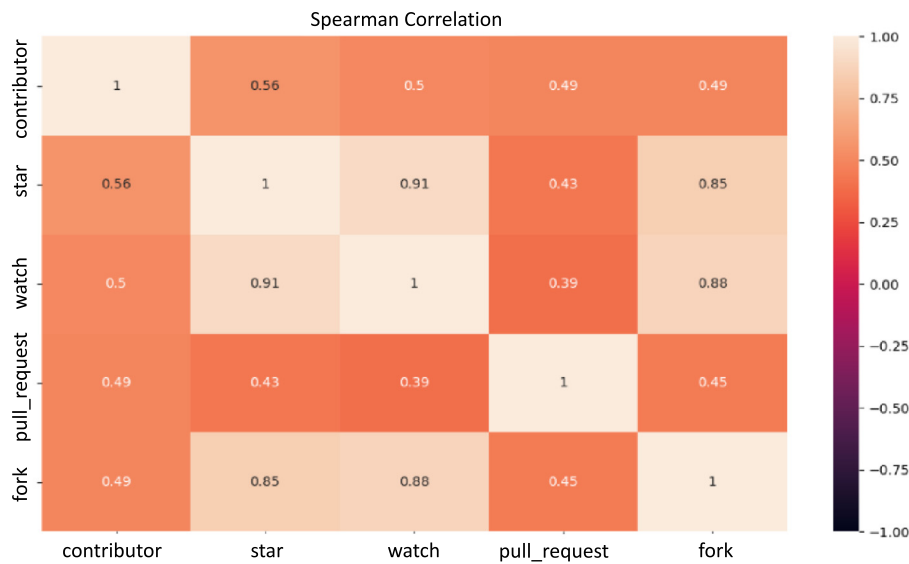


Fig. 3. Correlation matrix between the number of forks, pull requests, watches, stars, contributors. We observe that fork, watch, and star have a high correlation.

markdown for organizing their readme files. Therefore, we consider the formatting types that are commonly used in markdown. For instance, we collect the number of code blocks, indents, lists, and images. Kimble (1992) recommends using lists for clear writing. Tian et al. (2015) show the images displayed in a mobile app have a positive correlation with the popularity of the mobile app. Second, we consider the factors related to the information presented in the content (e.g., installation instruction, license, support link) to help visitors to understand the repository and get started on it. We collect the number of links ( $R\_NumLink$ ), badges ( $R\_NumBadge$ ), words presented in the headings of a readme file ( $R\_heading$ ), and license ( $P\_License$ ). We expect that a readme file with sufficient information could help visitors get started on the repository and improve its popularity. For instance, license information clearly indicates to visitors how to use a repository (Fan et al., 2021; Laurent, 2004). Links are usually used to point to references for external supporting materials. To collect those factors, for each readme file, we collected the heading information, code blocks, indents, images, lists, and badges. We used the regular expression shown in Table 3 with the Python library *re*. For heading information, we extracted each word in the heading, and applied lemmatization and stemming using Python libraries *spacy* and *nltk* for preprocessing. For other factors, we collected the number of occurrences in each readme file. Note

that since a badge is also an image, we further check if a link contains “<https://img.shields.io>” to determine if the link represents a badge.

For the update dimension, we consider the history of a readme file and examine how it was updated. In this study, we consider the number of updates that a readme file receives ( $R\_NumUpdate$ ) and update frequency ( $R\_UpdateInterval$ ). Aggarwal et al. (2014) found that frequent updates on the documentation of a project have a correlation with its popularity. We expect that frequent updates help maintain a repository stay up to date, therefore improving the popularity of the repository. Note that we count the initialization of a readme file when calculating  $R\_UpdateInterval$  in order to avoid zero in the denominator. To collect the update information, we use Github API (GitHub, 2022b) to collect all commits that involve changes on readme files for each studied repository and extracted the message and the time stamps for each commit. We then calculate each factor accordingly.

There are some other confounders that potentially impact the popularity of GitHub repositories. Therefore, we also consider factors related to the repository's meta information, i.e., the age ( $P\_Age$ ), the major programming language ( $P\_ProgramLanguage$ ), the size ( $P\_size$ ), license ( $P\_License$ ), the popularity ( $P\_Max/Min/Avg/Median\_Tag$ ) as controlling factors when building the model in RQ2, since they are related to the project itself and cannot be

**Table 2**

Studied factors that are related to three dimensions readme content, readme update, and repository meta.

	Factor name	Explanation (data type)	Rationale
Readme content	<i>R_NumCodeBlock</i>	The number of code blocks in the readme file (numeric).	Good presentation of readme file of a GitHub Repositories probably could provide a good impression to visitors and improve its popularity. For instance, the list is recommended for clear writing (Kimble, 1992; Fan et al., 2021).
	<i>R_NumIndent</i>	The number of indents in the readme file (numeric).	
	<i>R_NumList</i>	The number of lists in the readme file (numeric).	
	<i>R_NumLink</i>	The number of links in the readme file (numeric).	Providing links and images to further documentation and support channels could help visitors to get started and improves the popularity of the repository.
	<i>R_NumImage</i>	The number of images in the readme file (numeric).	
	<i>R_NumBadge</i>	The number of badges in the readme file (numeric).	The badge indicates the quality and properties of a repository and may have a correlation with the popularity of the repository.
	<i>R_heading</i>	The words appear in headings of the readme file (text).	
Readme update	<i>R_UpdateInterval</i>	The average update interval between updates (calculated as <i>total elapsed days since creation/number of updates</i> ) (numeric).	Frequently update on documentation of a project have positive correlation with its popularity (Aggarwal et al., 2014).
	<i>R_NumUpdate</i>	Number of updates for the readme file (numeric).	
	<i>R_SinceLastUpdate</i>	Number of days since last update in the readme file (numeric).	The time of being inactive of readme file of a project may have negative correlation with its popularity.
Repository meta	<i>P_Age</i>	Age of a repository since its creation in days (numeric).	An older repository tends to get more attention from the community and probably has more stars, forks, etc.
	<i>P_size</i>	The size of a repository in KB. (numeric).	Larger size of a repository might indicate more features or better functionality, which might be correlated with its popularity (Tian et al., 2015).
	<i>P_Max/Min/Avg/Median_Tag</i>	The max, min, average, and median value of tag popularity scores for a repository. Tag popularity score is calculated as the number of repositories in GitHub that have a tag (numeric).	The tags of a repository indicate its application domain and related topics. The repositories with popular tags may get more attention from the community and probably have more stars, forks, etc (Zhou et al., 2020a; Borges et al., 2016).
	<i>P_License</i>	Type of license the repository has (categorical).	Users can clearly understand the license of a repository and facilitates its popularity (Fan et al., 2021; Laurent, 2004).
	<i>P_ProgramLanguage</i>	The majority of programming languages used in the repository (categorical).	Prior study by Borges et al. (2016) observed that the one of main factors that impact the number of stars of GitHub projects is the programming language.

modified when editing readme file by the repository owner. Note that for popularity, we actually have four factors. We calculate the factors based on the meta information we collected for each repository. We present the studied factors in Table 2.

#### 4. Results of the research questions

4.1. RQ1: Is there a relationship between the studied readme file related factors and the popularity of GitHub repositories?

**Motivation** The readme file of a GitHub repository provides the first expression of the repository to visitors and plays an important role for visitors to get to know the repository. A readme file

**Table 3**

Regular expressions that are used to extract different pieces of information from readme files.

Factor	Regular expression
<i>R_heading</i>	<code>#+ .+\\s</code>
<i>R_NumCodeBlock</i>	<code>`{3}([Ss]?)`{3}.*</code>
<i>R_NumIndent</i>	<code>^(?:[ ] [4])\\t.*(R\$)+</code>
<i>R_NumLink</i>	<code>\\[^(\\] r\\n +\\ \\((\\^\\) r\\n +\\)</code>
<i>R_NumImage</i>	<code>!\\[^(\\] r\\n +\\ \\((\\^\\) r\\n +\\)</code>
<i>R_NumList</i>	<code>(`{3})(- + * 0-9 .)+ \\[Ss]`{3}?r\\n\\r\\n</code>
<i>R_NumBadge</i>	<code>!\\[^(\\] r\\n +\\ \\((\\^\\) r\\n +\\)</code>

may have a correlation with the popularity of the GitHub repository. In this RQ, we compare the readme file related features

between popular and non-popular GitHub repositories. By understanding this, GitHub project owners may leverage our findings to further improve their projects.

**Approach** For each of the readme file related numeric features (see details in Section 3.3), we apply the Wilcoxon rank-sum test (Mann and Whitney, 1947) to examine whether the differences in these features between popular and non-popular GitHub repositories are statistically different or not. We also compute Cliff's *delta* to examine the effect size of the difference between the two groups (Cliff, 1993). Note that we interpret  $|d|$  as follows: a  $|d|$  less than 0.147, between 0.147 and 0.33, between 0.33 and 0.474, and larger than 0.474 is considered as a negligible (N), small (S), medium (M), and large (L) effect size following prior study (Romano et al., 2006). The reason we select the Wilcoxon rank-sum test and cliff's *delta* since they are non-parametric and do not require any assumption on the distribution of two samples.

To understand if there exist differences in the content between the two groups, we compute frequent words that appear in the headings of readme files. We choose to analyze words appearing in headings, since the headings usually summarize the major topics presented in a readme file. We also applied stemming to reduce the size of vocabulary and convert similar words into their original form (e.g., we convert “install”, “installed”, “installation” into “install”). After the pre-processing, we compare the top 15 frequent words mined from the headings of readme files of the two groups.

**Results: All studied numeric readme file-related features have statistically significant differences between popular and non-popular repositories with non-negligible effect size.** Table 4 shows the comparison of the studied numeric features between popular and non-popular GitHub repositories, including the mean, median, minimum, and maximum values of each factor. We observe that the two groups of repositories have statistically significant differences ( $p$ -value  $< 0.05$ ) with a non-negligible effect size for all the studied features, in which *R\_NumLink*, *R\_NumUpdate*, *R\_NumUpdate*, *R\_Length*, and *R\_SinceLastUpdate* have large effect size. Therefore, in general, popular and non-popular GitHub repositories have significant differences in terms of the presented structure, length, and update activities of readme files.

For instance, popular Github repositories (12,396.1 characters) have a four times longer readme file than non-popular repositories (3,115.8 characters) on average. A longer readme file usually suggests more detailed information about the repositories. Meanwhile, popular repositories receive more updates (mean value is 17.8) compared with non-popular repositories (mean value is 6.7). Popular repositories have a mean value of 6.2 code blocks, 10 indents, and 3.4 lists, which are significantly more than non-popular repositories with 2.3 blocks, 5.2 incidents, and 0.02 lists on average. That is said, popular repositories are usually better presented than non-popular ones. For instance, the repository *AX-IOX*<sup>3</sup> is a promise-based HTTP client for the browser and node.js. It has 91k stars (as of the date of data collected), which is a very popular GitHub repository. It contains a table of contents to provide the index of the readme file. The readme file describes the features, detailed installing instructions, usage examples, lists of resources, which is very long.

**Popular and non-popular repositories share similar headings when examining the top 15 most frequent words appearing in headings of readme files.** Table 5 presents the top 15 frequent words occurring in the headings of readme files. We notice that 11 out of the top 15 words are overlapped between these two groups. Such observation indicates that the readme

files of popular and non-popular repositories share similar sections. For instance, sections related to license (i.e., “license”), how to use/start the repository (i.e., “install”, “start”, “build”, “example”, “usage”), contribution (“contribute”), support (“support”) commonly shared in the readme files of both groups. More importantly, 4 out of the top 5 words between the two groups are overlapped, i.e., “license”, “install”, “contribute”, and “start”. For instance, tensorflow<sup>4</sup> hosts the repository for tensorflow the readme file contains information about installation, license, contribution guidelines, patching guidelines, Continuous build status, resources, and courses. In other words, both popular and non-popular introduce how to use the repository, license, and contributor with top priority. This finding echoes the previous study by Prana et al. (2019), in which they observed that 88.5% readme files have a section to introduce how to use the repository, and 52.9% of the readme files have information related to contributor and license.

**Table of contents appears more frequently in readme file of popular repositories than non-popular repositories.** We also notice that the word “content” appears frequently in popular repositories, while not in non-popular ones. We manually examined several readme files containing “content” and find such readme files usually have a heading “table of contents” to indicate the content included in the readme files. We find that it appears in popular repositories (71 out of 756 popular repositories), while only 11 out of 810 in non-popular repositories.

*All studied readme file related features except heading content are statistically significantly different between popular and non-popular repositories with non-negligible effect size. Popular and non-popular repositories share similar headings. While the table of contents appears more frequently in the readme file of popular repositories than non-popular repositories.*

4.2. RQ2: What are the most important readme file related factors that differentiate popular and unpopular GitHub repositories?

**Motivation** Among the studied factors, some factors may be more important than others. In this research question, we investigate the most important factors in differentiating two groups of repositories. Our findings probably can provide insights to repository owners, helping them improve their readme files.

**Approach** To compute the importance of the studied factors, we build a random forest classification model (Cutler et al., 2012) to differentiate popular and non-popular GitHub repositories. The factors that contribute more explanation power are considered to be more important. Below, we elaborate on each step of the model construction, model validation, and variable importance analysis.

**Correlation & redundant factor analysis.** Before constructing the model, we perform correlation analysis to remove correlated factors. Highly correlated factors can cause multicollinearity problems in our model (Farrar and Glauber, 1967), and turned out to bias the model interpretation. Thus, we perform a correlation analysis to remove highly correlated factors using a variable clustering analysis technique by following prior studies (Rajbahadur et al., 2019; Wang et al., 2018; Fan et al., 2021). First, we use the Spearman rank correlation test to measure the correlation between factors and remove highly-correlated factors (using a cut-off value of 0.7). For each group of the highly-correlated factors, we keep one factor as the representative for the group and then performed a redundancy analysis to remove redundant

<sup>3</sup> <https://github.com/axios/axios>

<sup>4</sup> <https://github.com/tensorflow/tensorflow>

**Table 4**  
P-value and cliff's |D| for the readme file related features.

Feature	Mean/Median/Min/Max		P-value & Cliff's  d
	Popular	Non-popular	
<i>R_NumCodeBlock</i>	6.2/2/0/349	2.3/0/0/73	<0.05 & 0.26 (S)
<i>R_NumIndent</i>	10.0/2/0/366	5.2/0/0/283	<0.05 & 0.15 (S)
<i>R_NumImage</i>	4.3/2/0/171	0.8/0/0/14	<0.05 & 0.41 (M)
<i>R_NumLink</i>	43.5/16/0/1153	5.3/1/0/266	<0.05 & 0.67(L)
<i>R_NumList</i>	3.4/0/0/171	0.02/0/0/5	<0.05 & 0.39 (M)
<i>R_NumBadge</i>	1.1/0/0/53	0.14/0/0/8	<0.05 & 0.26 (S)
<i>R_Length</i>	12,396.1/6,459/14/178,059	3,115.8/1,403/11/60,947	<0.05 & 0.63 (L)
<i>R_UpdateInterval</i>	254.7/128/34/4,599	665.9/336/6/4,535	<0.05 & 0.9 (L)
<i>R_SinceLastUpdate</i>	693.4/433/82/3,958	2,401.6/2,471/171/3,714	<0.05 & 0.91(L)
<i>R_NumUpdate</i>	17.8/20/1/30	6.7/3/1/28	<0.05 & 0.686 (L)

**Table 5**

Top frequently appearing words in the headings of readme files in popular and non-popular repositories. The top 5 words for popular and non-popular repositories are highlighted.

Frequent word	Popular		Non-popular	
	Frequency	Rank	Frequency	Rank
license	0.313	<b>1</b>	0.15	<b>2</b>
install	0.29	<b>2</b>	0.19	<b>1</b>
contribute	0.274	<b>3</b>	0.11	<b>4</b>
start	0.181	<b>4</b>	0.095	<b>5</b>
document	0.165	<b>5</b>	0.055	15
content	0.161	6		
build	0.156	7	0.087	6
support	0.149	8	0.055	14
example	0.143	9	0.065	9
code	0.136	10		
usage	0.131	11	0.14	<b>3</b>
feature	0.131	12	0.063	10
project	0.128	13		
get	0.118	14	0.056	13
learn	0.115	15		
run	0.105		0.072	7
develop			0.0697	8

factors following prior studies (Rajbahadur et al., 2019; Wang et al., 2018; Fan et al., 2021). We ended up with 14 factors (see remaining factors in Table 6).

**Model construction.** Most of the factors which are numeric and categorical could be directly derived from the repository and used as a feature when building the model. However, one factor *R\_heading* is text-based. Because our objective is for explanation and not for prediction, we use permutation variable importance (see below for details) to examine which factors are important. If we directly treated each occurred word in the headings as a feature in the model, it would be difficult to determine the importance of *R\_heading* as a whole, since permutation variable importance can only provide the importance for each individual feature (i.e., each word in the heading) rather than the importance of a collection of the words in headings. Therefore, we require a special process on *R\_heading*. Following previous studies (Shihab et al., 2013; Mizuno and Hata, 2010), we use a classifier to convert the heading text into a single numeric factor. More specially, we build a classifier to discriminate between popular and non-popular repositories using the heading text of repositories in the training set. The headings in repositories are divided into tokens using write space, where each token represents a single word. We did not stem the words or remove stop words since words appearing in the headings are usually important. Once the classifier (i.e., Classifier 1) is built, we get a probability score for each repository based on its heading text. Then, we build another classifier (i.e., classifier 2) using all factors in Table 2 except *R\_heading* together with the probability generated by classifier 1 to discriminate between popular and non-popular repositories.

In this study, we select random forest as our classifiers (for both classifiers 1 and 2). Since our goal here is to conduct

model explanation (i.e., understanding which factors are important) rather than prediction. Therefore, we exclude deep learning models since they are hard to explain (Castelvecchi, 2016; Schwartz-Ziv and Tishby, 2017). In general, a model with a higher accuracy indicates more reliable explanation results (Jia et al., 2021). We select random forest since random forest is generally highly accurate (Ghotra et al., 2015; Rajbahadur et al., 2017) and is also robust to outliers since it summarizes the classification results of multiple trees that are learned differently (Cutler et al., 2012). We actually compared the performance of different commonly used classification models from various families on our task, i.e., logistic regression (AUC is 0.90), SVM (AUC is 0.92), decision tree (AUC is 0.89), AdaBoost(0.97), k-nearest neighbors (0.91), random forest (AUC is 0.98). We found that random forest achieved the best performance, therefore we selected random forest for our study. We use the Python package *sklearn.ensemble.RandomForestClassifier* as the implementation of our random forest model. Note that in this study, we consider the factors related to the repository (i.e., *P\_Age*, *P\_size*, *P\_ProgramLanguage*, *P\_License*, and *P\_Max/Min/Avg/Median\_Tag*) as the controlling factors.

**Model validation.** We use the area under the Receiver Operating Characteristic Curve (i.e., ROC) to assess the explanatory power of the built random forest (Rajbahadur et al., 2017). A random classifier has an AUC (i.e., the area under the curve) of 0.5, while the AUC for a perfect classifier is equal to 1. In practice, most of the classification models have an AUC between 0.5 and 1. In general, an AUC of 0.5 suggests no discrimination (i.e., unable to distinguish between popular and non-popular repositories), 0.7 to 0.8 is considered acceptable, 0.8 to 0.9 is considered excellent, and more than 0.9 is considered outstanding (Mandrekar, 2010). To ensure that our model is not overfitted and the results are reliable, we perform 10-fold cross validation to validate the performance of the built model.

**Variable importance analysis.** In this study, we use permutation feature importance (Permutation) to measure the importance of each feature. Permutation is one of the oldest and most popularly used approaches in both machine learning (typically for random forest) and software engineering communities (Rajbahadur et al., 2021, 2017; Janitza et al., 2013; Jiarpakdee et al., 2020). It is defined to be the decrease in a model performance score when a single feature value is randomly shuffled. We use the Python implementation *sklearn.inspection.permutation\_importance*. We use partial dependence plot (PDP) to show the marginal effect a feature has on the predicted outcome of our built model (Friedman, 2001).

**Results: Our model explains our dataset well and has reliable performance.** We find that the average AUC of our built classifier is 0.98. An AUC large than 0.7 is generally considered acceptable (Hosmer Jr. et al., 2013; Zhou et al., 2020a). The results suggest that our built model explains the relationship between the studied factors and the popularity of a repository well.



**Table 6**

Variable importance for studies factors that are ordered from the most important to least important.

Factor	Variable importance (median)	p-value
<i>P_Avg_Tag</i>	0.49	2.5e−34
<i>P_License</i>	0.22	2.5e−34
<i>R_NumList</i>	0.16	2.5e−34
<i>P_ProgramLanguage</i>	0.11	1.55e−19
<i>R_NumUpdate</i>	0.089	1.9e−06
<i>R_UpdateInterval</i>	0.086	2.7e−16
<i>R_Length</i>	0.047	0.0049
<i>R_NumLink</i>	0.038	0.0003
<i>R_heading</i>	0.022	0.54
<i>R_NumIndent</i>	0.01	0.90
<i>R_NumCodeBlock</i>	0.01	0.62
<i>P_size</i>	0.011	7.6e−08
<i>R_NumImage</i>	0.01	0.22
<i>R_NumBadge</i>	0.007	0.007

After controlling factors related to the repository (e.g., *P\_Avg\_Tag* and *P\_License*), *R\_NumList*, *R\_NumUpdate*, and *R\_UpdateInterval* are statistically significantly important factors that discriminate popular and non-popular repositories. From Table 6, we observe that several repository related factors (i.e., controlling factors) are statistically significantly important to the built model (*p*-value). For instance, *P\_Avg\_Tag* is the most important factor to discriminate between two groups of repositories. *P\_License* and *P\_ProgramLanguage* rank as the second and fourth most important. Besides the controlling factors, we observe that *R\_NumList*, *R\_UpdateInterval*, *R\_NumLink*, *R\_NumUpdate*, *R\_Length* are statistically significantly to the model. In other words, the presentation and frequent updates are correlated to the popularity of a repository. For instance, CodeXGLUE<sup>5</sup> is a benchmark published by Microsoft for various software engineering tasks (e.g., code summarization and code search). The readme file contains various links to external supporting resources, such as large pre-trained models, other datasets, and related papers, which are helpful for developers to understand and reuse the benchmark, and develop new techniques based on the benchmark.

There is a positive correlation between the frequency and number of readme updates and the likelihood of a repository being popular. As shown in Fig. 4, the likelihood of a repository being classified as popular is positively correlated with *R\_NumUpdate* and negatively correlative with *R\_UpdateInterval*. For *R\_NumUpdate*, the likelihood increases sharply from 0 to 7 and gets relatively stable afterward. In terms of *R\_UpdateInterval*, the likelihood decreases sharply from 0 to ~600 and gets stable afterward. Note that smaller *R\_UpdateInterval* indicates more frequent updates. Such observations indicate that there is a positive correlation between the frequency and number of readme updates and the likelihood of a repository being popular.

There is a positive correlation between the number of lists and links and the popularity of a repository. Fig. 4 shows a similar trend for *R\_NumLink* and *R\_NumList*, sharply increasing from 0 to a certain value (i.e., 3 for *R\_NumList* and ~30 for *R\_NumLink*). Such observations suggest that proper markdown format and providing a proper amount of external support materials probably have a positive correlation with the popularity of repositories.

In summary, We find that readme related factors have a high correlation with the popularity of a repository. Therefore, we suggest repository owners should pay more attention to factors related to readme files, such as frequently updating the readme up to date and using proper markdown format to better organize the content in readme files.

After controlling repository-specific factors (e.g., *P\_Avg\_Topic*, *R\_NumList*, *R\_NumLink*, *R\_Length*, *R\_NumUpdate* and *R\_UpdateInterval* are statistically significantly important factors that discriminate popular and non-popular repositories. There is a positive correlation between the frequency, number of readme updates, the number of lists and links, and the likelihood of a repository being popular.

#### 4.3. RQ3: What content is updated in readme files between popular and non-popular GitHub repositories?

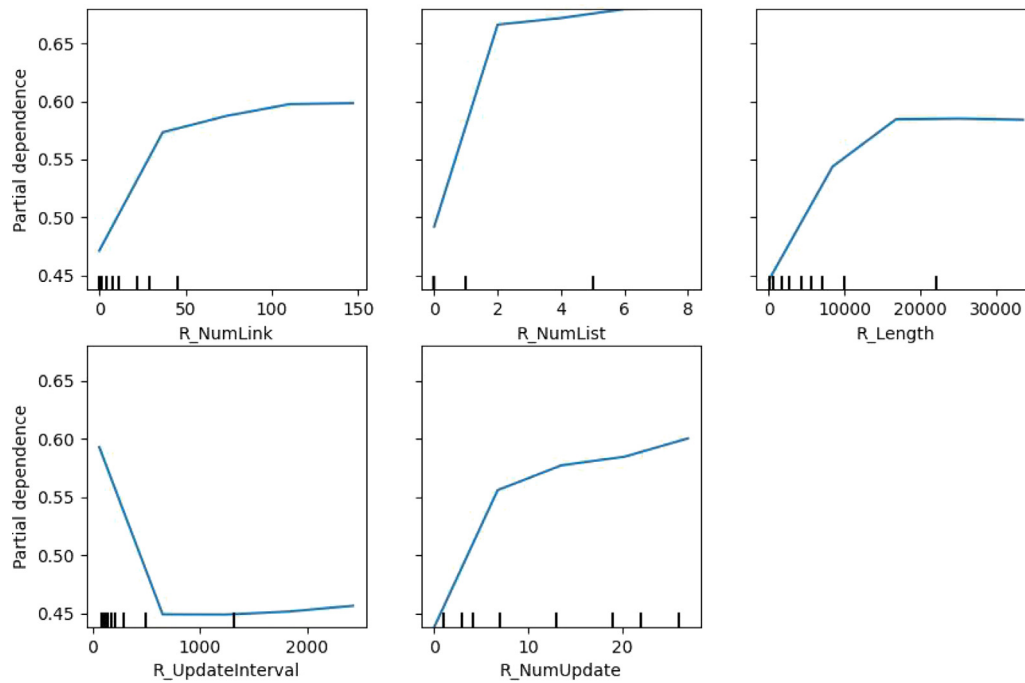
**Motivation** As observed in RQ2, factors related to readme update dimension are significantly important for the built model. In particular, *R\_UpdateInterval* is the most important factor after controlling repository-specific factors. Therefore, in this section, we explore the update activities in readme files to understand what and when were updated in readme files and compare these activities between popular and non-popular repositories.

**Approach** To further understand what content was updated in readme files of popular and non-popular repositories, we need to manually examine the updates in the commits. Note that since the number of stars of a repository changes over time, to understand the updates in readme files for popular repositories, we only extract the commits after the repositories become popular (i.e., after receiving at least 100 stars in our case). To do so, we first used star history data provided by Bytebase<sup>6</sup> to obtain the star history of each popular repository, and then identify the moment when the repository achieved at least 100 stars. We then extracted all commits of popular repositories that involved readme file updates after they became popular. Since there are too many commits and it is infeasible for us to manually examine all of them. Therefore, we randomly sampled 100 commits for popular and non-popular repositories, respectively. We use the coding scheme created by prior study (Prana et al., 2019) to label the content updated in readme files. To determine if the existing codes are sufficient for our study, the first author examined the 50 sampled commits with 25 from each group using open code, respectively. After this pilot examination, additional categories are derived. To alleviate the bias from human labeling, the first two authors then label the 50 sampled commits collaboratively. During this phase, the categories are revised and re-fined. Finally, we derived two additional categories “Badge” and “Presentation”. Then the first two authors independently applied the categories to the rest of sampled commits and any disagreements until a consensus is reached. During this phase, no new categories were introduced. The inter-rater agreement of this coding process has a Cohen’s kappa of 0.81. Table 7 presents the final categories we used to label the updated content in the readme files and example headings.

**Results: A larger portion of popular repositories were updated than non-popular repositories.** Although we observe that update activities are an important factor for discriminating popular and non-popular repositories, we observe that 34.6% (1728 out of 5,000) of the studied repositories have never been updated since creation. The situation is better in popular repositories, in which only 9% of them have never been updated, compared with non-popular repositories (35%). In terms of *R\_UpdateInterval* and *R\_NumUpdate* as shown in RQ1, popular repositories’ readme files are updated more frequently (254.7 days per update) than non-popular repositories (665.9 days per update).

<sup>5</sup> <https://github.com/microsoft/CodeXGLUE>

<sup>6</sup> <https://star-history.com/>



**Fig. 4.** The estimated probability when the values of  $R\_NumLink$ ,  $R\_NumList$ ,  $R\_UpdateInterval$ ,  $R\_Length$ , and  $R\_NumUpdate$  change. Y axis is the probability of a repository being classified as a popular one.

**Table 7**

Categories of updated content in the readme file.

Category	Definition	Example headings in readme files
What	Update the introduction of the repository.	Introduction, project background
Why	Update the description of advantages of the repository.	Advantages of the repository, comparison with related work
How	Update the information about how to use the project.	Getting started/quick start, how to run, installation, how to update, configuration, setup, requirements, dependencies, languages, platforms, demo, downloads, errors and bugs
When	Update the description of the status and plan of the repository.	Project status, versions, project plans, roadmap
Who	Update the information about who the project gives credit to.	Project team, community, mailing list, contact, acknowledgment, licence, code of conduct
References	Update the references to external information related to the repository.	API documentation, getting support, feedback, more information, translations, related projects, donation
Contribution	Update the information about how to contribute to the repository.	Contribution guidelines
Badge	Update the badge information for the repository.	NA
Presentation	Fix grammar issues and change markdown format.	NA

Most of the updates (47%) were performed to update references in popular repositories, while most updates were performed to update the content of “how” in non-popular repositories (34%). Fig. 5 compare the categories of content updated in readme files between popular and non-popular repositories.

On one hand, we observe that categories *How*, *References*, and *Presentation* are the top three most frequent categories in both popular and non-popular repositories. Those three categories account for 77% and 72% of popular and non-popular repositories, respectively. More importantly, we observe both in popular and

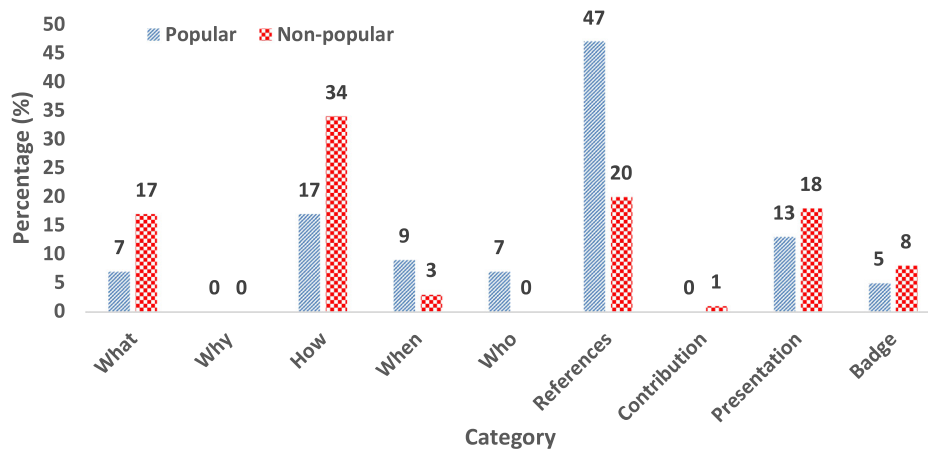


Fig. 5. Comparison of updated content in the readme files between popular and non-popular repositories.

non-popular projects, developers make significant efforts to update references, such as API documentation and related projects. For example, one commit was made to update a broken URL of API documentation to resolve an issue raised by developers.<sup>7</sup> This finding is compatible with prior studies on software knowledge platforms in which the authors found that a remarkable proportion of links are obsolete (Zhang et al., 2019; Liu et al., 2020). For instance, Zhang et al. (2019) found that 11.9% of the 5.5 million links in Stack Overflow answers are broken. Developers in both groups of repositories edit frequently the information on how to use the repositories, which is expected. One possible explanation is that almost all the repositories have such information and they are easy to become outdated (Prana et al., 2019).

**A remarkable portion of updates on readme files were performed to improve the presentation and update badges for both popular and non-popular repositories.** Developers also make efforts to improve the presentation of readme files (13% and 18% for popular and non-popular repositories). For instance, in a commit,<sup>8</sup> the contributor changed the format of section titles to make them more consistent. Based on our results in Section 4.2, the factors related to the presentation of readme files are important factors for discriminating between popular and non-popular repositories. In other words, the presentation of readme files is important for the popularity of a repository. 5% and 8% of updates on readme files are related to badge (i.e., category *Badge*) for popular and non-popular repositories, respectively. Various badges could be added to the readme file of a repository to indicate the standard/property of a repository, e.g., a 100% coverage badge could be added to indicate that 100% of the code in the repository is covered by its test cases. For instance, in a commit of openframeworks,<sup>9</sup> the contributor updated the badge status (e.g., changed the status of build on Linux 64 & Arm from “errored” to “passing”). They also changed the name of the badge for each platform to make the badge more specific (e.g., from “build” to “build-linux64-and-arm”). We compare the ratio of repositories in two groups that have at least one badge. Interestingly, we find more popular repositories have at least one badge (30%) than non-popular repositories (6.7%).

*Most of the updates (47%) were performed to update references in popular repositories, while most updates were performed to update the content of “how” in non-popular ones (34%). A remarkable portion of updates on readme files were performed to improve the presentation and update badges for both popular and non-popular repositories.*

## 5. Discussion

In this section, we discuss the implications of our findings and threats to validity.

### 5.1. Implications of our findings

**We suggest the GitHub repository contributors use proper markdown to organize the content of the readme file properly.** In RQ2, we observe that *R\_NumLink* and *R\_NumList* are significantly important factors and have a positive correlation with the popularity of a repository. We also observe that the table of contents are appeared more frequently in the readme file of popular repositories than in non-popular repositories, especially when the readme file is long. Therefore, we suggest contributors should pay attention to using proper markdown to organize the content of readme files.

**We suggest the GitHub repository contributors make continuous efforts to keep readme files up to date.** In RQ2, we observe that *R\_NumUpdate* and *R\_UpdateInterval* are statistically significantly important factors in the built model. Although our study does not reveal causation, at least we observe the frequency/number of readme updates positively correlates with the likelihood of a repository being popular. In other words, updates on the readme file are important. In RQ3, we observe that a remarkable proportion of updates were performed to update references and content related to “how”. Prior studies already prove that a good quality readme file can reduce the learning curve of users to get started and contribute to a repository (Steinmacher et al., 2014). Therefore, contributors should pay attention to maintaining the repository and make all content up to date.

In summary, the quality of the readme file is positively correlated with the popularity of repositories. Developers should pay attention to improving their readme files.

### 5.2. Future research direction

**Tools are encouraged to be developed to support the automatic update for readme files.** In RQ3, we observe that most of the updates (47%) in popular repositories were performed to

<sup>7</sup> <https://github.com/cristaloleg/go-advice/commit/c78d6c06d5327d731f61bbe416eeec2c28e79586>

<sup>8</sup> <https://github.com/emberjs/ember.js/commit/fb7022b8f4af64534470b002fe91a4da9b9d96b5>

<sup>9</sup> <https://github.com/openframeworks/openFrameworks/commit/faba5895d2687588ee0512ed90fb5288c8f3de1d>

update references due to various reasons, such as broken links and obsolescence. We also observe that developers frequently update the information on how to use the repository due to their obsolescence. Developers make significant efforts to maintain the references in their readme files. Therefore, future research is encouraged to develop tools to support such updates. One solution is to develop approaches to detect such obsolete references and content in the readme file during pushing a commit and recommend proper update solutions similar to prior research in just-in-time commit message generation (Xu et al., 2019; Liu et al., 2018). We also observe developers update the readme file to improve presentation, e.g., improving the formatting. Recently, large language models (e.g., GPT-4) have been proven in various NLP tasks (Bubeck et al., 2023; OpenAI, 2023). Future research is encouraged to leverage large language models for improving the presentation of readme files.

**Future research is encouraged to generate a readme template for a repository automatically.** In RQ1, we observe that popular and non-popular repositories share similar headings (e.g., usage and installation). Certain topics are frequently contained both in the readme file of popular and non-popular repositories. In other words, the readme files of GitHub repositories usually follow patterns containing certain pieces of information. Therefore, there is potential to generate a readme template automatically when a repository description is given using historical readme files as the training data using machine learning techniques, such as generative models (Radford et al., 2018; Nie et al., 2018). We encourage future research to investigate this.

### 5.3. Threats to validity

**Internal Validity** One threat to internal validation relates to the categorization of the studied repositories as popular and non-popular ones, in which we use a hard-cut 0 and 100 as the threshold. However, different thresholds may lead to different observations. To alleviate this threat, we build a random forest model using different thresholds. Instead of using the number of stars directly, we use percentages following prior studies (Zhou et al., 2020b; Wang et al., 2018). More specifically, we sort the repositories based on their stars and consider the top  $k\%$  and bottom  $k\%$  of the repositories as popular and non-popular. We tested three values 10, 20, and 30 for  $k$ . The findings still hold. The findings still hold. The  $R\_NumList$ ,  $R\_NumLink$ , and  $R\_UpdateInterval$  are still the most important features after controlling repository-related features.

One internal threat is that in RQ1, we use the most frequent words appearing in the headings as a proxy to estimate the content of readme files. Some readme files probably do not have headings. However, we find the proportion of such cases is low, i.e., 6% in popular but 29% in non-popular. Another internal threat is that in RQ1, we look at the top 15 most frequent words appearing in the headings to compare the content of readme files between popular and non-popular repositories. The findings may vary when we select different thresholds. To alleviate the threat, we select different thresholds (i.e., top 20 and 30) and we find the overlap between the two groups is 70% and 50%, respectively.

Another threat is that in data preparation, we did not consider outliers and their impact, typically the impact of the extremely popular repositories. To alleviate this threat, we used the three-sigma rule of thumb to identify outliers (Pukelsheim, 1994). For popular repositories, we identify 17 repositories that have more than 86,054 stars as outliers. For non-popular repositories, since all the repositories have zero star and we do not identify any outliers. After removing those outliers, we re-ran the analysis in RQ1 and RQ2 and our findings still hold.

In RQ2, we study the association between the studied features and the likelihood of a repository being popular. However, we cannot conclude the causation between them. Future research is encouraged to study their causation of them.

In RQ3, we involved human labeling. To reduce the bias from human labeling, each commit was labeled by two of the authors, and discrepancies were discussed until a consensus was reached. We also showed that the level of inter-rater agreement of the qualitative study is high (i.e., the value of Cohen's kappa is 0.81).

**External Validity** One external threat is that it is not clear whether our findings still hold on to other repositories (e.g., closed source repositories and open source repositories on other platforms). Although we performed the study on the most popular repository hosting platform (i.e., GitHub), future research is encouraged to replicate our analysis on other repositories.

Another threat is regarding the factors that there might be additional factors that could be associated with the popularity of a repository. In this study, we consider 6 factors related to the content of readme files, 3 factors related to the update activities on readme files, and 8 factors related to the repository as our controlling factors. There might be other factors associated with the popularity of a repository. However, our results show that the explanatory power of our models is high (0.98) when using the studied factors. Future studies should consider more features.

## 6. Conclusion and future work

In this study, we analyze the readme files of 5000 GitHub repositories across different domains and languages. We investigate the relationship between readme file related factors and the popularity of GitHub repositories. We observe that all studied readme file related features (e.g., the number of links and images, updates frequency) except heading content are significantly different between popular and non-popular repositories with non-negligible effect size. Popular and non-popular repositories share similar headings. We also observe the number of lists and the frequency of updates are significantly important factors that discriminate between popular and non-popular repositories and they positively correlate with the likelihood of a repository being popular. In summary, the quality of the readme file (e.g., proper organization of content and ensuring content is up to date) is positively correlated with the popularity of repositories. The most of the updates were performed to update reference in popular repositories, while the most updates were performed to update the content of "how" in non-popular ones. A remarkable portion of updates on readme files were performed to improve the presentation and update badges for both popular and non-popular repositories. Developers should pay attention to improving the readme file, such as improving the organization of content and ensuring the readme is up to date.

As the one of the first studies on this direction, we provide certain actionable insights for contributors on GitHub and future research. We encourage future research to verify those insights.

### CRedit authorship contribution statement

**Tianlei Wang:** Conceptualization, Investigation, Methodology, Software, Writing – original draft. **Shaowei Wang:** Conceptualization, Methodology, Software, Writing – original draft, Writing – review & editing. **Tse-Hsun (Peter) Chen:** Supervision, Writing – original draft, Writing – review & editing.



## Declaration of competing interest

We wish to confirm that there are no known conflicts of interest associated with this publication and there has been no significant financial support for this work that could have influenced its outcome. We confirm that the manuscript has been read and approved by all named authors and that there are no other persons who satisfied the criteria for authorship but are not listed. We further confirm that the order of authors listed in the manuscript has been approved by all of us. We confirm that we have given due consideration to the protection of intellectual property associated with this work and that there are no impediments to publication, including the timing of publication, with respect to intellectual property. In so doing we confirm that we have followed the regulations of our institutions concerning intellectual property. We further confirm that any aspect of the work covered in this manuscript that has involved either experimental animals or human patients has been conducted with the ethical approval of all relevant bodies and that such approvals are acknowledged within the manuscript.

## Data availability

Data will be made available on request.

## References

- Aggarwal, K., Hindle, A., Stroulia, E., 2014. Co-evolution of project documentation and popularity within Github. In: Proceedings of the 11th Working Conference on Mining Software Repositories. pp. 360–363.
- Bao, L., Xia, X., Lo, D., Murphy, G.C., 2019. A large scale study of long-time contributor prediction for Github projects. *IEEE Trans. Softw. Eng.* 47 (6), 1277–1298.
- Borges, H., Hora, A., Valente, M.T., 2016. Understanding the factors that impact the popularity of Github repositories. In: 2016 IEEE International Conference on Software Maintenance and Evolution. ICSME, IEEE, pp. 334–344.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y.T., Li, Y., Lundberg, S., et al., 2023. Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*.
- Castelvecchi, D., 2016. Can we open the black box of AI? *Nat. News* 538 (7623), 20.
- Cliff, N., 1993. Dominance statistics: Ordinal analyses to answer ordinal questions. *Psychol. Bull.* 114 (3), 494.
- Coelho, J., Valente, M.T., Silva, L.L., Shihab, E., 2018. Identifying unmaintained projects in Github. In: Proceedings of the 12th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement. pp. 1–10.
- Cutler, A., Cutler, D.R., Stevens, J.R., 2012. Random forests. In: Ensemble Machine Learning. Springer, pp. 157–175.
- Fan, Y., Xia, X., Lo, D., Hassan, A.E., Li, S., 2021. What makes a popular academic AI repository? *Empir. Softw. Eng.* 26 (1), 1–35.
- Farrar, D.E., Glauber, R.R., 1967. Multicollinearity in regression analysis: The problem revisited. *Rev. Econ. Stat.* 92–107.
- Friedman, J.H., 2001. Greedy function approximation: A gradient boosting machine. *Ann. Stat.* 1189–1232.
- Gharehyazie, M., Ray, B., Filkov, V., 2017. Some from here, some from there: Cross-project code reuse in Github. In: 2017 IEEE/ACM 14th International Conference on Mining Software Repositories. MSR, IEEE, pp. 291–301.
- Gharehyazie, M., Ray, B., Keshani, M., Zavosht, M.S., Heydarnoori, A., Filkov, V., 2019. Cross-project code clones in Github. *Empir. Softw. Eng.* 24 (3), 1538–1573.
- Ghotra, B., McIntosh, S., Hassan, A.E., 2015. Revisiting the impact of classification techniques on the performance of defect prediction models. In: 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, Vol. 1. IEEE, pp. 789–800.
- GitHub, 2022a. About READMEs. <https://docs.github.com/en/repositories/managing-your-repositorys-settings-and-features/customizing-your-repository/about-readmes>.
- GitHub, 2022b. REST API. <https://docs.github.com/en/rest>.
- GitHub, 2022c. Where the world builds software. <https://github.com/about>.
- Google, 2020. Use BigQuery to query GitHub data. <https://codelabs.developers.google.com/codelabs/bigquery-github>.
- Han, J., Deng, S., Xia, X., Wang, D., Yin, J., 2019. Characterization and prediction of popular projects on Github. In: 2019 IEEE 43rd Annual Computer Software and Applications Conference, Vol. 1. COMPSAC, IEEE, pp. 21–26.
- Hassan, F., Wang, X., 2017. Mining readme files to support automatic building of Java projects in software repositories. In: 2017 IEEE/ACM 39th International Conference on Software Engineering Companion. ICSE-C, IEEE, pp. 277–279.
- Hosmer Jr., D.W., Lemeshow, S., Sturdivant, R.X., 2013. Applied Logistic Regression, Vol. 398. John Wiley & Sons.
- Ikeda, S., Ihara, A., Kula, R.G., Matsumoto, K., 2019. An empirical study of readme contents for Javascript packages. *IEICE Trans. Inf. Syst.* 102 (2), 280–288.
- Janitza, S., Strobl, C., Boulesteix, A.-L., 2013. An AUC-based permutation variable importance measure for random forests. *BMC Bioinform.* 14 (1), 1–11.
- Jia, Y., Frank, E., Pfahringer, B., Bifet, A., Lim, N., 2021. Studying and exploiting the relationship between model accuracy and explanation quality. In: Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part II 21. Springer, pp. 699–714.
- Jiang, J., Lo, D., He, J., Xia, X., Kochhar, P.S., Zhang, L., 2017. Why and how developers fork what from whom in GitHub. *Empir. Softw. Eng.* 22 (1), 547–578.
- Jiarpakdee, J., Tantithamthavorn, C., Treude, C., 2020. The impact of automated feature selection techniques on the interpretation of defect models. *Empir. Softw. Eng.* 25 (5), 3590–3638.
- Kimble, J., 1992. Plain English: A charter for clear writing. *Cooley L. Rev.* 9, 1.
- Koskela, M., Simola, I., Stefanidis, K., 2018. Open source software recommendations using Github. In: International Conference on Theory and Practice of Digital Libraries. Springer, pp. 279–285.
- Laurent, A.M.S., 2004. Understanding Open Source and Free Software Licensing: Guide to Navigating Licensing Issues in Existing & New Software. " O'Reilly Media, Inc."
- Liu, Y., Noei, E., Lyons, K., 2022. How readme files are structured in open source Java projects. *Inf. Softw. Technol.* 148, 106924.
- Liu, Z., Xia, X., Hassan, A.E., Lo, D., Xing, Z., Wang, X., 2018. Neural-machine-translation-based commit message generation: How far are we? In: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering. pp. 373–384.
- Liu, J., Xia, X., Lo, D., Zhang, H., Zou, Y., Hassan, A.E., Li, S., 2020. Broken external links on stack overflow. *arXiv preprint arXiv:2010.04892*.
- Mandrekar, J.N., 2010. Receiver operating characteristic curve in diagnostic test assessment. *J. Thoracic Oncol.* 5 (9), 1315–1316.
- Mann, H.B., Whitney, D.R., 1947. On a test of whether one of two random variables is stochastically larger than the other. *Ann. Math. Stat.* 50–60.
- McDonald, N., Goggins, S., 2013. Performance and participation in open source software on Github. In: CHI'13 Extended Abstracts on Human Factors in Computing Systems. pp. 139–144.
- Mizuno, O., Hata, H., 2010. An integrated approach to detect fault-prone modules using complexity and text feature metrics. In: Advances in Computer Science and Information Technology. Springer, pp. 457–468.
- Munaiah, N., Kroh, S., Cabrey, C., Nagappan, M., 2017. Curating Github for engineered software projects. *Empir. Softw. Eng.* 22 (6), 3219–3253.
- Nie, W., Naroditska, N., Patel, A., 2018. Relgan: Relational generative adversarial networks for text generation. In: International Conference on Learning Representations.
- OpenAI, 2023. GPT-4 Technical Report.
- Prana, G.A.A., Treude, C., Thung, F., Atapattu, T., Lo, D., 2019. Categorizing the content of Github readme files. *Empir. Softw. Eng.* 24 (3), 1296–1327.
- Pukelsheim, F., 1994. The three sigma rule. *Amer. Statist.* 48 (2), 88–91.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., 2018. Improving language understanding by generative pre-training.
- Rajbahadur, G.K., Wang, S., Ansaldi, G., Kamei, Y., Hassan, A.E., 2021. The impact of feature importance methods on the interpretation of defect classifiers. *IEEE Trans. Softw. Eng.*
- Rajbahadur, G.K., Wang, S., Kamei, Y., Hassan, A.E., 2017. The impact of using regression models to build defect classifiers. In: 2017 IEEE/ACM 14th International Conference on Mining Software Repositories. MSR, IEEE, pp. 135–145.
- Rajbahadur, G.K., Wang, S., Kamei, Y., Hassan, A.E., 2019. Impact of discretization noise of the dependent variable on machine learning classifiers in software engineering. *IEEE Trans. Softw. Eng.* 47 (7), 1414–1430.
- Romano, J., Kromrey, J.D., Coraggio, J., Skowronek, J., 2006. Appropriate statistics for ordinal level data: Should we really be using t-test and Cohen'sd for evaluating group differences on the NSSE and other surveys. In: Annual Meeting of the Florida Association of Institutional Research, Vol. 177. p. 34.
- Sharma, A., Thung, F., Kochhar, P.S., Sulistya, A., Lo, D., 2017. Cataloging Github repositories. In: Proceedings of the 21st International Conference on Evaluation and Assessment in Software Engineering. pp. 314–319.
- Shihab, E., Ihara, A., Kamei, Y., Ibrahim, W.M., Ohira, M., Adams, B., Hassan, A.E., Matsumoto, K.-i., 2013. Studying re-opened bugs in open source software. *Empir. Softw. Eng.* 18 (5), 1005–1042.
- Shwartz-Ziv, R., Tishby, N., 2017. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*.
- Steinmacher, I., Gerosa, M.A., Redmiles, D., 2014. Attracting, onboarding, and retaining newcomer developers in open source software projects. In: Workshop on Global Software Development in a CSCW Perspective.

- Tian, Y., Nagappan, M., Lo, D., Hassan, A.E., 2015. What are the characteristics of high-rated apps? A case study on free android applications. In: 2015 IEEE International Conference on Software Maintenance and Evolution. ICSME, IEEE, pp. 301–310.
- Treude, C., Middleton, J., Atapattu, T., 2020. Beyond accuracy: assessing software documentation quality. In: Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. pp. 1509–1512.
- Trockman, A., Zhou, S., Kästner, C., Vasilescu, B., 2018. Adding sparkle to social coding: An empirical study of repository badges in the NPM ecosystem. In: Proceedings of the 40th International Conference on Software Engineering. pp. 511–522.
- Wang, S., Chen, T.-H., Hassan, A.E., 2018. Understanding the factors for fast answers in technical Q&A websites. *Empir. Softw. Eng.* 23 (3), 1552–1593.
- Weber, S., Luo, J., 2014. What makes an open source code popular on Git Hub? In: 2014 IEEE International Conference on Data Mining Workshop. IEEE, pp. 851–855.
- Xu, S., Yao, Y., Xu, F., Gu, T., Tong, H., Lu, J., 2019. Commit message generation for source code changes. In: IJCAI.
- Zhang, Y., Lo, D., Kochhar, P.S., Xia, X., Li, Q., Sun, J., 2017. Detecting similar repositories on GitHub. In: 2017 IEEE 24th International Conference on Software Analysis, Evolution and Reengineering. SANER, IEEE, pp. 13–23.
- Zhang, H., Wang, S., Chen, T.-H., Zou, Y., Hassan, A.E., 2019. An empirical study of obsolete answers on stack overflow. *IEEE Trans. Softw. Eng.* 47 (4), 850–862.
- Zhou, J., Wang, S., Bezemer, C.-P., Hassan, A.E., 2020a. Bounties on technical Q&A sites: A case study of stack overflow bounties. *Empir. Softw. Eng.* 25 (1), 139–177.
- Zhou, J., Wang, S., Bezemer, C.-P., Zou, Y., Hassan, A.E., 2020b. Studying the association between bountysource bounties and the issue-addressing likelihood of Github issue reports. *IEEE Trans. Softw. Eng.*
- Zhu, J., Zhou, M., Mockus, A., 2014. Patterns of folder use and project popularity: A case study of GitHub repositories. In: Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement. pp. 1–4.

**Tianlei Wang** is an embedded software developer at Fortinet. He received his bachelor's degree in Computer Ccience from University of Manitoba.

**Shaowei Wang** is an assistant professor in the Department of Computer Science at the University of Manitoba. He leads the Software Management, Maintenance, and Reliability Lab (Mamba). His research interests include software engineering, machine learning, data analytics for software engineering, automated debugging, and secure software development. His work has been published in flagship conferences and journals such as ICSE, FSE, ASE, TSE, TOSEM, and EMSE. Certain work received best paper award and test of time award. He serves regularly as a program committee member of international conferences in the field of software engineering, such as ICSE, ASE, ICSME, SANER, ICPC, and MSR, and he is a regular reviewer for software engineering journals such as TOSEM, JSS, EMSE, and TSE. He is one of four recipients of the 2018 distinguished reviewer award for the Springer EMSE (SE's highest impact journal). He obtained his Ph.D. from Singapore Management University and his B.Sc. from Zhejiang University. More information at: <https://sites.google.com/view/mambalab>.

**Tse-Hsun (Peter) Chen** is an Associate Professor in the Department of Computer Science and Software Engineering at Concordia University, Montreal, Canada. He leads the Software Performance, Analysis, and Reliability (SPEAR) Lab, which focuses on conducting research on performance engineering, program analysis, log analysis, production debugging, and mining software repositories. His work has been published in flagship conferences and journals such as ICSE, FSE, TSE, EMSE, and MSR. He serves regularly as a program committee member of international software engineering conferences, such as FSE, ASE, ICSME, SANER, and ICPC. He is also a regular reviewer for software engineering journals such as JSS, EMSE, and TSE, and serves as a guest editor for special issues in EMSE. Dr. Chen obtained his B.Sc. from the University of British Columbia, and M.Sc. and Ph.D. from Queen's University. Besides his academic career, Dr. Chen also worked as a software performance engineer at BlackBerry for over four years. Early tools developed by Dr. Chen were integrated into industrial practice for ensuring the quality of large-scale enterprise systems. More information at: <http://petertsehsun.github.io/>.