

Understanding the Factors that Impact the Popularity of GitHub Repositories

Hudson Borges, Andre Hora, Marco Túlio Valente
ASERG Group, Department of Computer Science (DCC)
Federal University of Minas Gerais (UFMG), Brazil
E-mail: {hsborges, hora, mtov}@dcc.ufmg.br

Abstract—Software popularity is a valuable information to modern open source developers, who constantly want to know if their systems are attracting new users, if new releases are gaining acceptance, or if they are meeting user's expectations. In this paper, we describe a study on the popularity of software systems hosted at GitHub, which is the world's largest collection of open source software. GitHub provides an explicit way for users to manifest their satisfaction with a hosted repository: the stargazers button. In our study, we reveal the main factors that impact the number of stars of GitHub projects, including programming language and application domain. We also study the impact of new features on project popularity. Finally, we identify four main patterns of popularity growth, which are derived after clustering the time series representing the number of stars of 2,279 popular GitHub repositories. We hope our results provide valuable insights to developers and maintainers, which could help them on building and evolving systems in a competitive software market.

Index Terms—GitHub; Software Popularity; Open Source software; Social coding.

I. INTRODUCTION

GitHub is the world's largest collection of open source software, with around 9 million users and 17 million public repositories.¹ In addition to a git-based version control system, GitHub integrates many features for social coding. For example, developers can fork their own copy of a repository, work and improve the code locally, and then submit a pull request to integrate their changes in the main repository. The key characteristics and challenges of this pull-based development model is recently explored in many studies [1]–[4]. However, GitHub also supports other typical features from social networks. For example, users can star a repository to manifest their interest or satisfaction with the hosted project. Consequently, the number of stars of a GitHub repository can be seen as a proxy of its popularity. Currently, the two most popular repositories on GitHub are FREECODECAMP/FREECODECAMP (a coding education software, which claims to have more than 300K users²) and TWBS/BOOTSTRAP (a library of HTML and CSS templates, which is used by almost 7M web sites³).

A deep understanding of the factors that impact the number of stars of GitHub repositories is important to software

developers because they want to know whether their systems are attracting new users, whether the new releases are gaining acceptance, whether their systems are as popular as competitor systems, etc. Unfortunately, we have few studies about the popularity of GitHub systems. The exceptions are probably an attempt to differentiate popular and unpopular Python repositories using machine learning techniques [5] and a study on the effect of project's popularity on documentation quality [6]. By contrast, popularity is extensively studied on other social platforms, like YouTube [7], [8] and Twitter [9], [10]. These studies are mainly conducted to guide content generators on producing successful social media content. Similarly, knowledge on software popularity might also provide valuable insights on how to build and evolve systems in a competitive market.

This paper presents an in-depth investigation on the popularity of GitHub repositories. We first collected historical data about the number of stars of 2,500 popular repositories. We use this dataset to answer four research questions:

RQ #1: How popularity varies per programming language, application domain, and repository owner? The goal is to provide an initial view about the popularity of the studied systems, by comparing the number of stars according to programming language, application domain, and repository owner (user or organization).

RQ #2: Does popularity correlate with other characteristics of a repository, like age, number of commits, number of contributors, and number of forks? This investigation is important to check whether there are factors that can be worked to increase a project's popularity.

RQ #3: How early do repositories get popular? With this research question, we intend to check whether gains of popularity are concentrated in specific phases of a repository's lifetime, specifically in early releases.

RQ #4: What is the impact of new features on popularity? This investigation can show if relevant gains in popularity happen due to new features (implemented in new releases).

In the second part of the paper, we identify four patterns of popularity growth in GitHub, which are derived after clustering the time series that describe the growth of the number of stars of the systems in our dataset. These patterns can help developers to understand how their systems have grown in the

¹<https://github.com/search/advanced>, verified on 04/04/2016

²<https://www.freecodecamp.com/about>, verified on 04/04/2016

³<http://trends.builtwith.com/docinfo/Twitter-Bootstrap>, verified 04/04/2016

past and to predict future growth trends. Finally, in the third part of the paper, we present a qualitative study with GitHub developers to clarify some findings and themes of our study. A total of 44 developers participated to the study.

The main contribution of this paper is an investigation of factors that may impact the popularity of GitHub repositories, including the identification of the major patterns that can be used to describe popularity trends. Although similar studies exist for social networks, to our knowledge we are the first to focus on the popularity of systems hosted in an ultra-large repository of open source code.

Organization: The rest of this paper is organized as follows. Section II describes and characterizes the dataset used in this study. Section III uses this dataset to provide answers to four questions about the popularity of GitHub repositories. Section IV documents four patterns that describe the popularity growth of GitHub systems. Section V reports the feedback of GitHub developers about three specific themes of our study. Section VI discusses threats to validity and Section VII presents related work. Finally, Section VIII concludes the paper and lists future work.

II. DATASET

The dataset used in this paper includes the top-2,500 public repositories with more stars in GitHub. We limit the study to 2,500 repositories for two major reasons. First, to focus on the characteristics of the highly popular GitHub systems. Second, because we investigate the impact of application domain on popularity, which demands a manual classification of the domain of each system.

All data was obtained using the GitHub API, which provides services to search public repositories and to retrieve specific data about them (e.g., stars, commits, contributors, and forks). The data was collected on March 28th, 2016. Besides retrieving the number of stars on this date for each system, we also relied on GitHub API to collect historical data about the number of stars. For this purpose, we used a service from the API that returns all star events of a given repository. For each star, these events store the date and the user who starred the repository. However, GitHub API returns at most 100 events by request (i.e., a page) and at most 400 pages. For this reason, it is not possible to retrieve all stars events of systems with more than 40K stars, as is the case of FREECODECAMP, BOOTSTRAP, ANGULARJS, D3, and FONT-AWESOME. Therefore, these five systems are not considered when answering the third and fourth research questions (that depend on historical data) and also on the study about common growth patterns (Section IV).

Table I shows descriptive statistics on the number of stars of the repositories in our dataset. The number of stars ranges from 2,150 (for CYBERAGENT/ANDROID-GPUIMAGE) to 97,948 stars (for FREECODECAMP/FREECODECAMP). The median number of stars is 3,441.

Age, Commits, Contributors, and Forks: Figure 1 shows the distribution of the age (in number of weeks), number of

TABLE I: Descriptive statistics on the number of stars of the repositories in our dataset of 2,500 popular GitHub systems

Min	1st Quartile	2nd Quartile	3rd Quartile	Max
2,150	2,682	3,441	5,331	97,948

commits, number of contributors, and number of forks for the 2,500 systems in the dataset. For age, the first, second, and third quartiles are 101, 169, and 250 weeks, respectively. For number of commits, the first, second, and third quartiles are 228, 608, and 1,721, respectively. For number of contributors, the first, second, and third quartiles are 17, 41, and 96, respectively;⁴ and for number of forks, the first, second, and third quartiles are 298, 533, and 1,045, respectively. Therefore, the systems in our dataset are mature and have many commits and contributors.

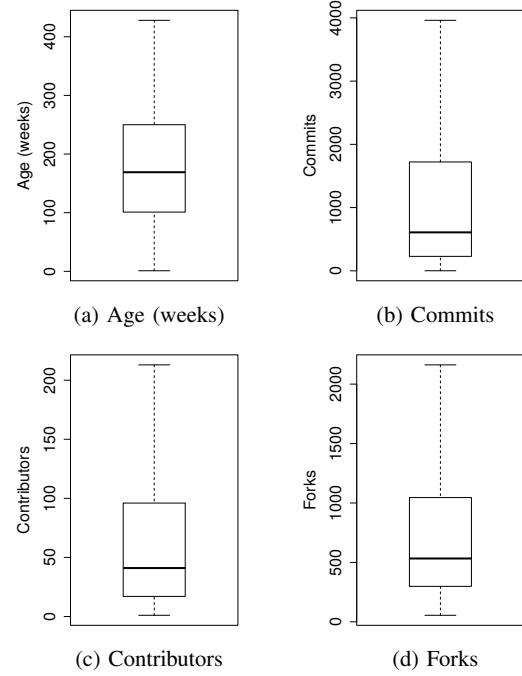


Fig. 1: Age, number of commits, number of contributors, and number of forks (outliers are omitted)

Programming Language: As returned by GitHub API, the language of a repository is the one with the highest percentage of source code, considering the files in the repository. Figure 2 shows the distribution of the systems per programming language. JavaScript is the most popular language (855 repositories, 34.2%), followed by Python (203 repositories, 8.1%), Java (202 repositories, 8.0%), Objective-C (188 repositories, 7.5%), and Ruby (178 repositories, 7.1%). Despite a concentration of systems in these languages, the

⁴We report data from contributors as retrieved by GitHub API. This data may be different from the one presented on the project's page at GitHub, which only counts contributors with GitHub account.

dataset includes systems in 53 languages, including Groovy, R, Julia, and XSLT (all with just one repository).

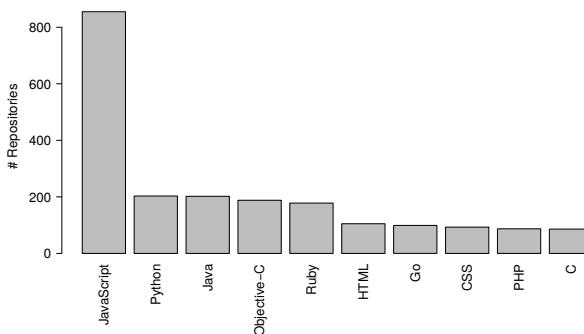


Fig. 2: Top-10 languages by number of repositories

Owner: We also provide results grouped by repository owner. In GitHub, a repository can be owned by a user (e.g., TORVALDS/LINUX) or by an organization (e.g., FACEBOOK/REACT). In our dataset, 1,263 repositories (50.5%) are owed by organizations and 1,237 repositories (49.5%) by users.

Application Domain: In the study reported in this paper, results are also grouped by application domain. However, different from other source code repositories, like SourceForge, GitHub does not include information about the application domain of a project. For this reason, we manually classified the domain of each system in our dataset. Initially, the first and third authors of this paper inspected the description of the top-200 repositories to provide a first list of application domains. After this initial classification, the first author inspected the short description (and in many cases the GitHub page and the project's page) of the remaining 2,300 repositories. During this process, he also marked the repositories with dubious classification decisions. These particular cases were discussed by the first and second authors, to reach a consensus decision. The spreadsheet with the proposed classification is publicly available at <https://goo.gl/73Sbvz>.

The systems are classified in the following six domains:

- Application software: systems that provide functionalities to end-users, like browsers and text editors (e.g., WORDPRESS/WORDPRESS and ADOBE/BRACKETS).
- System software: systems that provide services and infrastructure to other systems, like operating systems, middleware, servers, and databases (e.g., TORVALDS/LINUX and MONGODB/MONGO).
- Web libraries and frameworks (e.g., TWBS/BOOTSTRAP and ANGULAR/ANGULAR.JS).
- Non-web libraries and frameworks (e.g., GOOGLE/GUAVA and FACEBOOK/FRESCO).
- Software tools: systems that support software development tasks, like IDEs, package managers, and compilers (e.g., HOMEBREW/HOMEBREW and GIT/GIT).

- Documentation: repositories with documentation, tutorials, source code examples, etc. (e.g., ILUWATAR/JAVA-DESIGN-PATTERNS).

Figure 3 shows the number of systems in each domain. The top-3 domains are web libraries and frameworks (837 repositories, 33%), non-web libraries and frameworks (641 repositories, 25%), and software tools (470 repositories, 18%).

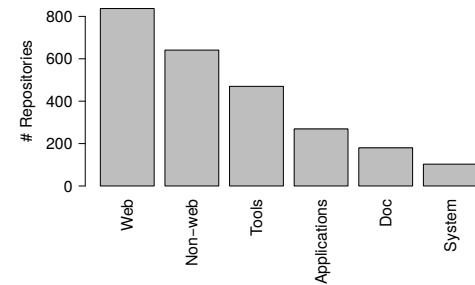


Fig. 3: Number of repositories by domain

III. RESULTS

In this section, we use the described dataset to answer the four research questions listed in the paper's introduction.

RQ #1: *How popularity varies per programming language, application domain, and repository owner?*

Figure 4 shows the distribution of the number of stars for the top-10 languages with more repositories. The top-3 languages whose repositories have the highest median number of stars are: JavaScript (3,697 stars), Go (3,549 stars), and HTML (3,513 stars). The three languages whose repositories have the lowest median number of stars are PHP (3,245 stars), Java (3,224 stars), and Python (3,099 stars). By applying the Kruskal-Wallis test to compare multiple samples, we find that the distribution of the number of stars per language is different (*p-value* = 0.001). Thus, we can consider that programming language may impact on system popularity.

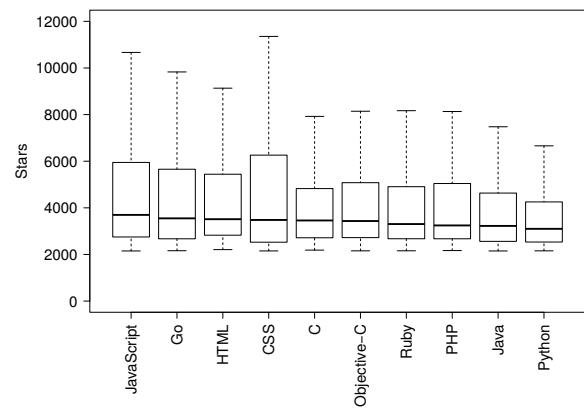


Fig. 4: Stars by programming language (considering only the top-10 languages with more repositories)

Figure 5 shows the distribution of the number of stars for the repositories in each application domain. The median number of stars varies as follow: systems software (3,807 stars), web libraries and frameworks (3,596 stars), documentation (3,547 stars), software tools (3,538 stars), applications (3,443 stars), and now-web libraries and frameworks (3,204 stars). By applying the Kruskal-Wallis test, we find that the distribution of the number of stars by domain is different ($p\text{-value} < 0.001$). Therefore, application domain is also an important factor that may impact on system popularity.

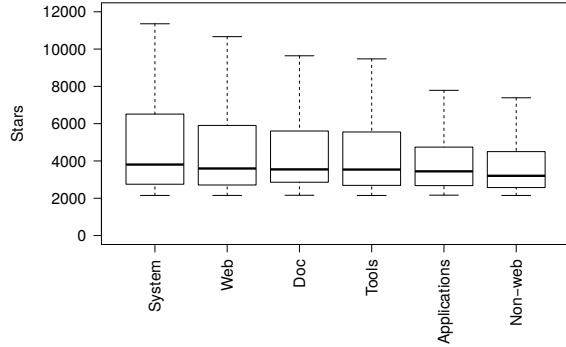


Fig. 5: Popularity by application domain

Finally, Figure 6 shows how popularity varies depending on the repository owner (i.e., user or organization). The median number of stars is 3,622 stars for repositories owned by organizations and 3,298 stars for repositories owned by users. By applying the Mann-Whitney test, we detect that indeed these distributions are different ($p\text{-value} < 0.001$). We hypothesize that repositories owned by organizations—specially major software companies and free software foundations—have more funding and resources, which somehow explains their higher popularity.

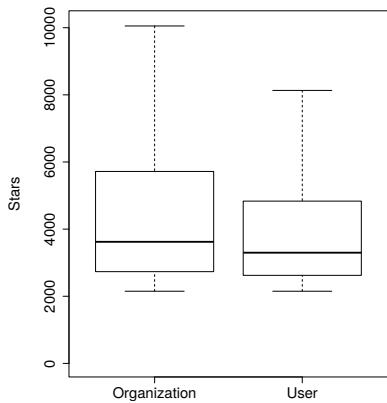


Fig. 6: Popularity by repository owner

Summary: The top-5 languages with more stars are JavaScript, Python, Java, Objective-C, and Ruby (Figure 2). However, the top-5 languages whose systems have the highest median number of stars are JavaScript, Go, HTML, CSS, and C (Figure 4).

The top-3 application domains whose repositories have more stars are systems software, web libraries and frameworks, and documentation. Repositories owned by organizations are more popular than the ones owned by individuals.

RQ #2: Does popularity correlate with repository's age, number of commits, number of contributors, and number of forks?

Figure 7 shows scatterplots correlating the number of stars with the age (in number of weeks), number of commits, number of contributors, and number of forks of a repository. First, the plots suggest that stars are not correlated with the repository's age (Figure 7a). We have old repositories with few stars and new repositories with many stars. For example, APPLE/SWIFT has only five months and 28,105 stars, while MOJOMBO/CHRONIC has more than 8 years and 2,440 stars. Essentially, this result shows that repositories gain stars at different speeds. We ran Spearman's rank correlation test and the resulting correlation coefficient is close to zero ($\rho = 0.0757$ and $p\text{-value} < 0.001$).

The scatterplot in Figure 7b suggests that stars are weakly correlated with number of commits ($\rho = 0.249$ with $p\text{-value} < 0.001$). Similarly, as presented in Figure 7c stars are weakly correlated with contributors ($\rho = 0.341$ with $p\text{-value} < 0.001$). In this figure, a logarithm scale is used in both axes; the line represents the identity relation: below the line are the systems with more contributors than stars. Interestingly, two systems indeed have more contributors than stars: RASPBERRYPI/LINUX (17,766 contributors and 2,739 stars) and LINUXBREW/LINUXBREW (7,304 contributors and 2,241 stars). This happens because they are forks of highly successful repositories (TORVALDS/LINUX and HOMEBREW/BREW, respectively). The top-3 systems with more stars per contributor are SHADOWSOCKS/SHADOWSOCKS (12,287 stars/contributor), OCTOCAT/SPOON-KNIFE (9,944 stars/contributor), and WG/WRK (7,923 stars/contributor). All these systems have just one contributor. The three systems with less stars per contributor are ANDROID/PLATFORM_FRAMEWORKS_BASE (2.28 stars/contributor), FFmpeg/FFmpeg (2.39 stars/contributor), and DEFINITELYTYPED/DEFINITELYTYPED (2.68 stars/contributor).

Finally, Figure 7d shows plots correlating a system popularity and its number of forks. As visually suggested by the figure, there is a strong positive correlation between stars and forks ($\rho = 0.549$ and $p\text{-value} < 0.001$). For example, TWBS/BOOTSTRAP is the second repository with the highest number of stars and the second one with more forks. ANGULAR/ANGULAR.JS is the third repository in number of stars and the third one with more forks. In Figure 7d, we can also see that only nine systems (0.36%) have more forks than stars. As examples, we have a repository that just provides a tutorial for forking a repository (OCTOCAT/SPOONKNIFE) and a popular puzzle game (GABRIELECIRULLI/2048), whose success motivated many forks with variations of the original

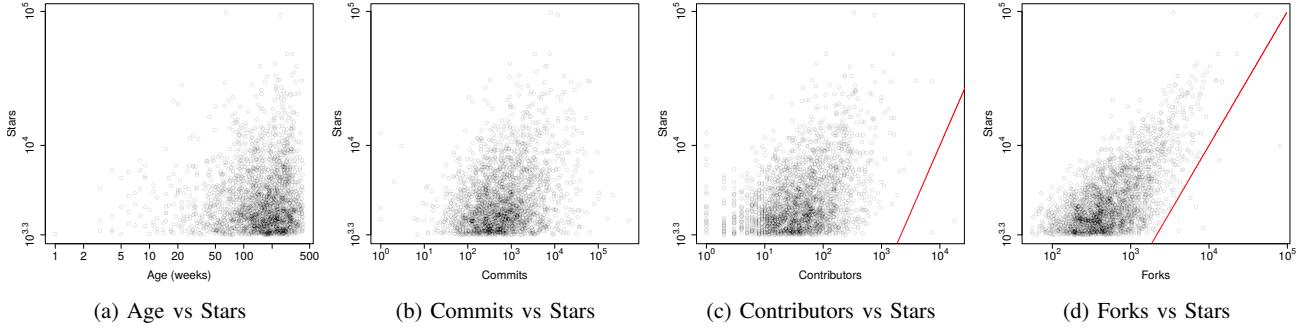


Fig. 7: Correlation analysis. In subfigures (c) and (d), the line is the identity relation

implementation. Since the game can be downloaded directly from the web, we hypothesize that it receives most users' feedback in the web and not on GitHub.

Summary: *There is no correlation between numbers of stars and the repository's age; however, there is a weak correlation with commits and contributors. Moreover, a strong correlation with forks was found.*

RQ #3: How early do repositories get popular?

Figure 8 shows the cumulative distribution of the fraction of time a repository takes to receive at least 10%, at least 50%, and at least 90% of its stars. Specifically, the y-axis shows the fraction of repositories that achieved 10%, 50%, and 90% of their stars in a period of time that does not exceed the fraction of time shown in the x-axis. Around 40% of the repositories receive 10% of their stars very early, in the first days after the initial release (label A, in Figure 8). We hypothesize that many of these initial stars come from early adopters, who start commenting and using novel open source software immediately after they are out. After this initial burst of popularity, the growth of half of the repositories tend to stabilize. For example, half of the repositories take 51% of their age to receive 50% of their stars (label B); and half of the repositories take 91% of their age to receive 90% of their total number of stars (label C).

Summary: *Repositories have a tendency to receive more stars right after their first public release. After this period, for half of the repositories the growth rate tends to stabilize.*

RQ #4: What is the impact of new features on popularity?

In this research question, we investigate the impact of new features on the popularity of GitHub repositories. The goal is to check whether the implementation of new features (resulting in new releases of the projects) contribute to a boost in popularity. Specifically, we selected 834 repositories from our dataset (33.3%) that follow a semantic versioning convention to number releases.⁵ In such systems, versions

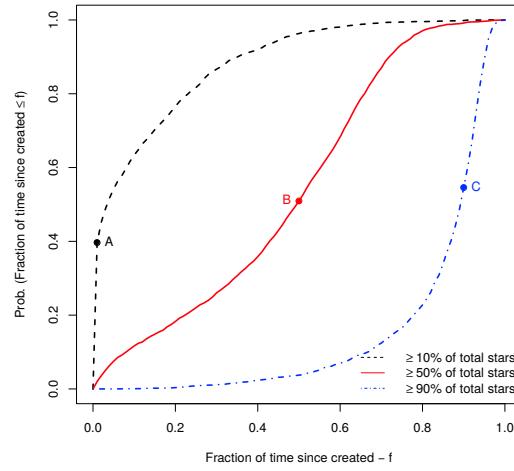


Fig. 8: Cumulative distribution of the fraction of time a repository takes to receive 10%, 50%, and 90% of its stars

are identified by three integers, in the format $x.y.z$, with the following semantics: increments in x denote major releases, which can be incompatible with old versions; increments in y denote minor releases, which add functionality in a backwards-compatible manner; and increments in z denote patches implementing bug fixes. In our sample, we identified 580 major releases and 4,343 minor releases.

First, we counted the fraction of stars received by each repository in the week following all releases (major or minor) and just after major releases. As mentioned, the goal is to check the impact of new releases in the number of stars. Figure 9 shows the distribution of these fractions. When considering all releases, the fraction of stars gained in the first week after the releases is 1.1% (first quartile), 3.2% (second quartile), and 10.2% (third quartile). For the major releases only, it is 0.5% (first quartile), 1.4% (second quartile), and 4.3% (third quartile). SO-FANCY/DIFF-SO-FANCY (a visualization for git diffs) is the repository with the highest fraction of stars received after releases. The repository has 53 days and 4,402 stars. Since it has a fast releasing rate (one new release

⁵<http://semver.org>

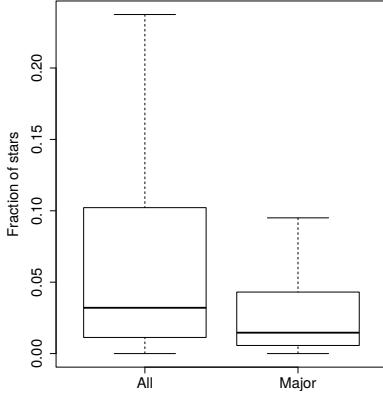


Fig. 9: Fraction of stars gained in the first week after all releases and just after the major releases

per week, on average), it gained almost of its stars (89.1%) in the weeks after releases.

We computed a second ratio for each repository: *fraction of stars in the week following all releases or just major releases (FS) / fraction of time represented by these weeks (FT)*. When $FS/FT > 1$, the repository gains proportionally more stars after the releases. Figure 10 shows boxplots with the results of FS/FT for all repositories. When considering all releases, we have that FS/FT is 0.80 (first quartile), 1.25 (second quartile), and 1.98 (third quartile). For major releases only, we have that FS/FT is 0.81 (first quartile), 1.53 (second quartile), and 2.98 (third quartile).

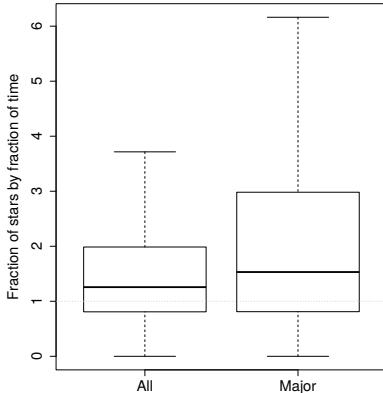


Fig. 10: Fraction of stars in the week following all releases (or just the major releases) / fraction of time represented by these weeks

Figure 11 shows the median values of FS/FT computed using stars gained after n weeks ($1 \leq n \leq 4$). This ratio decreases, both for major and for all releases. Therefore, although there is some gains of stars after releases, they tend to decrease after few weeks.

Summary: There is an acceleration in the number of stars gained just after releases. For example, half of the repositories

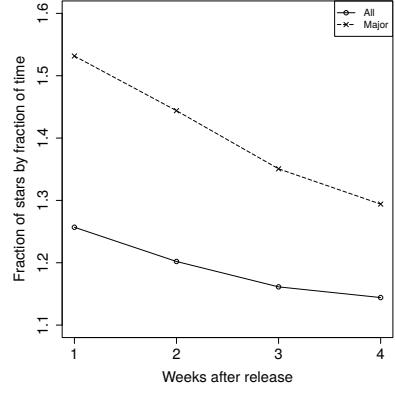


Fig. 11: Fraction of stars by fraction of time (median values), computed using different time intervals (in weeks)

gain at least 53% more stars in the week following major releases, than in the other weeks (see Figure 10). However, because repositories usually have much more weeks without releases than with releases, this phenomenon is not sufficient to generate a major concentration of popularity gain after releases. For example, 75% of the systems gain at most 4.3% of their stars in the week following major releases (see Figure 9).

IV. POPULARITY GROWTH PATTERNS

In this section, we investigate common patterns of popularity growth concerning the GitHub repositories in our dataset. To this purpose, we use the KSC algorithm [11]. This algorithm clusters time series with similar shapes using a metric that is invariant to scaling and shifting. The algorithm is used in other studies to cluster time series representing the popularity of YouTube videos [12] and Twitter [9]. Like K-means [13], KSC requires as input the number of clusters k .

Because the time series provided as input to KSC must have the same length, we only consider data regarding the last 52 weeks (one year). Due to this restriction, we exclude 216 repositories (8.6%) that have less than 52 weeks. We use the β_{CV} heuristic [14] to define the best number k of clusters. β_{CV} is defined as the ratio of the coefficient of variation of the intracluster distances and the coefficient of variation of the intercluster distances. The smallest value of k after which the β_{CV} ratio remains roughly stable should be selected. This means that new added clusters affect only marginally the intra and intercluster variations [8]. In our dataset, the values of β_{CV} stabilize for $k = 4$ (see Figure 13).

A. Proposed Growth Patterns

Figure 12 shows plots with the time series in each cluster. The time series representing the clusters' centroids are presented in Figure 14. The time series in clusters C1, C2, and C3 suggest a linear growth, but at different speeds. On the other hand, the series in cluster C4 suggest repositories with a sudden growth on the number of stars. We refer to these clusters as including systems with *Slow*, *Moderate*, *Fast*, and *Viral* Growth, respectively.

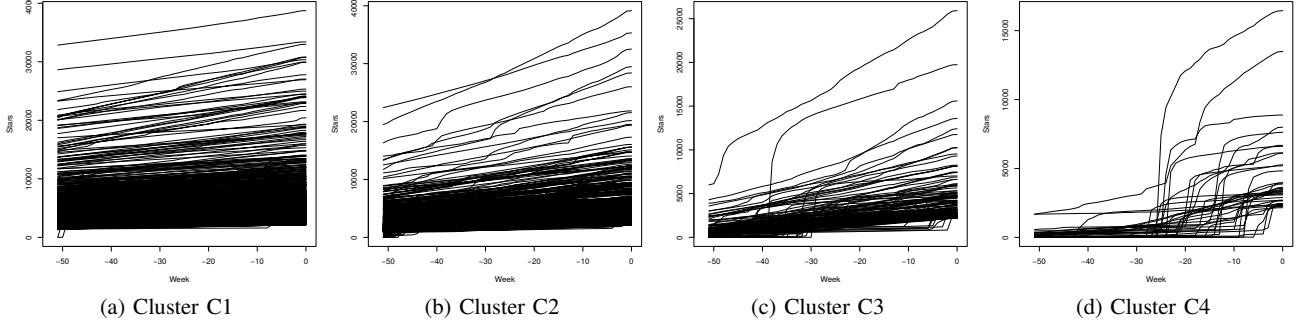


Fig. 12: Clusters of time series representing the growth of the number of starts of 2,279 GitHub repositories

TABLE II: Popularity Growth Patterns

Pattern	Cluster	# Repositories	% Growth	Top-3 Repositories
Slow	C1	1,497 (65.7%)	27.3	JQUERY/JQUERY, H5BP/HTML5-BOILERPLATE, and METEOR/METEOR
Moderate	C2	614 (26.9%)	94.0	FACEBOOK/REACT, ROBBYRUSSELL/OH-MY-ZSH, and AIRBNB/JAVASCRIPT
Fast	C3	131 (5.7%)	469.2	ATOM/ELECTRON, GOOGLE/MATERIAL-DESIGN-LITE, and VUEJS/VUE
Viral	C4	37 (1.6%)	2,673.8	NYLAS/N1, LETSENCRYPT/LETSENCRYPT, and JWAGNER/SMARTCROP.JS

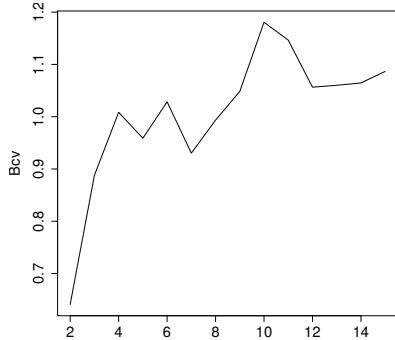


Fig. 13: β_{CV} for $2 \leq k \leq 15$

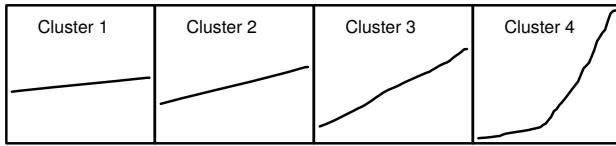


Fig. 14: Time series representing the centroids of each cluster

Slow growth is the dominant pattern, including 65.7% of the repositories in our sample, as presented in Table II. The table also shows the number of repositories in each cluster and the percentage of stars gained by the cluster's centroids in the period under analysis (52 weeks). The speed in which the stars are gained by repositories on cluster C1 is the lowest one (27.3% of new stars in one year). Moderate growth is the second pattern with more repositories (26.9% of the repositories and 94% of new stars in one year). 5.7% of the repositories have a fast growth (469.2% of new stars in the analyzed year).

The last cluster (Viral Growth) describes repositories with

a massive growth in their number of stars in a short period of time. It is a less common pattern, including 1.6% of the repositories. Figure 15 shows two examples of systems with a viral growth: NYLAS/N1 (an email client, with a peak of more than 7,300 stars in a single week) and SOUNDNODE/SOUNDNODE-APP (a desktop client for SoundCloud, which received almost 1,400 stars in a single week).

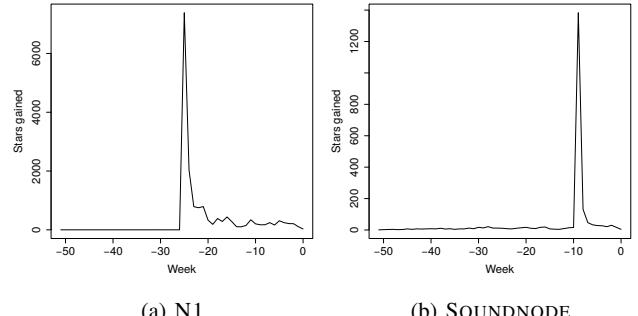


Fig. 15: Examples of viral growth

B. Growth Patterns vs Repositories Properties

Figure 16 shows the percentage of systems following the proposed growth patterns, for the top-10 programming languages in number of repositories, application domains, repository owners, and age. The three languages with the highest percentage of systems with slow growth are Ruby (92%), CSS (82%), and HTML (79%). By contrast, the languages with the highest percentage of systems with fast growth are Go (7.6%) and Java (7.6%). Go is a new language that is attracting a lot of interest.⁶ Regarding Java, 61 out of 95 repositories

⁶<https://www.thoughtworks.com/radar/languages-and-frameworks/go-language>, verified on 04/07/2016.

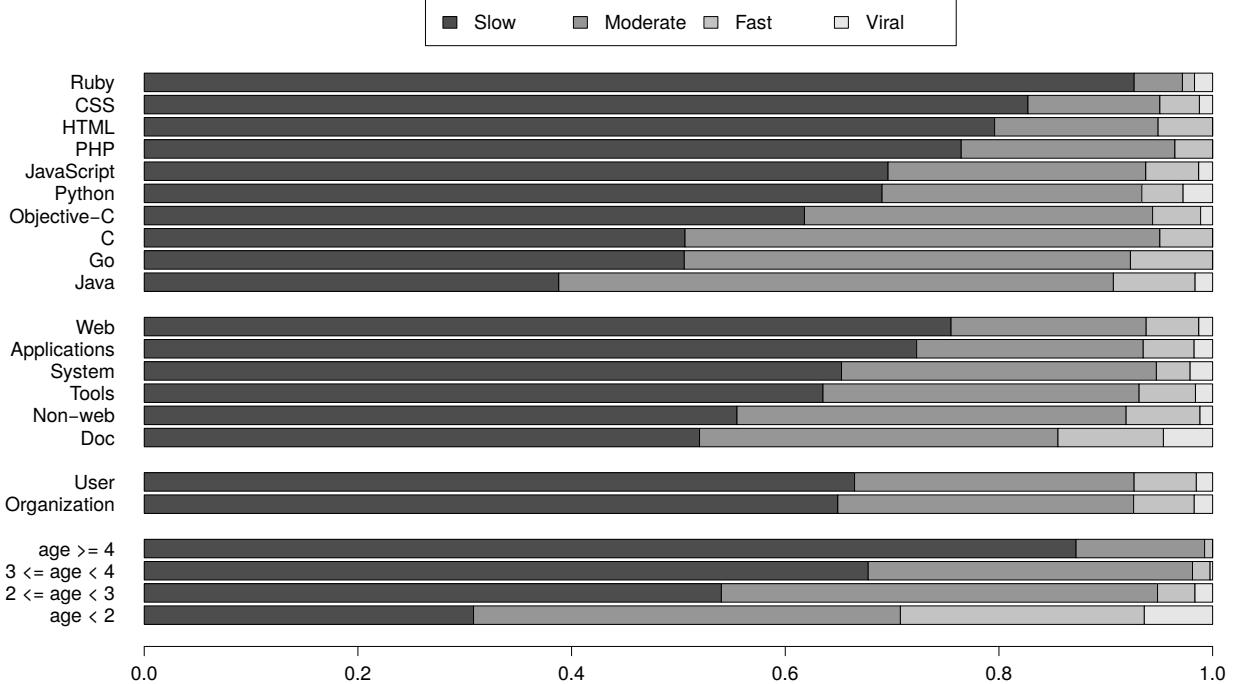


Fig. 16: Percentage of systems following the proposed growth patterns, for the most popular programming languages, application domains, and repository owners

with fast growth are Android applications. JavaScript is the language with the highest number of repositories with viral growth (10 repositories), followed by C++ (5 repositories) and Python (5 repositories). In relative terms, 2.7% of the Python systems have a viral growth, followed by 1.6% of the systems implemented in Ruby. When we group the systems by application domain, 75% of the web libraries and frameworks have a slow growth. Interestingly, the two domains with the highest percentage of systems following a fast growth are documentation (9.8%) and non-web libraries and frameworks (6.9%). Regarding the repository owners, there is no substantial difference between users and organizations. For slow growth, the percentage of systems is 66.4% and 64.9%, for users and organizations, respectively. For fast growth, the percentage is 5.8% and 5.6%, respectively. Finally, the last bars in Figure 16 show that old repositories tend to present a slow growth. The percentage of such repositories ranges from 30.8% (age < 2 years) to 87.2% (age ≥ 4 years).

As mentioned, we found a high percentage of web frameworks and libraries—especially the ones implemented in Ruby, CSS, and HTML—with a slow growth. We hypothesize two main reasons to explain this result. First, web libraries and frameworks are the dominant applications in our dataset of popular applications (837 repositories, 33%). This implies in a high competition, with many systems disputing the same users. For example, we found a list of JavaScript MVC-based frameworks with slow growth, in-

cluding systems like KNOCKOUT/KNOCKOUT, SPINE/SPINE, QUIRKEY/SAMMY, and SPROUTCORE/SPROUTCORE. These systems have to compete with “blockbusters”, like ANGULAR/ANGULAR.JS, which is certainly a challenging task. The second reason is that there are many highly popular web frameworks and libraries in our dataset. For example, among the top-10% repositories in number of stars, 42.8% are web libraries and frameworks. We cannot assume that these systems will present the same growth rates of less popular ones. For instance, if ANGULAR/ANGULAR.JS starts to grow at 469.2% per year (the growth rate observed for the centroid of the repositories with fast growth) it will have almost 1.5M stars in two years.

V. FEEDBACK FROM DEVELOPERS

We contacted the main developers of some GitHub repositories to clarify the results presented in the previous sections. Specifically, we surveyed developers about three themes: (a) the impact on popularity of repositories owned by users (Section V-A); (b) the main characteristics of successful releases (Section V-B); (c) the reasons for the peaks of popularity observed in systems with viral growth (Section V-C). The surveys were performed by means of follow-up emails.

A. Impact on Popularity of Repositories Owned by Users

In Section III (RQ #1), we found that repositories owned by organizations are more popular than the ones owned by individuals. For example, among the top-100 most popular

repositories, only 30 repositories are owned by users. The developers of 17 of such systems have a public mail address in their GitHub profile. We sent a short survey to these developers and received responses from five of them (29.8%). In this survey, we asked two questions. First, we asked the developers about possible plans to migrate their repositories to an organization account. All developers answered negatively this question. Two developers mentioned they want to explicitly appear as the repository owner, like in this answer:

“I worked hard to create the project, and having it under my personal username is necessary to have proper credit for it.”

To complement the first question, we asked the developers if they agree that migrating the repositories to an organization account would help to attract more users. Four developers (80%) answered negatively to this question and only one participant provided the following answer:

“It depends on what organization it is. If it’s a well known org I’m sure it helps, otherwise I don’t think it makes a difference.”

Therefore, although it seems “easier” to organizations to reach the top positions of GitHub popularity ranking, some systems owned by individual developers also reach these positions. These developers usually do not want to move to organizational accounts, basically to keep full control and credit for their repositories.

B. Characteristics of Successful Releases

To reveal the characteristics of the most successful releases in our dataset (see RQ #4, Section III), we perform a survey with the main developers of 60 releases with the highest fraction of stars gained on the week after the release (and whose developers have a public mail address on their GitHub profile). We received answers from 25 developers, which corresponds to a response ratio of 41.6%. First, we asked the developers about the type of features implemented in these releases. As presented in Table III, the releases usually include both functional and non-functional requirements (14 answers), followed by releases with mostly functional requirements (9 answers). We did not receive answers about releases including non-functional requirements. Two developers provide other types of answers (“complete rewrite” and “maintenance release”, respectively).

TABLE III: Features implemented in successful releases

Features	Answers
Both functional and non-functional	14
Mostly functional	9
Other answers	2
Mostly non-functional	0

We also asked the developers to explain how the features implemented in these releases were selected (answers including multiple items are possible in this question). As presented in Table IV, the features usually come from ideas of the repository’ maintainers (23 answers) and from user’s suggestions (11 answers).

TABLE IV: How the features are selected?

Features selected from	Answers
Ideas of the repository maintainers	23
Users suggestions	11
Features of similar projects	6
Other answers	3

C. Reasons for Viral Growth

To expose the reasons for viral growth, we sent a message to the main developer of 22 systems with viral growth and who have a public mail address on their GitHub profile. In the message, we asked the developers to explain the peaks observed in the number of stars of their repositories. We received answers from 14 developers, which corresponds to a response ratio of 63%. As presented in Table V, 11 developers (78.5%) linked the peaks to posts in social media sites, mostly Hacker News.⁷ For example, we received the following answer:

“I posted about this project on HackerNews. It quickly got a lot of attention and remained on the front page of HackerNews (a very high traffic tech site) for over 24 hours. It subsequently made it onto the github.com/explore as one of the top starred repositories for around a week. Because the repo was highlighted in these two high-profile locations for so much time, it received an incredible amount of traffic, which translated to a considerable number of stars.”

TABLE V: Sources of popularity

Source	Answers
Social media sites (e.g., HackerNews)	11
Blogs and news sites (e.g., infoq.com)	3
Other answers (e.g., private mailing list)	4

VI. THREATS TO VALIDITY

Number of stars as a proxy for popularity: In the paper, we consider that stars are proxies for a project popularity, as common in studies about the popularity of social media content [7]–[10]. However, a developer can star a repository for other reasons, for example, when she in fact finds problems in the system and wants to create a bookmark for later access and analysis.

Dataset. GitHub has millions of repositories. We build our dataset by collecting the top-2,500 repositories with more stars, which represents a small fraction in comparison to the GitHub’s universe. However, our goal is exactly to investigate the popularity of the most starred repositories. Furthermore, most GitHub repositories are forks and have very low activity [15], [16].

Application domains. Because GitHub does not classify the hosted applications in domains, we performed this classification manually. Therefore, it is subjected to errors and

⁷<https://news.ycombinator.com/>

inaccuracies. To mitigate this threat, the dubious classification decisions were discussed by two paper's authors.

Growth patterns. The selection of the number of clusters is a key parameter in algorithms like KSC. To mitigate this threat, we employed a heuristic that considers the intra/intercluster distance variations [14]. Furthermore, the analysis of growth patterns was based on the stars obtained on the last year. The stars before this period are not considered, since the KSC algorithm requires time series with the same length.

VII. RELATED WORK

Several studies examine the relationship between popularity of mobile apps and their code properties [17]–[25]. Yuan et al. investigate 28 factors along eight dimensions to understand how high-rated Android applications are different from low-rated ones [22]. Their result shows that external factors, like number of promotional images, are the most influential factors. Guerrouj and Baysal explore the relationships between mobile apps' success and API quality [26]. They found that changes and bugs in API methods are not strong predictors of apps' popularity. Ruiz et al. examine the relationship between the number of ad libraries and app's user ratings [20]. Their results show that there is no relationship between the number of ad libraries in an app and its rating. Linares-Vásquez et al. investigate how the fault- and change-proneness of Android API elements relate to applications' lack of success [19]. They state that making heavy use of fault- and change-prone APIs can negatively impact the success of these apps.

Other studies examine source code repositories in order to understand what makes a project popular. Weber and Luo attempt to differentiate popular and unpopular Python projects on GitHub using machine learning techniques [5]. They found that in-code features are more important than author metadata features. Zho et al. study the frequency of folders used by 140 thousands GitHub projects and the results suggest that the use of standard folders (e.g., doc, test, examples) may have an impact on project popularity [27]. Bissyande et al. analyze the popularity, interoperability, and impact of various programming languages, using a dataset of 100K open source software projects [28]. Aggarwal et al. study the effect of social interactions on GitHub projects' documentation [6]. They conclude that popular projects tend to attract more documentation collaborators. By analyzing usage of Java APIs, Mileva states that popularity trend is a method for displaying the users preferences and for predicting their future [29].

Finally, other studies analyze the relationship between popularity and software quality. Sajnani et. al. study the relationship between component popularity and component quality in Maven [30], finding that, in most cases, there is no correlation. Capra et. al. evaluate the effect of firms' participation on communities of open source projects and conclude that firms' involvement improves the popularity, but leads to lower software quality [31].

To our knowledge, we are the first study to track popularity over time on social code sharing sites, like GitHub.

However, there are similar studies in other contexts, like App Stores [23], video sharing sites [32], and social platforms [10]. Chatzopoulou et al. [32] analyze popularity of YouTube videos by looking at properties and patterns metrics. They report that many of the popularity metrics are highly correlated. In our study, we also report correlations between stars and other popularity metrics (e.g., forks). Lehmann et al. [9] analyze popularity peaks of hashtags in Twitter. They found four usage patterns restricted to a two-week period centered on the peak time whereas the popularity patterns presented in this study are based on the last year data.

VIII. CONCLUSION

In this paper, we first studied the popularity of GitHub repositories aiming to answer four research questions. We concluded that three most common domains on GitHub are web libraries and frameworks, non-web libraries and frameworks, and software tools. However, the three domains whose repositories have more stars are systems software, web libraries and frameworks, and documentation. Additionally, we found that repositories owned by organizations are more popular than the ones owned by individuals (RQ #1). We also reported the existence of a strong correlation between stars and forks, a weak correlation between stars and commits, and a weak correlation between stars and contributors (RQ #2), confirming the importance of a large base of contributors to the success of open source software [33]. We concluded that repositories have a tendency to receive more stars right after their first public release. After this period, for half of the repositories the growth rate tends to stabilize (RQ #3). In other words, bursts of popularity do not explain the popularity growth of most repositories. We showed that there is an acceleration in the number of stars gained just after releases (RQ #4), which confirms the importance of developers constantly evolving and improving their systems.

We identified four patterns of popularity growth, which were derived after clustering the time series that describe the number of stars of the systems in our dataset. We found that slow growth is the most common pattern (65.7%) and that very few systems present a viral behavior (1.6%). Slow growth is more common in case of overpopulated application domains (as web libraries and frameworks) and for old repositories.

As future work, we plan to investigate repositories that are not popular yet and compare them with the popular ones. We also plan to correlate repository and language popularity to provide relative measures of popularity. For example, if we restrict the analysis to developers from a given language, a Scala repository can be considered more popular than a JavaScript one, although having less stars. Moreover, we plan to investigate models for predicting software popularity, which can be used for example to warn developers when signs of stagnation are detected in their repositories.

ACKNOWLEDGMENTS

This research is supported by FAPEMIG and CNPq.

REFERENCES

- [1] G. Gousios, M. Pinzger, and A. van Deursen, "An exploratory study of the pull-based software development model," in *36th International Conference on Software Engineering (ICSE)*, 2014, pp. 345–355.
- [2] G. Gousios, A. Zaidman, M.-A. Storey, and van Arie Deursen, "Work practices and challenges in pull-based development: the integrator's perspective," in *37th IEEE International Conference on Software Engineering (ICSE)*, 2015, pp. 358–368.
- [3] Y. Yu, H. Wang, V. Filkov, P. Devanbu, and B. Vasilescu, "Wait for it: determinants of pull request evaluation latency on GitHub," in *12th Working Conference on Mining Software Repositories (MSR)*, 2015, pp. 367–371.
- [4] G. Gousios, M.-A. Storey, and A. Bacchelli, "Work practices and challenges in pull-based development: the contributor's perspective," in *38th International Conference on Software Engineering (ICSE)*, 2016, pp. 1–12.
- [5] S. Weber and J. Luo, "What makes an open source code popular on GitHub?" in *13th IEEE International Conference on Data Mining Workshop (ICDW)*, 2014, pp. 851–855.
- [6] K. Aggarwal, A. Hindle, and E. Stroulia, "Co-evolution of project documentation and popularity within GitHub," in *11th Working Conference on Mining Software Repositories (MSR)*, 2014, pp. 360–363.
- [7] M. Ahmed, S. Spagna, F. Huici, and S. Niccolini, "A peek into the future: predicting the evolution of popularity in user generated content," in *6th International Conference on Web Search and Data Mining (WSDM)*, 2013, pp. 607–616.
- [8] F. Figueiredo, J. M. Almeida, G. A. alves, and F. Benevenuto, "On the dynamics of social media popularity," *ACM Transactions on Internet Technology (TOIT)*, vol. 14, no. 4, pp. 1–23, 2014.
- [9] J. Lehmann, G. alves, J. J. Ramasco, and C. Cattuto, "Dynamical classes of collective attention in Twitter," in *21st International Conference on World Wide Web (WWW)*, 2012, pp. 251–260.
- [10] Z. Ma, A. Sun, and G. Cong, "On predicting the popularity of newly emerging hashtags in Twitter," *Journal of the American Society for Information Science and Technology*, vol. 64, no. 7, pp. 1399–1410, 2013.
- [11] J. Yang and J. Leskovec, "Patterns of temporal variation in online media," *Time*, vol. 468, no. 1, pp. 177–186, 2011.
- [12] F. Figueiredo, "On the prediction of popularity of trends and hits for user generated videos," in *6th international Conference on Web Search and Data Mining (WSDM)*, 2013, pp. 741–746.
- [13] J. A. Hartigan, *Clustering algorithms*. John Wiley & Sons, Inc., 1975.
- [14] D. A. Menasce and V. Almeida, *Capacity Planning for Web Services: Metrics, Models, and Methods*, 1st ed. Prentice Hall, 2001.
- [15] E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. M. German, and D. Damian, "The promises and perils of mining GitHub," in *11th Working Conference on Mining Software Repositories (MSR)*, 2014, pp. 92–101.
- [16] ———, "An in-depth study of the promises and perils of mining GitHub," *Empirical Software Engineering*, pp. 1–37, 2015.
- [17] D. Datta and S. Kajanan, "Do app launch times impact their subsequent commercial success? an analytical approach," in *International Conference on Cloud Computing and Big Data (CloudCom-Asia)*, 2013, pp. 205–210.
- [18] B. Fu, J. Lin, L. Li, C. Faloutsos, J. Hong, and N. Sadeh, "Why people hate your app: making sense of user feedback in a mobile app store," in *19th International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2013, pp. 1276–1284.
- [19] L.-V. squeez, G. Bavota, B.-C. rdenas, D. M. Penta, R. Oliveto, and D. Poshyvanyk, "API change and fault proneness: a threat to the success of Android apps," in *9th Foundations of Software Engineering (FSE)*, 2013, pp. 477–487.
- [20] I. J. Mojica Ruiz, M. Nagappan, B. Adams, T. Berger, S. Dienst, and A. E. Hassan, "Impact of ad libraries on ratings of Android mobile apps," *IEEE Software*, vol. 31, no. 6, pp. 86–92, 2014.
- [21] G. Lee and T. Raghu, "Determinants of mobile apps' success: evidence from the app store market," *Journal of Management Information Systems*, vol. 31, no. 2, pp. 133–170, 2014.
- [22] Y. Tian, M. Nagappan, D. Lo, and A. E. Hassan, "What are the characteristics of high-rated apps? a case study on free Android applications," in *31st International Conference on Software Maintenance and Evolution (ICSME)*, 2015, pp. 1–10.
- [23] L. Guerrouj, S. Azad, and P. C. Rigby, "The influence of app churn on app success and StackOverflow discussions," in *22nd International Conference on Software Analysis, Evolution and Reengineering (SANER)*, 2015, pp. 321–330.
- [24] F. Palomba, L.-M. Vasquez, G. Bavota, R. Oliveto, D. M. Penta, D. Poshyvanyk, and D. A. Lucia, "User reviews matter! tracking crowd-sourced reviews to support evolution of successful apps," in *International Conference on Software Maintenance and Evolution (ICSME)*, 2015, pp. 291–300.
- [25] L. Corral and I. Fronza, "Better code for better apps: a study on source code quality and market success of Android applications," in *2nd International Conference on Mobile Software Engineering and Systems (MOBILESoft)*, 2015, pp. 22–32.
- [26] L. Guerrouj and O. Baysal, "Investigating the Android apps' success: an empirical study," in *24th International Conference on Program Comprehension (ICPC)*, 2016, pp. 1–4.
- [27] J. Zhu, M. Zhou, and A. Mockus, "Patterns of folder use and project popularity: A case study of GitHub repositories," in *8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ISEM)*, 2014, pp. 30:1–30:4.
- [28] T. F. Bissyande, F. Thung, D. Lo, L. Jiang, and L. Reveillere, "Popularity, interoperability, and impact of programming languages in 100,000 open source projects," in *37th Annual International Computer Software and Applications Conference (COMPSAC)*, 2013, pp. 303–312.
- [29] Y. M. Mileva, "Mining the evolution of software component usage," Ph.D. dissertation, Saarland University, 2012.
- [30] H. Sajnani, V. Saini, J. Ossher, and C. V. Lopes, "Is popularity a measure of quality? an analysis of Maven components," in *30th Software Maintenance and Evolution (ICSME)*, 2014, pp. 231–240.
- [31] E. Capra, C. Francalanci, F. Merlo, and C. Rossi-Lamastra, "Firms involvement in open source projects: a trade-off between software structural quality and popularity," *Journal of Systems and Software*, vol. 84, no. 1, pp. 144 – 161, 2011.
- [32] G. Chatzopoulou, C. Sheng, and M. Faloutsos, "A first step towards understanding popularity in YouTube," in *30th IEEE International Conference on Computer Communications Workshops (INFOCOM)*, 2010, pp. 1–6.
- [33] A. Mockus, R. T. Fielding, and J. D. Herbsleb, "Two case studies of open source software development: Apache and Mozilla," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 11, no. 3, pp. 309–346, 2002.