

Pull Request Decisions Explained: An Empirical Overview

Xunhui Zhang, Yue Yu*, Georgios Gousios, and Ayushi Rastogi

Abstract—*Context:* The pull-based development model is widely used in open source projects, leading to the emergence of trends in distributed software development. One aspect that has garnered significant attention concerning pull request decisions is the identification of explanatory factors. *Objective:* This study builds on a decade of research on pull request decisions and provides further insights. We empirically investigate how factors influence pull request decisions and the scenarios that change the influence of such factors. *Method:* We identify factors influencing pull request decisions on GitHub through a systematic literature review and infer them by mining archival data. We collect a total of 3,347,937 pull requests with 95 features from 11,230 diverse projects on GitHub. Using these data, we explore the relations among the factors and build mixed effects logistic regression models to empirically explain pull request decisions. *Results:* Our study shows that a small number of factors explain pull request decisions, with that concerning whether the integrator is the same as or different from the submitter being the most important factor. We also note that the influence of factors on pull request decisions change with a change in context; *e.g.*, the area hotness of pull request is important only in the early stage of project development, however it becomes unimportant for pull request decisions as projects become mature.

Index Terms—pull-based development, pull request decision, distributed software development, GitHub



1 INTRODUCTION

THE PULL-BASED development model is an important paradigm for global collaboration in open source projects. In this model [1], contributors (*also known as requesters and submitters*) submit their proposed code changes to a base repository by creating a pull request from their cloned repository for the reviewers to inspect. The integrator (*also known as the closer and the merger*) evaluates the proposed changes and decides whether to accept or reject the pull request. However, this process is made complex by additional actors and mechanisms. For instance, during the review, anyone can discuss the feature(s), correctness, etc., of the pull request. Moreover, DevOps tools that automatically check code adaptability and provide results to contributors and integrators exist.

Many studies on understanding pull-based development have emerged in recent years to improve developer contributions, balance integrators' workloads, optimize review processes, etc. There are studies on pull request decisions [2], their latency [3], reviewer recommendations [4], [5], the duplication of pull requests [6], [7], the automatic generation of pull request descriptions [8], and the prioritization of pull request lists [9], among others. This study focuses on explaining pull request decisions.

* is the corresponding author

- X. Zhang and Y. Yu are with the National Key Lab of Parallel Distribution, National University of Defense Technology, Changsha, Hunan 410073, China. E-mail: {zhangxunhui, yuyue}@nudt.edu.cn.
- A. Rastogi is with the Faculty of Science and Engineering, the University of Groningen, The Netherlands. A part of the work was performed while the author was affiliated to TU Delft. E-mail: a.rastogi@rug.nl.
- G. Gousios is with the TU Delft. E-mail: g.gousios@tudelft.nl.

Many studies have made strides in explaining pull request decisions by introducing new factors in the past decade. Some examples of these factors are continuous integration (CI) [10], [11], geographical location [12], and bot usage [13], [14]. Relatedly, a few studies have presented a list of factors that can influence pull request decisions. One outstanding work along this line of Gousios et al. [1] provided a list of developer, project, and pull request characteristics. Tsay et al. [15] split factors into two categories, *i.e.*, social- and technical-related factors. A more recent study by Dey et al. [16] combined many such factors (50) to rank their importance for prediction.

While several studies have contributed individual pieces to understand pull request decisions, a systematic synthesis of the body of knowledge to explain such decisions is missing. If new mechanisms emerge and a new set of factors occurs. Researchers need to decide which factors are more critical when selecting control variables for an empirical study to find their impact on pull request decisions. However, there lack relevant studies to tell them how to make choices. Also, understanding factors' influence in different contexts is essential for researchers to select projects and factors. From developers' perspectives, when creating predictive tools, it is also important to consider the impact of different contexts. *E.g.*, how to choose factors if reviewers comment during the review process? What factors should be considered if a pull request uses CI tools? Factors, if properly selected, not only maintain accuracy but also significantly improve the efficiency of decision prediction. Therefore, our current work presents an empirical investigation explaining pull request decisions from GitHub in terms of the factors known to influence them. Particularly, we explore the following two research questions:

RQ1 How do these factors influence pull request decisions?

RQ2 *How do the factors influencing pull request decisions change with a change in context?*

First, we conduct a systematic literature review (SLR) to identify a comprehensive list of factors known to influence pull request decisions. Then, we create a large and diverse dataset of pull requests and factors (or their indicators) that can be mined from archival software data. Finally, we build models (mixed effects logistic regression models) that suggest the relations between each factor and pull request decisions in general, specific scenarios (e.g., when pull requests use CI), and different contexts (e.g., the time when pull requests are closed).

This paper makes the following contributions to software engineering research and practice:

- 1) We present a curated dataset of 11,230 projects on GitHub with 95 factors and 3,347,937 pull requests. Our dataset is diverse in terms of the number of contributors, programming language, and activities (see Table 1). It also covers the entire project lifecycle as a representation of diversity in time. Future researchers can use and extend our large and rich dataset¹ to conduct deeper investigations and use scripts to replicate the results.²
- 2) We present a synthesis of the factors identified in the literature, indicating their significance and direction.
- 3) We show the importance of these factors in explaining pull request decisions and how these decisions change with a change in context.

The rest of the paper is organized as follows. In Section 2, we explain our research design. In Section 3, we present the results. In Section 4, we conduct a case study about affiliation-related factors. We discuss the implications in Section 5 and present the threats in Section 6. In Section 7, we describe the related work of this study. In Section 8, we present our conclusions and directions for future work.

2 STUDY DESIGN

The framework of our study is shown in Figure 1, which mainly comprises four parts presenting the steps to empirically explain pull request decisions. First, we gather a comprehensive list of the factors known to influence pull request decisions (see the SLR part in Figure 1). Next, we collect data from diverse collaboratively developed software projects on GitHub to use as proxies for the factors identified above (see the Data Collection part in Figure 1). Then, we transform the data and transfer them into a form usable for analysis (see the Data Preprocessing part in Figure 1). Finally, we model the data to answer our research questions, starting with an exploratory data analysis (see the Statistical Modeling part in Figure 1).

2.1 Systematic literature review

To collect all factors known to influence pull request decisions, we conducted a systematic literature review (see the SLR part in Figure 1(a)), which was based on the guidelines from Kitchenham et al. [17].

Our search strategy was to identify all scientific articles relating to pull request decisions. We selected two widely used search terms, “pull request” and “pull based”, which are often used interchangeably as pull request models, pull-based development, and similar variants. We combined the two search terms with a logical “OR” operator (i.e., “pull request” OR “pull based”) defining our search space. We searched for (“pull request” OR “pull based”) on Google Scholar, ACM Digital Library, IEEEExplore, Web of Science and Ei Compendex, resulting in a total of 3,941 papers. We ran the query on April 17th, 2020. We identified 1,000 papers from Google Scholar, 1,433 from ACM Digital Library, 352 from IEEEExplore, 487 from Web of Science, and 669 papers from Ei Compendex. We performed an additional step of searching Google Scholar for papers published only in 2020. (Here, we only consider 2020 because we can get all relevant papers through the backward snowballing process [17]. Therefore, we don’t have to perform searches for each year.) This step was necessary since Google Scholar retrieves only the top 1,000 results, which means that it is likely to miss many articles [18], [19]. The additional search (also conducted on April 17th, 2020) resulted in 610 more papers, leading to a total of 4,551 papers for backward snowballing.

To identify the factors influencing pull request decisions, the first author manually analyzed the title and abstract of each paper and selected all studies presenting all the factors influencing pull request decisions that can be inferred by mining software archives. The search resulted in 19 papers after excluding papers for the following reasons:

- they were written in languages other than English (45 papers)
- they were duplicates (1,181 papers)
- they were initial versions of the papers when extended versions were available (12 papers)
- they presented factors not applicable to GitHub (5 papers); e.g., *a study on Firefox and Mozilla core projects shows that “bug severity” and “bug priority” influence patch acceptance [20]. These attributes do not exist on GitHub*
- they were not related to pull request decisions (3,277 papers)
- they were related to pull request decisions but difficult to reproduce (4 papers), e.g., *using medical equipment to track the eyes of reviewers [21]*
- they included factors not generalizable to a wider range of software projects on GitHub (4 papers), e.g., *labels [22] that vary across communities*
- they presented different operationalizations of related concepts (3 papers); e.g., *emotions can be measured directly as joy, love, sadness, and anger; indirectly via valence, arousal, and dominance [23]; and abstractly based on polarity [24]. We choose one of three representations of emotions, i.e., polarity. As another example, Calefato et al. [25] measured trust using agreeableness, one of the five personality traits used by Iyer et al. [2]. Thus, we chose five personality traits*
- they presented factors not measurable quantitatively (1 paper), i.e., *the features relating to pull request decisions found in a qualitative study [26]*

Next, we identified other relevant articles by considering the references of the 19 selected seed articles. We applied

1. <https://zenodo.org/record/4837134#.YLEWY3isdW>
 2. https://github.com/zhangxunhui/TSE_pull-based-development

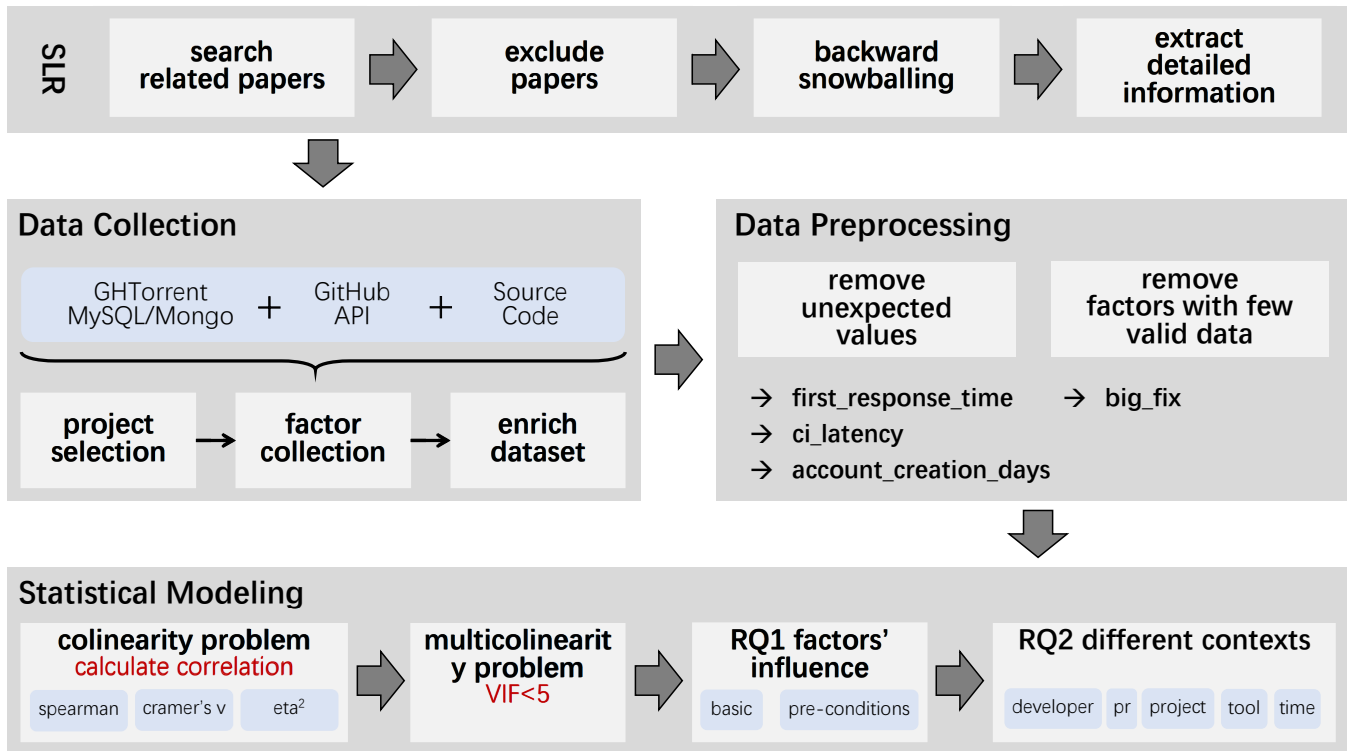


Fig. 1: Framework of this paper

173 the backward snowballing method [17] twice, meaning that
 174 we selected (a) the references of the 19 articles and (b) the
 175 references of the references. After two rounds, we did not
 176 find any new related papers. This process resulted in 7 new
 177 papers, bringing the total to 26 papers presenting the factors
 178 related to pull request decisions.

179 An overview of the 94 features (the factor *same_user* was
 180 not considered in previous studies) found in the systematic
 181 literature review is shown in Table 2, which lists the sym-
 182 bolic representations of the features in columns 1 and 3, fol-
 183 lowed by their descriptions in columns 2 and 4, respectively.
 184 All the features are classified as developer, project, and
 185 pull request characteristics. Furthermore, Table 10 shows
 186 the relations between each of the factors and pull request
 187 decisions, as identified in the 26 selected research articles.

188 For the accuracy and validity of the data extraction
 189 process, the first and the last author did the whole process
 190 together. First, in the paper screening phase, the first author
 191 got the initial results. Then the first author and the last
 192 author met to discuss the paper with uncertainty and finally
 193 reached an agreement. *E.g.*, the paper [26] was a relevant
 194 study on pull request decisions, but as a qualitative study,
 195 it lacked a measure of certainty about the relevant factors,
 196 so we removed the paper. After that, in the factor extraction
 197 stage, the first author extracted the initial factors, including
 198 the name of the factor, the related description, the category
 199 to which it belongs (pull request, project, or developer), and
 200 the description of related findings, forming a list. The first
 201 and the last author then met to discuss and agree on the
 202 information in the list, which consisted of the following
 203 steps.

204 1) For relevant factors with unclear descriptions, reach an

- agreement, *e.g.*, factor *pushed_delta* (see Table 2). 205
 2) Remove factors that are not applicable for GitHub, *e.g.*, 206
 bug severity. 207
 3) Remove factors that are difficult to reproduce, *e.g.*, eye 208
 tracking of reviewers. 209
 4) Confirm the category to which factors belong. 210
 5) The last author maintained a list of relevant factors in 211
 advance based on the research experience and checked 212
 during the meeting to see if they all appeared in the list 213
 provided by the first author. 214

After the above process, we finally identified the 94 relevant 215
 factors. 216

2.2 Data collection 217

We collected data on a variety of software projects hosted 218
 on GitHub as a proxy for the factors identified above. The 219
 dataset used for this study came from our prior work [27], 220
 featuring 96 factors collected from 11,230 projects. Further- 221
 more, we enriched the dataset with missing factors and 222
 values (see the Data Collection part in Figure 1). 223

Our initial dataset [27] was built on the publicly available 224
 GHTorrent MySQL data dump dated June 1st, 2019.³ It 225
 features 96 factors relating to pull requests, developers, 226
 or projects (derived from 76 research articles published 227
 between 2009 and 2019) for 11,230 software projects. The 228
 screening steps of GitHub projects are summarized as fol- 229
 lows: 230

- 1) Filter forked or deleted repositories based on 231
 GHTorrent.³ 232

3. <http://ghtorrent-downloads.ewi.tudelft.nl/mysql/mysql-2019-06-01.tar.gz>

- 2) Filter repositories that do not have any pull requests in the last three months.
- 3) Select projects from six programming languages (as against 4 programming languages in the case of Gousios et al.'s [28] dataset). The extended JavaScript and Go languages are the most popular programming languages on Github⁴ and the fastest growing programming languages in recent years, respectively.⁵
- 4) Select all projects with at least 33 submitted pull requests. These projects constitute the top 3% of all projects in terms of pull request count (as against the top 1% in the case of Gousios et al.'s [28] dataset). The top 3% here is chosen to make some extensions based on Gousios et al.' dataset [28]. With the development of Github, a large number of open source projects have emerged. In addition to the most active open source projects, we also want to include a wide range of projects, including small and relatively less active projects. After discussion, we have chosen the top 3% of projects.
- 5) Split projects according to the tertile thresholds of the number of developers in the project, *i.e.*, small-sized teams (low tertile) with 12 or fewer developers, medium-sized teams (middle tertile) with 13 and up to 30 developers, large-sized teams (high tertile) with more than 30 developers. Randomly select 4,000 projects in each class.
- 6) Remove the data holding project "everypolitician/everypolitician-data", which is extremely large, and we lack the ability to collect related factors.
- 7) After discussion among authors, remove projects with less than 20 closed pull requests related to their default branch to ensure enough data for the subsequent steps required in research.

After the above steps, 11,230 projects remained, which offers a total of 3,347,937 closed pull requests (meaning a decision has been made) submitted to the repository's default branch.

TABLE 1: Description of project diversity

category	type	project count	percentage
language	JavaScript	3,879	34.5%
	Python	3,055	27.2%
	Java	1,823	16.2%
	Ruby	1,243	11.1%
	Go	913	8.1%
	Scala	317	2.8%
project size	small ≤ 12 developers	3,711	33%
	mid ≤ 31 developers	3,634	32.4%
	large > 31 developers	3,885	34.6%
project activity	min = 33 pull requests	-	-
	25% ≤ 55 pull requests	2,843	25.3%
	50% ≤ 106 pull requests	2,796	24.9%
	75% ≤ 261 pull requests	2,791	24.9%
	max = 38,953 pull requests	-	-

Our initial dataset is futuristic and emphasizes generalizability - a design choice for a wide range of explorations [27]. Moreover, our dataset has 12 times more projects and 10 times more pull requests than Gousios et al.'s [28] dataset and is more diverse than any of the datasets of prior studies focusing on pull request decisions, which have, until now, largely focused on the most popular projects.

From Table 1, we can see that the diversity of selected projects is mainly manifested in three aspects, *i.e.*, covering 6 languages, containing different numbers of contributors, and including projects with different activity levels (the number of pull requests ranges from 33 to more than 30 thousand). Our dataset has features that are applicable to projects outside GitHub and has additional features that are likely to influence pull request development - an extrapolation of existing features.

For our analysis, we selected data related to the factors identified by our systematic literature review from the initial dataset. We noticed that 14 factors identified by our systematic literature review did not exist in the initial dataset, so we added these missing features. Table 2 presents a complete list of the factors known to influence pull request decisions on GitHub. Factors marked as * are additions to those of the initial dataset [27].

Finally, we enriched our dataset by filling in missing values wherever possible based on GHTorrent⁶, GitHub API and source code of repository. For example, the initial dataset used the tool by Vasilescu et al. [29] to infer country information. The resulting dataset, however, had a large number of missing values. We applied several steps, such as using *country_code* information and *pycountry* package⁷ to extract country names. In this way, we were able to derive the country information of an additional 546,682 contributors (1,473,008 previously), 747,204 integrators (1,580,256 previously) and 796,083 same-country participants (1,081,668 previously). The expanded country information can be seen on GitHub.⁸ To verify the validity of the data, we randomly selected 100 developers with predicted country information. Then, the first author manually checks the accuracy according to the developer's GitHub homepage and the given external site. Only two developers made a mistake in their predictions, and another two developers' country information could not be judged based on the existing knowledge. Therefore, the precision of the extracted country information $\approx 96\%$.

We added a factor, *same_user*, that did not exist in prior studies (marked as • in Table 2). While the information on the same user is not useful itself, it adds meaning to factors such as *same_country*, *same_affiliation*, and personality-difference-related factors (*e.g.*, *open_diff*), which make sense only when the contributor and integrator are not the same users. In our dataset, we found that 43.6% of the pull requests were integrated by submitters (85.7% of them were core contributors, and 14.3% were external contributors). Compared to directly committing to code repositories, pull-based development is becoming a standard collaborative model in which not only external contributors but also core

4. <https://octoverse.github.com/#top-languages-over-the-years>
 5. <https://hub.packtpub.com/why-golan-is-the-fastest-growing-language-on-github/>

6. <https://ghntorrent.org/>
 7. <https://pypi.org/project/pycountry/>
 8. https://github.com/zhangxunhui/TSE_pull-based-development/blob/master/country_info.csv

TABLE 2: Comprehensive list of the factors known to influence pull request decisions on GitHub

Factor	Description	Factor	Description
Developer Characteristics			
first_pr	first pull request? yes/no	prior_review_num	# of previous reviews in a project
core_member	core member? yes/no	first_response_time	# of minutes from pull request creation to the reviewer's first response
contrib_gender	gender? male or female	contrib_country	country of residence
same_country	same country contributor/integrator? yes/no	prior_interaction	# of interactions with a project in the last three months
same_affiliation	same affiliation contributor/integrator? yes/no	contrib/inte_affiliation	contributor/integrator affiliation
contrib/inte_X	contributor/integrator personality traits (<i>open</i> : openness; <i>cons</i> : conscientious; <i>extra</i> : extraversion; <i>agree</i> : agreeableness; <i>neur</i> : neuroticism)	perc_contrib/inte_X_emo	% of contributor/integrator (<i>neg</i> : negative/ <i>pos</i> : positive) emotion in comments
X_diff	absolute difference in the personality traits of the contributor and the integrator	contrib/inte_first_emo	emotion in contributor/integrator's first comment
social_strength	fraction of team members interacted with in the last three months	contrib_follow_integrator	contributor followed integrator before pull request creation? yes/no
followers	# of followers at pull request creation time	same_user	same contributor and integrator? yes/no
prev_pullreqs	# of previous pull requests	account_creation_days	# of days from the contributor's account creation to pull request creation
contrib_perc_commit *	% of the contributor's previous commit	requester_succ_rate	past pull request success rate
Project Characteristics			
sloc	executable lines of code	team_size	# of active core team members in the last three months
language	programming language	open_issue_num	# of open issues
project_age	# of months from project to pull request creation	open_pr_num	# of open pull requests
pushed_delta	# of seconds between two latest pull requests	fork_num	# of forks
pr_succ_rate	pull request acceptance rate of project	test_lines_per_kloc	# of test lines per 1K lines of code
stars	# of stars	integrator_availability *	latest activity of the two most active integrators
test_cases_per_kloc	# of test cases per 1K lines of code	asserts_per_kloc	# of assertions per 1K lines of code
perc_external_contribs	% of external pull request contributions		
Pull Request Characteristics			
churn_addition	# of added lines of code	churn_deletion	# of deleted lines of code
bug_fix	fixes a bug? yes/no	description_length	length of pull request description
test_inclusion	test case existing? yes/no	comment_conflict	keyword "conflict" exists in comments? yes/no
hash_tag	"#" tag exists? yes/no	num_participants	# of participants in pull request comments
lifetime_minutes	# of minutes from pull request creation to latest close time	part_num_code	# of participants in pull request and commit comments
ci_exists	uses CI? yes/no	ci_build_num	# of CI builds
ci_latency	# of minutes from pull request creation to the first CI build finish time	perc_neg_emotion	% of negative emotion in comments
num_code_comments *	# of code comments	perc_pos_emotion	% of positive emotion in comments
test_churn	# of lines of test code changed (added + deleted)	num_code_comments_con *	# of contributor's code comments
ci_test_passed	all CI builds passed? yes/no	ci_first_build_status	CI first build result
ci_failed_perc	% of CI builds failed	ci_last_build_status	CI last build status
num_commits	# of commits	src_churn	# of lines changed (added + deleted)
files_added	# of files added	files_deleted	# of files deleted
files_changed	# of files touched	Friday_effect *	pull request submitted on a Friday? yes/no
reopen_or_not *	pull request is reopened? yes/no	commits_on_files_touched	# of commits on files touched
has_comments *	pull request has a comment? yes/no	num_comments	# of comments
has_participants *	has a participant? yes/no	core_comment *	has a core member comment? yes/no
contrib_comment *	has a contributor comment? yes/no	inte_comment *	has an integrator comment? yes/no
has_exchange *	has contributor and integrator comments? yes/no	other_comment *	has noncontributor/core team comment? yes/no
num_comments_con *	# of contributor comments	at_tag	"@" tag exists? yes/no

NOTE: Factors marked as * are additions of our study to the latest MSR Data Showcase pull request dataset [27], while • are additions to previous studies. All metrics are relative to a referenced pull request in a project. Factors that change over time (e.g., core team) are measured using the previous three months of development activities in a project. The related paper information and the nature of each factor can be seen in Table 10.

329 members are interested. Therefore, it is necessary to add this
 330 factor and study its influence on pull request decisions.

331 For factor *bug_fix*, we followed Fan et al.'s [30] method
 332 in finding the tag for determining whether the pull request
 333 is a bug fix or not. In their method, they manually found the
 334 most used tags for bug-prone and non-bug-prone issues.
 335 (The tags are listed in Table 3.) Therefore, we first check
 336 whether the pull request has a tag marking its type. If not,
 337 we link the pull request to an issue [31]. If the pull request
 338 fixes an issue, we check the related issue's tag to see whether
 339 the pull request fixes a bug or not. To ensure data accuracy,
 340 we did not use a prediction model to predict the type of pull
 341 request.

TABLE 3: Bug and non-bug tags

Category	Tags
Bug	"bug"; "defect"; "type:bug"
Non-bug	"enhancement"; "feature"; "question"; "feature request"; "documentation"; "improvement"; "docs"

2.3 Data preprocessing

Our exploration of the resulting dataset (manually and using data distribution graphs) showed some unexpected data values for factors such as *first_response_time*, *ci_latency*, *account_creation_days* and *project_age*. It is important to fix them for reliable inferences (see the Technical Report [32] for examples).

- *first_response_time* has *negative* values for some pull requests. One possible reason is that our metric considers the discussion under a pull request and the comments under the related code. Since some comments exist before pull request creation, our data show negative values. We fix this issue by excluding pull requests with negative values (0.4%).
- *ci_latency* has *negative* values for some pull requests. CI latency measures the time from pull request creation to CI build finish time. In some cases, however, commits exist prior to pull request creation, and the time of first build recorded is earlier than the creation time of a pull request. We fix this problem by removing such pull requests (1.5%).
- *account_creation_days* and *project_age* have *negative* values, which happens in special cases where the creation time of a user account on GHTorrent is different from that on Github API. Here too, we remove such cases (0.1%).
- *bug_fix* has 99.3% empty values. We remove this factor, which otherwise can adversely affect the analysis.

For the time-related factors, we verified the accuracy of the remaining data by randomly selecting 100 records. We found that the inconsistency between the GHTorrent MySQL version and GitHub API resulted in the accuracy of *first_response_time*, *account_creation_days*, *project_age*, and *ci_latency* at about 98%, 97%, 96%, and 94%, respectively. We have added this part to the Threats to Validity section.

2.4 Statistical modeling

Presenting a comprehensive analysis of the factors influencing pull request decisions, we build generic models comprising all the factors and models representing specific cases. We also build models within different contexts. However, first, we explore relationships among the factors identified above.

Our preliminary exploration into the relationship among factors started with calculating the correlations among all the factors. We calculated the Spearman correlation coefficient (ρ) for continuous factors [1], Cramér's V (Φ_c) for categorical factors [33], and partial Eta-squared (η^2) for the correlation between continuous and categorical factors [34]. We consider $\rho > 0.7$ [1], $\Phi_c > \frac{0.5}{df}$ [35] and $\eta^2 > 0.14$ [35] as strong correlations.

A list of strongly correlated factors is presented in Table 4, in which the strongly correlated factors are separated from the other factors by a dotted line. For a complete list of correlations between each pair of factors, refer to our technical report [32].

Next, we built mixed effects logistic regression models to empirically explain the factors influencing pull request decisions. The models used the project identifier as a random effect, indicating similarity among the pull requests of a project [36]. All other factors had fixed effects. The resulting model indicated the significance of a factor and direction of its association with a pull request decision (accept or reject). We used the *glmer* function of the *lme4* [37] package in R to model pull request decisions.

To build an explanatory model, we included all factors that could be meaningfully added together, did not present

TABLE 4: Choices and corresponding reasons for strongly correlated factors

Correlated factors	Selected factor	Reason
test_lines_per_kloc		
test_cases_per_kloc	<i>test_lines_per_kloc</i>	previous study
asserts_per_kloc		
src_churn		
churn_addition	<i>src_churn</i>	frequency
churn_deletion		
num_comments		
at_tag		
num_participants	<i>num_comments</i>	frequency
num_comments_con		
core_member		
perc_external_contribs	<i>core_member</i>	frequency
social_strength		
requester_succ_rate		
stars		
fork_num	<i>stars</i>	frequency
inte_affiliation		
prev_pullreqs	<i>prev_pullreqs</i>	frequency
prior_interaction		
num_code_comments		
part_num_code	<i>num_code_comments</i>	frequency
num_code_comments_con		
open_pr_num	<i>open_pr_num</i>	frequency
fork_num		
ci_latency		
ci_build_num	<i>ci_latency</i>	promising performance
sloc		
language	<i>sloc</i>	promising performance
has_comments		
has_participants		
core_comment		
contrib_comment	<i>has_comments</i>	expressiveness
inte_comment		
has_exchange		
prior_review_num	<i>prior_review_num</i>	data availability
inte_affiliation		
open_issue_num	<i>open_issue_num</i>	data availability
inte_affiliation		
inte_cons	<i>inte_cons</i>	data availability
inte_affiliation		
inte_extra	<i>inte_extra</i>	data availability
inte_affiliation		
inte_agree	<i>inte_agree</i>	data availability
inte_affiliation		
same_country	<i>same_country</i>	discussion
contrib_country		
perc_contrib_pos_emo	<i>perc_contrib_pos_emo</i>	discussion
contrib_first_emo		
perc_inte_neg_emo	<i>perc_inte_neg_emo</i>	discussion
inte_first_emo		
perc_inte_pos_emo	<i>perc_inte_pos_emo</i>	discussion
inte_first_emo		
same_user		
inte_first_emo	<i>same_user</i>	discussion
inte_affiliation		
contrib_affiliation		
contrib_rate_author		
same_affiliation	<i>same_affiliation</i>	discussion
contrib_affiliation		
perc_neg_emotion		
perc_contrib_neg_emo	<i>perc_neg_emotion</i>	discussion
contrib_first_emo		
inte_first_emo		
perc_pos_emotion	<i>perc_pos_emotion</i>	discussion
inte_first_emo		
ci_failed_perc		
ci_test_passed	<i>ci_failed_perc</i>	discussion
ci_first_build_status		
ci_last_build_status		

407 the same or similar information as other factors, and were
 408 easy to interpret.

409 1) *Adding meaningful factors.* While adding factors to a
 410 model, we observed that 17 factors (postconditional
 411 factors in Table 5) did not make sense outside a specific
 412 context. For example, if the contributor and integrator
 413 were the same, then factors such as “personality dif-
 414 ference” did not exist and made no sense. We refer
 415 to such factors as “preconditional factors” and “postcon-
 416 ditional factors”. “Preconditional factors” are those that
 417 must exist for another factor to exist and make sense
 418 (e.g., *same_user* in the previous example). Conversely,
 419 “postconditional factors” are the factors in which their
 420 existence is conditional on preconditional factors (e.g.,
 421 *open_diff*). All the other factors are classified into the
 422 “others” category. A complete list of pre- and postcon-
 423 ditional factors is presented in Table 5.

424 2) *Factors presenting the same information.* Our prelimi-
 425 nary investigation showed that several factors identi-
 426 fied from the literature were strongly correlated with
 427 each other (see Table 4 for a list of strongly corre-
 428 lated factors). When two related factors were added
 429 to a model, they changed not only pull request deci-
 430 sions but also other factors, which could change
 431 the estimated effect of these factors on pull request
 432 decisions and their significance, also referred to as a
 433 multicollinearity problem [38]. To avoid multicollinear-
 434 ity, we selected one of the many strongly correlated
 435 factors. Our choice of the selection of a factor was
 436 influenced by its use in previous studies (e.g., [1] chose
 437 *test_lines_per_kloc*), frequency of occurrence in the litera-
 438 ture (e.g., *core_member* appeared most often), promising
 439 performance (indicating the likelihood of strong corre-
 440 lation with pull request decisions) (e.g., *sloc* significantly
 441 influences pull request decisions [12], while *language*
 442 does not have such a conclusion according to previous
 443 studies), expressiveness (e.g., *has_comments* is broader
 444 and more informative than *contrib_comment*), data avail-
 445 ability (e.g., *open_issue_num* has most nonempty val-
 446 ues), and otherwise in discussion with the last author
 447 (e.g., *perc_pos_emotion* is more representative for the
 448 whole review process than *inte_first_emo*; *same_country*
 449 takes the country relationship between the contributor
 450 and the integrator into consideration; *same_user* is the
 451 precondition for eight factors (see Table 5)). We also
 452 excluded factors with variance inflation factor (VIF) val-
 453 ues ≥ 5 , as such values could inflate variance, measured
 454 using the *vif* function of the *car* package in R [39]. In
 455 this way, we removed *num_code_comments* that were
 456 otherwise moderately correlated with *num_comments*
 457 ($\rho = 0.63$).

458 3) *Ease of interpretation.* Models perform better when fea-
 459 tures are approximately normal and in a comparable
 460 scale.⁹ We stabilized the variance in features by adding
 461 a value “1” and log-transforming the continuous vari-
 462 ables. Then, we transformed the features into a com-
 463 parable scale with a mean value of “0” and a standard
 464 deviation of “1”.

TABLE 5: Factors with dependency

postconditional factor	preconditional factor
<i>perc_pos_emotion</i>	
<i>perc_neg_emotion</i>	<i>has_comments</i>
<i>first_response_time</i>	
<i>perc_contrib_pos_emo</i>	<i>contrib_comment</i>
<i>perc_contrib_neg_emo</i>	
<i>perc_inte_neg_emo</i>	<i>inte_comment</i>
<i>perc_inte_pos_emo</i>	
<i>ci_latency</i>	<i>ci_exists</i>
<i>ci_failed_perc</i>	
<i>same_country</i>	
<i>same_affiliation</i>	
<i>contrib_follow_integrator</i>	
<i>open_diff</i>	<i>same_user</i>
<i>cons_diff</i>	
<i>extra_diff</i>	
<i>agree_diff</i>	
<i>neur_diff</i>	

2.4.1 *Factors influencing pull request decisions*

465 To explain pull request decisions, we intended to build a
 466 model with all the known factors. However, in practice, this
 467 is not possible. We noticed that the postconditional factors
 468 (see Table 5) did not make sense unless a precondition was
 469 met. For example, the factor *ci_latency* was meaningful only
 470 when the factor *ci_exists* was true. Here, *ci_exists* presents a
 471 precondition contingent on which factors, such as *ci_latency*,
 472 are meaningful, which are also referred to as postconditional
 473 factors. Table 5 presents a complete list of the dependent
 474 factors in our dataset. The remaining factors have no such
 475 dependency on other factors.

476 To understand how the identified factors influence pull
 477 request decisions, we built two types of models.

- 478 1) We built a *basic model* that comprised all the factors with
 479 no dependencies on each other and preconditional fac-
 480 tors. This model offered an overview without entering
 481 the details offered by the postconditional factors.
- 482 2) Next, we built models for the *special cases* relating to
 483 preconditions: developer, pull request, and tools as
 484 identified in Table 5.

- 485 • *developer*: when the contributor and the integrator are
 486 not the same users (*same_user=0*)
- 487 • *pull request*: when a pull request has comments
 488 (*has_comment=1*)
- 489 • *tool*: when a pull request uses the CI tool (*ci_exists=1*).
 490 Each of these special case models are built on a subset
 491 of the data used in the basic model that meets the
 492 precondition.

2.4.2 *Influence of context*

494 To explore the relevance of context in explaining pull re-
 495 quest decisions, we studied five scenarios relating to the
 496 developer, pull request, project, tools, and time. Figure 2
 497 presents a pictorial depiction of the five scenarios in rela-
 498 tion to the pull request decision and metrics. To study the
 499 influence of context, we trained the same model on different
 500 observations representing specific contexts.

- 501 • *developer characteristic*: We chose the factor *same_user*
 502 indicating whether a pull request is submitted and
 503 integrated by the same user. It is the most important de-
 504 veloper characteristic influencing pull request decisions

9. <https://medium.com/@sjacks/feature-transformation-21282d1a3215>

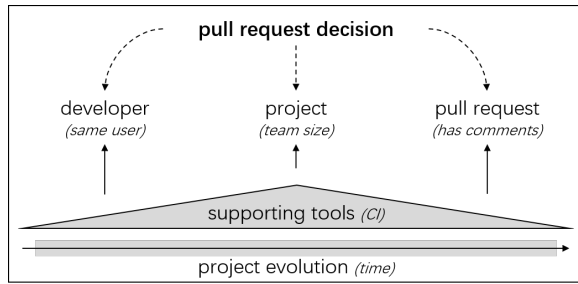


Fig. 2: Contexts in pull request decisions

(see the basic model in Table 6) and a precondition for a range of factors. We think that pull requests integrated by oneself behave differently than those integrated by others.

- *pull request characteristic*: We chose the factor *has_comments* as an indicator of a pull request characteristic influencing the decision [15]. It is one of the top five factors influencing decisions (see the basic model in Table 6) and a precondition for several factors, including *perc_pos_emotion* and *first_response_time* (see Table 5). This factor explores decisions for pull requests both with and without comments.
- *project characteristic*: We selected the factor *team_size* as an indicator of project characteristics such as project popularity and maturity. We assumed that teams of different sizes represented different contexts (as also seen in other studies [40], [41]). We studied three team sizes: small ($team_size \leq 4$), medium ($4 < team_size \leq 10$), and large ($team_size > 10$). Here we split the pull request according to the tertile of factor *team_size*.¹⁰
- *supporting tools*: We selected the factor *ci_exists* for its reported influence on pull request decisions [10] and relevance in our special case model (refer to Table 6). In addition, a previous study has shown that the usage of CI tools changes during the development of projects [43]. Therefore, we assumed that factors influence pull request decisions differently depending on whether they are pull requests using CI tools or those not using CI tools.
- *project evolution*: We studied temporal evolution to see if the process changed over time. We studied decision-making in three time periods: before June 1st, 2016, between June 1st, 2016, and June 1st, 2018, and between June 1st, 2018, and June 1st, 2019 (aka after June 1st, 2018). A pull request belonged to a time period when it was integrated. For this scenario, we included only projects (and their pull requests) active in all three time periods.¹¹

2.4.3 Interpretation of statistical models

The resulting mixed effects logistic regression models explain the influence of factors in models and their relative

10. A sensitivity analysis with threshold values (small size ranging from 2-6, large size ranging from 8-12) yielded similar results. See the technical report [42] for the detailed results.

11. A sensitivity analysis with threshold values (first period ranging from December 1st, 2015 to December 1st, 2016, third period ranging from December 1st, 2017 to December 1st, 2018) yielded similar results. See the technical report [42] for the detailed results.

relevance. Section 3 presents the findings from these mixed effects logistic regression models. Each model has two parts: an intercept and influence of a factor, expressed as follows:

$$\text{odds ratio}^{p\text{-value}}[\text{percentage variance}] \quad (1)$$

The *odds ratio* expresses the association between a factor and a pull request decision as “the increase or decrease in the odds of acceptance for a ‘unit’ increase of a factor” [15]. In this work, a “unit” of each factor was one standard deviation from the standardization of the log-transformed factors. The term *p value* indicates the statistical significance of a factor, which was indicated by asterisks: *** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$ [10], [12]. It represents the probability of the evidence against the null hypothesis, *i.e.*, “there is no association between each factor and pull request decisions.” Finally, the *percentage of explained variance* was used as a proxy for the relative importance of a factor. The variance explained by each factor is derived from ANOVA Type-II analysis [44]. When it is relative to the total amount of variance (the percentage of explained variance), the result can serve as a proxy for effect size, which means how much effect one factor has in explaining pull request decisions. This metric is similar to the percentage of total variance explained by least squares regression [39] and has been used in prior studies [45].

We reported the *goodness of fit* of each model using the area under the receiver operating characteristic curve (*AUC*) value (for training data), where an *AUC* value greater than 0.5 indicated the effectiveness of the model [12]. We also reported the predictive performance of related models using the weighted precision, weighted recall, and weighted f-score [46].

In practice, we split the pull requests in close time and used the first 90% of pull requests for training and the remaining 10% for testing. We measured the predictive performance of the basic model only to present the prediction effect of pull request decisions by integrating as many factors as possible and to explain factor performance in other situations, without reporting their prediction performance. The above metrics collectively indicated the predictive performance of both the baseline and logistic regression models for our highly imbalanced dataset [46].

3 RESULTS

This section presents how factors influence pull request decisions (answering **RQ1**) via a basic model, which comprises all the factors likely to influence pull request decisions, excluding those that cannot make it to the basic model. Next, we describe how the factors influencing pull request decisions change with a change in context (answering **RQ2**). We present five scenarios representing developer, pull request, project, tool, and time characteristics.

3.1 RQ1: How do factors influence pull request decisions?

3.1.1 Basic model

Our basic model in Table 6 (column 3) shows 46 factors known to influence pull request decisions arranged in non-increasing order of relative relevance. In comparison to a

random classifier (with weighted precision: 0.81, weighted recall: 0.79, weighted f-score: 0.80, and AUC_test: 0.50), our basic model performed better (with weighted precision: 0.89, weighted recall: 0.90, weighted f-score: 0.89, and AUC_test: 0.82), suggesting an improvement in our model in terms of decision making.

The five most important factors influencing pull request decisions are *same_user*, *lifetime_minutes*, *prior_review_num*, *has_comments* and *core_member*. Table 6 (column 3) shows that these top five factors (shown in dark gray) explain approximately 83% of the variance. This number reaches approximately 95% when considering the influence of the top 10 factors. The remaining 36 factors collectively explain 5% of the explained variance.

The most important factor influencing pull request decisions was *same_user* (with 31% variance). Moreover, *same_user* decreased the odds of acceptance of a pull request by 48% per unit when a pull request was integrated by the contributor. One possible explanation for this observation relates to the process of pull-based development. Due to the standardized process of such development, contributions should be reviewed and merged by others during the process. However, since all contributors could close their own pull requests, it was possible for them to find problems in their pull requests from others' comments or CI build results and close their own pull requests.

Through Table 8, we can see that many project related factors, including *open_pr_num*, *project_age*, *sloc*, were found to influence pull request decisions in related works significantly. However, through integrating various factors, we find that the project related factors did not contribute greatly to pull request decisions, as these factors explained only approximately 1% of the variance. However, the developer- and pull-request-related factors are more important, explaining 52% and 46% of the variance, respectively. See the dynamic treemap to compare the relative importance of factors in different categories visually.¹²

3.1.2 Special cases

Table 6 shows the results of the three special cases in the last three columns. Factors ranking the top 5 in each model (T_{1-5}) are shown in deep gray, and factors ranking in the top 6-10 in each model (T_{6-10}) are shown in light gray.

When the contributor and integrator were different users (*same_user*=0) (see column 4 in Table 6), we found that three additional factors had a small effect on pull request decisions. The only factor that made it into the top 10 factors was personality difference, namely, differences in agreeableness (*agree_diff*). The two other factors were differences in openness to experience (*open_diff*), also indicating differences in personality, and the same affiliation of the contributor and integrator (*same_affiliation*).

When there existed at least one comment (*has_comments*=1) (see column 5 in Table 6), positive emotion became relatively important, with a sizable effect (> 3% variance). This change can be attributed to the phenomenon that positive reactions during the code review process can lead to contributors' active

participation and increase the likelihood of pull request acceptance. However, negative emotion is not important in pull request decisions. A possible explanation for this is that different developers tend to act differently toward negative emotion. Therefore, negative emotion during discussion faces difficulty in effectively making the final decision. To verify our observation, we built models for pull requests that had at least one comment from a contributor (*contrib_comment*=1) or at least one comment from an integrator (*inte_comment*=1) [42]. We found that both *perc_contrib_pos_emo* and *perc_inte_pos_emo* explained more than 3% of the variance, which was much higher than that of negative emotion.

When pull requests used CI tools (*ci_exists*=1) (see column 6 in Table 6), factor *ci_failed_perc* stood out, explaining 18% of the variance, which implies that the build status of CI tools is important for review decisions, especially the percentage of build failures.

Pull request decisions is mostly explained by a few factors (5 to 10 factors) such that developer and pull request characteristics are more important than project characteristics.

The relation between contributor and integrator (*same_user*) is the most important factor influencing pull request decisions.

In special cases, when a pull request has comments, comment's positive emotion is linked to pull request acceptance. Likewise, when pull requests use CI tools, the percentage of failed CI builds become important for pull request decisions.

3.2 RQ2: How do the factors influencing pull request decisions change with a change in context?

3.2.1 Developer characteristic

Table 7 shows that in comparison to the pull requests submitted and integrated by the same user, when the contributor and integrator are not the same person, the variance explained by the experience of the integrator (*prior_review_num*) decreases from 31% (row 1, column 2 - same user: yes) to 0% (row 1, column 3 - same user: no). This finding implies that the integrator's experience plays a limited role when making decisions regarding others' contributions. However, this factor becomes very important for an integrator's own contributions. One way to explain this observation can be that external contributors, without review experience, generally do not have the right to merge the code. Experienced integrators, in contrast, are familiar with the management process, know when to merge a pull request, and have the ability to merge a pull request. In this way, differences in permission linked to integrators' experience can influence pull request decisions.

For the lifetime of pull requests (*lifetime_minutes*), the percentage of explained variance increased from 19% (row 2, column 2 - same user: yes) to 44% (row 2, column 3 - same user: no). A possible explanation for this observation is that when there is no response from a contributor for a long time, a pull request is more likely to be closed by the reviewer. However, when the pull request is reviewed by

12. https://github.com/zhangxunhui/TSE_pull-based-development/blob/main/treemap-basic-model.html

TABLE 6: Results of special cases. - means the factor is not included in the model. Color: deep gray represents factors with explained variance rank in the Top 5 and light gray represents factors rank in the Top 6-10.

Factor Index		Dependent variable: merged_or_not=1			
		basic model	same_user=0	has_comments=1	ci_exists=1
	(Intercept)	21.1***	32.5***	15.4***	26.1***
(1)	same_user	0.52***[31.17]	-	0.60***[21.53]	0.50***[21.27]
(2)	lifetime_minutes	0.61***[21.10]	0.52***[43.08]	0.53***[30.40]	0.50***[26.08]
(3)	prior_review_num	1.53***[13.53]	1.06** [0.20]	1.50***[13.94]	1.50***[8.01]
(4)	has_comments	0.63***[11.97]	0.52***[25.39]	-	0.64***[6.70]
(5)	core_member	1.29***[5.29]	1.02 [0.05]	1.30***[5.86]	1.32***[3.58]
(6)	num_commits	1.30***[4.49]	1.35***[6.67]	1.58***[13.65]	1.56***[7.43]
(7)	other_comment	1.21***[3.76]	1.27***[6.47]	1.12***[1.58]	1.24***[2.88]
(8)	ci_exists	1.16***[1.47]	1.25***[5.13]	1.11***[0.93]	-
(9)	hash_tag	1.12***[1.36]	1.06***[0.51]	1.10***[1.15]	1.13***[1.04]
(10)	files_added	0.91***[0.74]	0.96** [0.18]	0.91***[0.61]	0.90***[0.48]
(11)	prev_pullreqs	1.15***[0.73]	1.10***[0.51]	1.16***[0.99]	1.15***[0.43]
(12)	commits_on_files_touched	1.09***[0.59]	1.13***[1.51]	1.05***[0.22]	1.03***[0.04]
(13)	open_pr_num	0.82***[0.47]	1.16***[0.26]	0.94***[0.05]	0.87***[0.15]
(14)	account_creation_days	1.06***[0.41]	1.16***[2.52]	1.06***[0.41]	1.09***[0.53]
(15)	first_pr	0.95***[0.36]	0.99 [0.01]	0.96***[0.27]	0.96***[0.16]
(16)	test_churn	1.07***[0.27]	1.10***[0.59]	1.11***[0.59]	1.12***[0.42]
(17)	files_changed	0.92***[0.26]	0.91***[0.42]	0.94***[0.14]	0.97***[0.02]
(18)	project_age	1.11***[0.26]	1.06 [0.08]	1.08***[0.19]	1.21***[0.53]
(19)	reopen_or_not	0.97***[0.25]	0.99 [0.05]	0.98***[0.12]	0.98***[0.08]
(20)	contrib_open	1.06***[0.24]	1.05***[0.27]	1.05***[0.20]	1.07***[0.20]
(21)	stars	0.86***[0.22]	0.88***[0.22]	0.79***[0.69]	0.89***[0.10]
(22)	inte_open	1.06***[0.21]	1.10***[0.46]	1.10***[0.64]	0.98***[0.01]
(23)	description_length	1.04***[0.17]	1.02 [0.04]	1.01 [0.00]	1.01***[0.01]
(24)	pushed_delta	1.04***[0.15]	1.06***[0.39]	1.04***[0.22]	1.04***[0.12]
(25)	followers	1.04***[0.12]	0.92***[0.52]	1.02***[0.04]	1.03***[0.03]
(26)	contrib_cons	1.03***[0.07]	1.04** [0.16]	1.05***[0.17]	1.03***[0.03]
(27)	team_size	1.06***[0.06]	1.02 [0.00]	1.06***[0.07]	1.07***[0.06]
(28)	contrib_gender	0.98***[0.05]	0.93***[0.54]	0.97***[0.10]	0.98***[0.03]
(29)	files_deleted	0.98***[0.03]	0.99 [0.02]	0.96***[0.18]	0.96***[0.10]
(30)	pr_succ_rate	0.98***[0.03]	1.09***[0.73]	0.98***[0.05]	0.96***[0.06]
(31)	contrib_agree	0.98***[0.02]	0.99 [0.00]	0.97***[0.05]	0.98***[0.02]
(32)	contrib_extra	0.99***[0.02]	0.94***[0.29]	0.97***[0.07]	0.97***[0.04]
(33)	contrib_neur	1.02***[0.02]	1.07***[0.41]	1.01** [0.01]	1.00 [0.00]
(34)	inte_neur	1.02***[0.02]	1.00 [0.00]	1.04***[0.08]	0.99 [0.00]
(35)	num_comments	1.02***[0.02]	1.00 [0.00]	0.91***[0.88]	0.97***[0.04]
(36)	comment_conflict	1.01***[0.01]	1.00 [0.00]	1.02***[0.05]	1.01***[0.01]
(37)	friday_effect	1.01***[0.01]	1.01 [0.02]	1.02***[0.06]	1.02***[0.02]
(38)	inte_agree	1.02***[0.01]	0.89***[0.54]	0.98***[0.02]	1.02* [0.01]
(39)	inte_extra	1.01***[0.01]	1.02 [0.01]	1.01* [0.01]	1.06***[0.10]
(40)	open_issue_num	1.03***[0.01]	1.08 [0.07]	1.02 [0.00]	1.03 [0.00]
(41)	sloc	1.02***[0.01]	0.97 [0.04]	1.04***[0.04]	0.93***[0.06]
(42)	test_inclusion	1.02***[0.01]	1.00 [0.00]	1.01* [0.01]	1.02***[0.01]
(43)	inte_cons	1.01 [0.00]	1.04 [0.07]	1.00 [0.00]	0.99 [0.00]
(44)	integrator_availability	1.00 [0.00]	1.04** [0.17]	1.01** [0.01]	1.01 [0.00]
(45)	src_churn	1.00 [0.00]	1.00 [0.00]	1.05***[0.17]	1.07***[0.15]
(46)	test_lines_per_kloc	1.01 [0.00]	0.91***[0.32]	1.02***[0.02]	1.02* [0.00]
(47)	agree_diff	-	0.93***[0.55]	-	-
(48)	cons_diff	-	0.98 [0.03]	-	-
(49)	contrib_follow_integrator	-	1.01 [0.01]	-	-
(50)	extra_diff	-	0.99 [0.01]	-	-
(51)	neur_diff	-	0.98 [0.05]	-	-
(52)	open_diff	-	0.97* [0.09]	-	-
(53)	same_affiliation	-	1.06***[0.30]	-	-
(54)	same_country	-	1.02 [0.03]	-	-
(55)	perc_pos_emotion	-	-	1.18***[3.14]	-
(56)	perc_neg_emotion	-	-	0.96***[0.37]	-
(57)	first_response_time	-	-	1.01***[0.02]	-
(58)	ci_failed_perc	-	-	-	0.65***[18.28]
(59)	ci_latency	-	-	-	1.11***[0.67]
	Observations	1,765,730	91,874	839,505	954,386
	AUC_train	0.848	0.891	0.850	0.865

TABLE 7: Partial results in different contexts. Whole results are shown in Appendix A.

Gray color marks the factors that have more than 5% difference of explained variance in different contexts. Value before bracket means the odds ratio, value in bracket means the percentage of explained variance, - means the factor is not included in the model.

(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)
Dependent variable: <i>merged_or_not=1</i>												
	same user or not		has comments or not		ci exists or not		different team sizes			different periods		
	yes	no	yes	no	yes	no	small	mid	large	before 2016.6	2016.6-2018.6	after 2018.6
(Intercept)	10.4	34.4	13.1	42.4	20.4	13.3	24.9	20.7	15.9	6.9	16.6	7.1
(1) <i>prior_review_num</i>	2.86[31]	0.98[0]	1.51[14]	1.91[22]	1.53[14]	1.53[12]	1.59[11]	1.41[9]	1.57[19]	1.30[6]	1.63[14]	1.72[17]
(2) <i>lifetime_minutes</i>	0.66[19]	0.52[44]	0.61[30]	0.70[13]	0.60[22]	0.61[21]	0.54[24]	0.61[20]	0.67[17]	0.65[20]	0.57[21]	0.62[13]
(3) <i>core_member</i>	1.26[9]	1.13[1]	1.29[6]	1.33[6]	1.30[5]	1.26[5]	1.42[6]	1.28[5]	1.19[3]	1.27[6]	1.34[5]	1.29[3]
(4) <i>num_commits</i>	1.23[4]	1.46[10]	1.49[11]	0.98[0]	1.32[5]	1.25[4]	1.36[5]	1.31[5]	1.26[4]	1.18[2]	1.32[4]	1.36[5]
(5) <i>commits_on_files_touched</i>	1.06[0]	1.11[1]	1.05[0]	1.18[2]	1.06[0]	1.13[1]	1.10[0]	1.12[1]	1.05[0]	1.30[7]	0.99[0]	0.99[0]
(6) <i>has_comments</i>	0.68[10]	0.50[27]	-	-	0.65[10]	0.52[25]	0.57[13]	0.64[10]	0.66[12]	0.63[15]	0.62[10]	0.55[14]
(7) <i>same_user</i>	-	-	0.56[29]	0.42[42]	0.51[33]	0.59[23]	0.49[24]	0.49[36]	0.55[33]	0.57[31]	0.46[33]	0.46[29]
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
Observations	950,985	1,010,937	1,152,714	809,208	1,611,277	350,645	601,460	703,396	701,900	512,707	585,401	274,121
AUC_train	0.862	0.874	0.837	0.872	0.843	0.884	0.877	0.843	0.837	0.850	0.867	0.879

705 the contributor himself/herself, he/she knows exactly what
 706 is happening, and the related decision making is thus not in-
 707 fluenced as much by the lifetime of a pull request. Likewise,
 708 for the number of commits (*num_commits*), the percentage
 709 of explained variance increased from 4% (row 4, column 2
 710 - same user: yes) to 10% (row 4, column 3 - same user: no).
 711 It is likely that during the interaction, the integrator will
 712 ask the contributor to modify the contribution, increase the
 713 number of commits, and then make decisions according to
 714 these changes.

715 When comments were present (*has_comments*), the ex-
 716 plained variance increased when a pull request was inte-
 717 grated by another person in comparison to oneself from 10%
 718 (row 6, column 2 - same user: yes) to 27% (row 6, column
 719 3 - same user: no). This result can be explained by the fact
 720 that when integrating pull requests submitted by others, it
 721 is common for the integrator to understand the contribution
 722 by communicating with the contributor.

723 For whether the contributor is a core developer
 724 (*core_member*), we find a notable difference in the influence
 725 of this factor on the pull request decisions in the set of self-
 726 integrated pull requests (row 3, column 2 - same user: yes)
 727 and the other-integrated pull requests (row 3, column 3 -
 728 same user: no). Although this factor is positively correlated
 729 with the pull request decisions in both cases (odds ratio>1),
 730 *i.e.*, pull requests submitted by core developers are more
 731 likely to be accepted than those submitted by external
 732 contributors; the explained variance reduces from 9% to
 733 1%. This indicates that whether the contributor is a core
 734 developer becomes less important than other factors for pull
 735 requests integrated by others.

Whether the contributor and integrator is the same person or not influences pull request decisions the most.

If the contributor and integrator is the same, pull request decisions depend on the contributor's relationship to the target project (*prior_review_num* and *core_member*).

When the contributor and integrator are different, pull request decisions depend on the interaction between contributor and integrator (*has_comments*, *lifetime_minutes*) and the intermediate results during the process (*num_commits*).

3.2.2 Pull request characteristic

737 When a pull request did not have comments, the percentage
 738 of explained variance of *same_user* increased from 29% (row
 739 7, column 4 - has comments: yes) to 42% (row 7, column 5 -
 740 has comments: no). This situation illustrates that the factor
 741 *same_user* is more associated with pull request decisions for
 742 those without comments. To investigate the reason, we cal-
 743 culated the merging rate of pull requests in four situations
 744 (see Table 8).
 745

TABLE 8: Pull request merge rate for *has_comments* and *same_user* cross situations

	<i>has_comments=true</i>	<i>has_comments=false</i>
<i>same_user=true</i>	74.5%	88.3%
<i>same_user=false</i>	82.1%	93.3%

746 From the table, we can find that for pull requests without
 747 comment, the merge rate increases for both cases of fac-
 748 tor *same_user*. However, we find that the merge rate even
 749 reaches 93% when *same_user=false*. Such high probability
 750 may be why this factor plays a decisive role in explaining
 751 pull request decisions when there is no comment.

752 Regarding integrator experience (*prior_review_num*), the
 753 explained variance increased from 14% (row 1, column 4 -
 754 has comments: yes) to 22% (row 1, column 5 - has comments:
 755 no). It is likely that when there are no comments, there
 756 are cases in which developers close or merge their own
 757 pull requests. In comparison to core members, external
 758 developers do not have the right to merge. This restricted
 759 permission linked to the integrator's review experience can
 760 potentially influence the pull request decision.

761 For the lifetime of a pull request (*lifetime_minutes*)
 762 and the number of commits included in a pull request
 763 (*num_commits*), when there exist comments, the integrator
 764 tends to make the decision based on the contributor's re-
 765 sponse speed and how he/she modifies the contribution
 766 according to the integrator's suggestions. This can be a
 767 reason why there exists a higher percentage of variance in
 768 situations where comments exist.

When there is no communication between the contributor and reviewers, factors indicating the affiliation of a contributor to the project - whether the contributor and the integrator are the same (*same_user*) and review experience (*prior_review_num*), are important in influencing pull request decisions. When there is communication between the contributor and reviewers, factors representing the activeness of the interaction (*lifetime_minutes*, *num_commits*) have a bigger influence on pull request decisions.

770 3.2.3 Project characteristic

771 As team size increased, the variance explained by the exper-
 772 ience of the integrator (*prior_review_num*) initially decreased
 773 from 11% (row 1, column 8 - team size: small) to 9% (row
 774 1, column 9 - team size: mid) and then increased from 9%
 775 (row 1, column 9 - team size: mid) to 19% (row 1, column 10
 776 - team size: large).

777 When considering whether pull requests were submitted
 778 and integrated by the same user (*same_user*), the change
 779 trend was the opposite, increasing from 24% (row 7, column
 780 8 - team size: small) to 36% (row 7, column 9 - team size:
 781 mid) and then decreasing from 36% (row 7, column 9 - team
 782 size: mid) to 33% (row 7, column 10 - team size: large).

783 These two types of change indicate that for pull re-
 784 quests targeting teams of different sizes, the importance of
 785 *prior_review_num* and *same_user* changed nonlinearly. How-
 786 ever, we have no explanation for this observation.

As team size increases, integrator's experience (*prior_review_num*) and whether submitter and integrator are the same (*same_user*) have a V-shaped and inverted V-shaped relations to pull request decisions respectively.

788 3.2.4 Supporting tools

789 When not using CI tools, the percentage of variance ex-
 790 plained by comments (*has_comments*) was 25% (row 6, col-
 791 umn 7 - ci exists: no), which was higher than that of pull
 792 requests using CI tools (10%) (row 6, column 6 - ci exists:

793 yes). This result can be explained by the fact that when
 794 there are no CI tools, contributors can obtain feedback only
 795 from reviewers. Therefore, whether comments exist matters
 796 greatly in pull request decisions. When using CI tools,
 797 contributors can first obtain responses from CI outcomes,
 798 which can help with making decisions.

799 For factor *same_user*, its explained variance decreases
 800 from 33% (row 7, column 6 - ci exists: yes) to 23% (row
 801 7, column 7 - ci exists: no). According to the previous
 802 study [11], teams using CI tools are more effective at merg-
 803 ing pull requests submitted by core members. Therefore, we
 804 think that the existence of CI tools leads contributors to be
 805 more able to make judgments about their own contributions
 806 through the build outcome.^{13,14}

The use of CI tools leads to significant changes in the influence of two factors on pull request decisions, i.e., whether the pull request contains comments and whether the contributor and the reviewer are the same people. When using CI tools, the availability of CI build results makes the comments less important in explaining pull request decisions, while the influence of contributor and integrator's relationship becomes stronger.

807 3.2.5 Project evolution

808 Before June 2016, the experience of the integrator
 809 (*prior_review_num*) explained just 6% (row 1, column 11 -
 810 period: before 2016.6) of the variance, which increased to
 811 17% after June 2018 (row 1, column 13 - period: after 2018.6).
 812 We calculated the experience of integrators corresponding
 813 to pull requests at different periods of project development,
 814 as shown in Figure 3. We find that the gap between inte-
 815 grators' experience for merged and unmerged pull requests
 816 gradually increases as projects become mature. This is why
 817 the variance explained by factor *prior_review_num* gradually
 818 increases. This indicates that the integrator's experience
 819 gradually becomes an important indicator of pull request
 820 decisions as the project evolves.

821 For the area hotness of contributions (*com-
 822 mits_on_files_touched*), before June 2016, it had a moderate
 823 effect on the decision-making of pull requests, which
 824 explained 7% of the variance (row 5, column 11 - period:
 825 before 2016.6), and increased the odds of acceptance by 30%
 826 per unit. However, as projects became mature, the variance
 827 explained decreased to 0% (row 5, column 13 - period:
 828 after 2018.6). For the three periods, we also calculated the
 829 mean value of *commits_on_files_touched* (before 2016.6: 40,
 830 2016.6-2018.6: 33, and after 2018.6: 28), which shows that
 831 the contributions in the early stage of the project were more
 832 concentrated. In other words, as projects become larger and
 833 more mature, contributions are more widely distributed,
 834 and the area hotness of pull requests can hardly contribute
 835 to the merging of pull requests for mature projects.

836 For the lifetime of pull requests (*lifetime_minutes*), the
 837 explained variance decreased from 20% (row 2, column 11
 838

13. <https://github.com/react-boilerplate/react-boilerplate/pull/2256>

14. <https://github.com/mggg/GerryChain/pull/290>

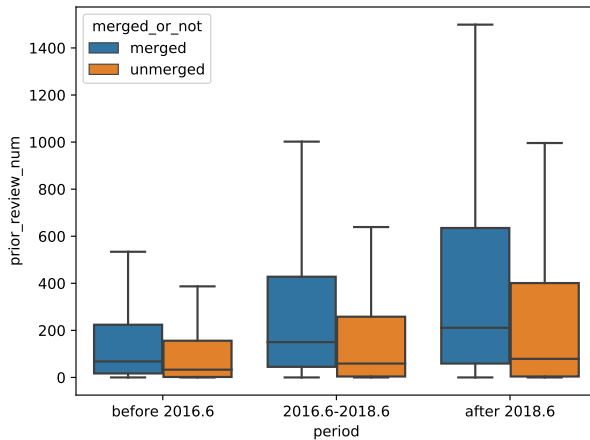


Fig. 3: The comparison between integrators' experience

839 - period: before 2016.6) to 13% (row 2, column 13 - period:
 840 after 2018.6). Although this factor negatively influenced the
 841 pull request merging in all three time periods, the effect size
 842 decreased. We calculated the changes in the pull request
 843 lifetime median value as projects evolve. It is found that
 844 the overall processing time of pull requests increases significantly
 845 (before 2016.6: 802min, 2016.6-2018.6: 1,188min, and
 846 after 2018.6: 1,316min). There are many possible reasons for
 847 this situation. *E.g.*, at the beginning of a project, the develop-
 848 ment team is small, and the pull requests that have been
 849 left unprocessed for a long time are likely to be rejected. As
 850 the project develops, more pull requests are left unprocessed
 851 (before 2016.6: 58, 2016.6-2018.6: 112, and after 2018.6: 174).
 852 The reviewers have their processing order, so the overall
 853 processing time of pull requests grows, but the impact on
 854 the decision becomes smaller. Also, we think as projects
 855 become mature, the use of various supporting mechanisms
 856 in the review process becomes stabilized, *e.g.*, the use of CI
 857 tools [10], the request of reviews [47], etc. These mechanisms
 858 lead to the increase of pull request lifetime. However, the
 859 standardized processes reduce the impact of processing time
 860 on the final result. There may not be a single reason for the
 861 change in results. Still, the result reveals that pull request
 862 processing time on decision-making decreases as the project
 863 develops.

As a project evolves, the integrator's experience (*prior_review_num*) becomes more and more important for pull request decisions, while the area hotness of contribution (*commits_on_files_touched*) no longer influences the decision making. Compared to the early stages of project evolution, the influence of pull request lifetime (*lifetime_minutes*) on pull request decisions decreases.

865 4 CASE STUDY

866 Since companies' contribution is relatively high in the open
 867 source world [48], the strategy, decision making, and partic-
 868 ipation patterns of different companies in open source

vary greatly [49]. The participation of companies in open
 source projects also impacts the inflow and retention of
 external contributors [50]. Therefore, we also consider it
 interesting to analyze the impact of affiliation-related factors
 on pull request decisions. Therefore, we added the analysis
 of affiliation-related factors.

We first merge developer accounts and ignore those with
 more than one affiliation (this may be due to developers' af-
 filiation). When considering the merge rate (Figure 4,5,6), we
 only consider the pull requests submitted and integrated by
 different users, as factor *same_user* significantly influences
 pull request decisions and acts as the precondition of factor
same_affiliation.

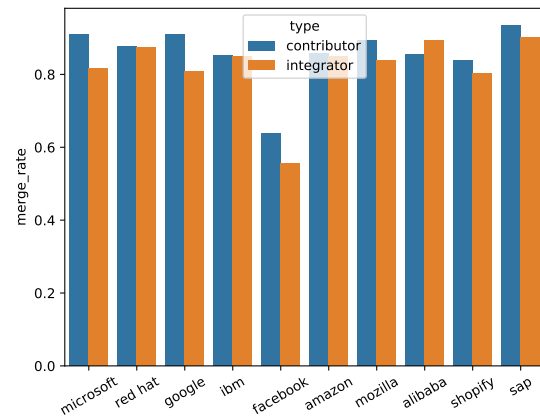


Fig. 4: Merge rate of top 10 affiliation when acting as contributor and integrator respectively

Different affiliations have different contribution intensi-
 ties regarding the number of submitted and integrated pull
 requests [49]. Figure 4 shows that the merge rate for dif-
 ferent affiliations varies a lot. For Facebook, its related pull
 requests' merge rate is much lower than other affiliations.
 This may be related to differences in policies or the way
 contributions are handled by different companies.

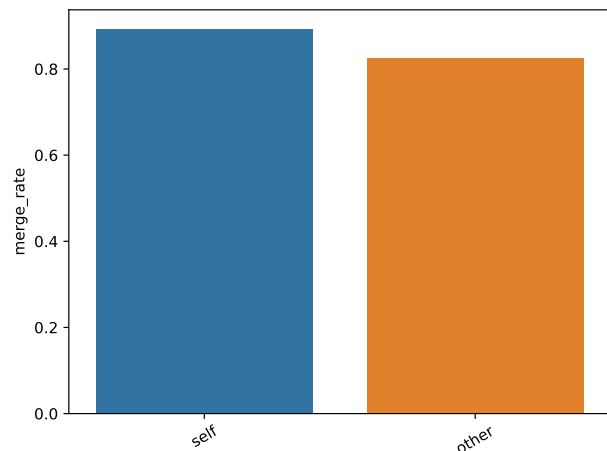


Fig. 5: Overall merge rate for affiliations integrating their own contributions (self) or contributions from other affiliations (other)

889 Second, we consider the effect of whether the pull
 890 request submitter and the integrator are from the same
 891 affiliation on pull request decisions. In the overall case, the
 892 merging probability is higher for pull requests submitted by
 893 their colleagues than those by developers from other affili-
 894 ations (see Figure 5), which is in line with our perception.
 895 However, from the result of **RQ2** (Table 5 *same_user=0*), we
 896 found that when considering together with other factors, the
 897 factor *same_affiliation*, although significantly associated with
 898 pull request decisions, is less effective (explaining only 0.3%
 899 variance).

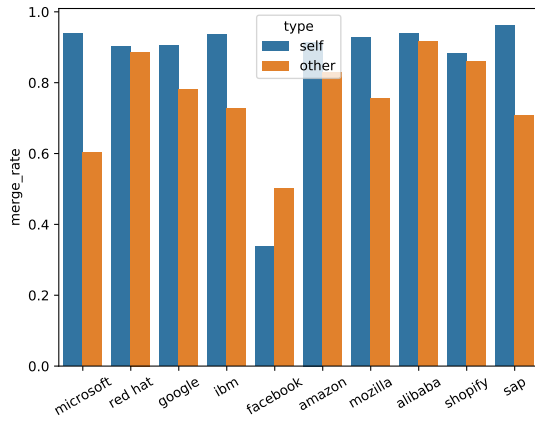


Fig. 6: Merge rate of different affiliations when integrating their own contributions (self) or contributions from other affiliations (other)

900 Our statistical analysis of each company reveals differ-
 901 ences in the way companies treat their own contributions
 902 and external contributions (see Figure 6). For Facebook, the
 903 probability of merging external contributions is even higher
 904 than that of merging internal contributions. We think that
 905 the policy and openness of different companies lead to the
 906 different treatment of external contributions.

5 DISCUSSION

5.1 Pull request decisions explained

909 Our study shows that there is no one answer to our research
 910 questions. Instead, there are generic answers and specific
 911 answers for the context represented, given the dependen-
 912 cies among factors. Generally, whether a pull request is
 913 submitted and integrated by the same person, its lifetime,
 914 experience of the integrator, presence of comments, and
 915 coreness of the contributor play decisive roles in pull re-
 916 quest decisions. When comments in pull requests exist, the
 917 positive emotion for communication influences pull request
 918 decisions. When pull requests use CI tools, the percentage
 919 of build failure influences the decision.

920 Interestingly, the influence of the factors changes with a
 921 change in context:

922 *Developer characteristic (same user or not):* Compared to
 923 pull requests integrated by different persons, when pull
 924 requests are submitted and integrated by the same person,
 925 the importance of the integrator's experience and the con-
 926 tributor's coreness increase for pull request decisions, while

the importance of the pull request lifetime and the included
 number of commits decreases (Section 3.2.1).

Pull request characteristic (has comments or not): When pull
 requests have comments, the lifetime and the number of
 commits included are more important compared to pull
 requests without any comment. In contrast, the importance
 of the integrator's experience and whether the contributor
 and integrator are the same person are less important when
 comments exist (Section 3.2.2).

Project characteristic (different team sizes): The importance
 of the integrator's experience and whether the contributor
 and the integrator are the same person for pull request
 decisions changes nonlinearly for teams of different sizes
 (Section 3.2.3).

Tool (CI exists or not): The use of CI tools decreases the
 importance of comment existence, but the importance of
 whether the contributor and the integrator are the same
 person increases for pull request decisions (Section 3.2.4).

Project evolution (different periods): The importance of the
 integrator's experience in pull request decisions increases
 as projects evolve, while the importance of area hotness and
 the lifetime of the contribution decreases (Section 3.2.5).

5.2 Relations to the literature

5.2.1 Discussion of previous conclusions

949 Referring to the literature (summarized in Table 10), rel-
 950 atively speaking, project-related factors are less discussed
 951 than pull-request- and developer-related factors. To this end,
 952 our study contributes in that not only have few project
 953 characteristics been explored in the literature, but they have
 954 been considered relatively less important (explains 2% of
 955 the variance) than developers (explains 52% variance) and
 956 pull request characteristics (explains 46% variance). Our
 957 study further provides evidence that human factors are as
 958 important or more important than technical factors [51].

960 When comparing the findings of previous studies with
 961 each other and those of our study, we found that in most of
 962 the cases, the results were consistent. Only four factors had
 963 opposite findings regarding the direction of influence, *i.e.*,
 964 *files_changed*, *project_age*, *team_size* and *num_commits*. One
 965 potential explanation that has emerged from our study is
 966 that all these factors are relatively less important for pull
 967 request decisions, which can potentially explain the differ-
 968 ences in the findings. Alternatively, this can simply be due
 969 to the differences in the dataset used. Interestingly, many
 970 factors that are widely studied across related works, *e.g.*,
 971 *core_member* and *src_churn*, indicating that these factors are
 972 likely to influence the decision, are not as important for pull
 973 request decisions.

974 For the factor *num_commits*, which is relatively impor-
 975 tant, ranking in the top 10 across models (Table 6), we
 976 focus on this factor to uncover the reasons for conflict
 977 findings between previous studies. Yu et al. [10] found a
 978 positive effect (the likelihood of pull requests being accepted
 979 increases as the number of commits increases), while other
 980 studies [52], [53], [54] found a negative effect. Our results
 981 are consistent with Yu et al. and argue that the number of
 982 commits cannot simply indicate the contribution size. At
 983 the time of submission, the number of commits indicates
 984 the contribution's size to some extent. However, as the pull
 985

request review process continues, contributors will modify their contributions based on the review feedback and thus complete more commits to facilitate the merging of contributions. Accordingly, we collect the number of commits contained in a pull request at both open time and close time, investigate their effects on pull request merging separately, and find that the number of commits at commit time is negatively correlated with pull request merging. At the same time, it shifts to a positive correlation at close time.¹⁵ Therefore, when a pull request is submitted, the number of commits represents the size of the contribution [52], [53], [54]. However, commits during the review process represent the changes made by the contributor according to the reviewers' comments, thus increasing the likelihood of pull request acceptance [10].

5.2.2 Findings in general context

Considering all pull requests without distinguishing between contexts, the top 5 factors for explaining pull request decisions are: whether the contributor and integrator are the same people (*same_user*), the lifetime of pull requests (*lifetime_minutes*), the experience of the integrator (*prior_review_num*), whether there exists comment (*has_comments*), whether the contributor is the core member (*core_member*).

- 1) *same_user*. The association of this factor reflects the decision propensity of self-integration in the pull-based development model, *i.e.*, a preference for self-rejected rather than self-approved. As you can see from the related work [55], the self-approved patch is defect-prone. To address this situation, future researchers need to consider whether to change the pull-based development model, *e.g.*, for self-approved contributions, generate a warning to other developers in the community.
- 2) *lifetime_minutes*. In related works [52], [56], they only discussed the direction of the association of this factor with pull request decisions. We found that, compared to other factors, the correlation between lifetime and pull request decisions is relatively high. In the future, when exploring the influence of factors on pull request decisions, the lifetime should be considered as an essential control variable.
- 3) *prior_review_num*. This factor is not considered to have a significant association with pull request decisions in related work [57]. However, our result shows that it is significantly important, which ranks the third when considering other factors in an overall perspective. The conflict of conclusion here is not to negate the past research but offers a view applicable at a large scale, as Baysal et al. only did a case study on two projects.
- 4) *has_comments*. Many previous studies focused on the association between the number of comments and pull request decisions [2], [10], [12], [15], [24]. Although there were studies focused on comment existence [53], [58], there is no discussion on its importance and comparing these two factors. Our result finds that the existence of comments is relatively important and can

replace the number of comments in explaining pull request decisions.

- 5) *core_member*. For this factor, compared to previous studies [2], [10], [15], [24], [59], [60], we not only conclude a positive correlation of consistency but also find that the factor has a sizable effect when compared with all the other factors. Unlike the top 4 factors, this factor is present at the time of pull request submission. Therefore, this factor has an irreplaceable effect on predicting pull request decisions at the open time of pull requests.

5.2.3 Findings in different contexts

Under different contexts, we find the relative importance of the influence of postconditional factors. In previous studies, while Iyer et al. [24] found that both positive emotion and negative emotion significantly affect pull request decisions, our results, on the other hand, found that only positive emotion had a sizable effect when considering all factors. It also illustrates that when there exist comments, effectively tapping the hidden positive emotion in comments is important for predicting the final states of pull requests. Also, for pull requests using CI tools [10], the pass of CI builds positively and significantly influences the merging of pull requests. However, our model verifies its relative importance compared to other factors, *i.e.*, the decisions of pull requests are heavily influenced by the outcome of CI builds, which is the third most important factor in explaining pull request decisions.

While having comments leads to a lower probability of merging pull requests, it is needed to differentiate according to the characteristics of the commenter. We found that if there exist comments from others (*other_comment*), *e.g.*, end-users or external developers, the pull request is more likely to be merged (Section 3.1.1). Different from Golzadeh et al. [61], we validated on a much larger dataset and consider different kinds of projects instead of just Cargo ecosystem.

The importance of factors changes and varies significantly as the context changes. And these findings have not been explicitly discussed in previous studies. We find that the number of commits has a sizable effect on the decision-making of pull requests containing comments. For those without comments, the effect is relatively small. This leads to the fact that when studying factors' association with pull request decisions, the impact of the number of commits on pull request decisions should be fully considered when there is no comment. Similarly, for pull requests that do not use CI tools, more significant consideration needs to be given to the weight of the comment. As the project develops, the importance of the factors changes. Among them, the influence of contribution's area hotness (*commits_on_files_touched*) on pull request decisions should be considered for the early stage of the project. And as projects become mature, the experience of integrators becomes important.

5.3 Implications

Our findings have implications for research and practice. Unlike related work, we construct a model from a more comprehensive perspective by collecting measurable factors from all pull request decision-related papers to explain the

15. https://github.com/zhangxunhui/TSE_pull-based-development/blob/main/technical_report.pdf

1100 association and relative importance of factors with pull request
1101 decisions. The discussion of different contexts reveals
1102 the influence of context on the relevance of factors, which
1103 guides future related studies to select appropriate control
1104 variables when empirically analyzing pull request decisions
1105 in global or different contexts. Some findings from the study
1106 also provide theoretical support for future research and the
1107 optimization of pull-based development models. Next, we
1108 will discuss the implications in detail.

1109 **5.3.1 For research**

1110 For future research, this paper can give some guidance. For
1111 example:

1112 *When conducting research on pull request decisions*, re-
1113 searchers can find usable findings from our paper for both
1114 a general overview and specific contexts (see Section 3.1).
1115 *E.g.*, when studying the association of new factors with pull
1116 request decisions, different factors should be considered as
1117 control factors for different situations, and here we give the
1118 recommended list (see Table 9) (the set of factors with more
1119 than 1% of explained variance in various situations). For
1120 other contexts, our dataset and scripts can be used to find
1121 the factors that rank high on the explanation of pull request
1122 decisions in the corresponding contexts as control variables.

1123 Since the impact of a factor on the decision may vary at
1124 different periods of the pull request (*e.g.*, *num_commits* - Sec-
1125 tion 5.2), we think that future research and the construction
1126 of evaluation tools need to consider the impact of changing
1127 factor dynamics.

1128 *When conducting research related to pull-based development*,
1129 researchers can find useful data and conclusions. *E.g.*, when
1130 studying how CI tools influence the code review process,
1131 researchers can easily find that in an overall perspective, the
1132 usage of CI tools increases the likelihood of pull request
1133 acceptance (Section 3.1.1), and the outcome of CI builds
1134 significantly influence the decision making with large effect
1135 (Section 3.1.2). However, there still exist exceptional cases,
1136 *e.g.*, merge without passing CI builds. Thus, subsequent
1137 studies can be conducted based on our data and findings.

1138 **5.3.2 For practice**

1139 The results of our study can provide open source con-
1140 tributors and maintainers with many recommendations for
1141 practices to follow. For example:

1142 *For pull request contributors*, if they want to increase the
1143 chances of having their contributions being accepted, they
1144 should respond to criticism from stakeholders on time, as
1145 the lifetime significantly influences pull request decisions
1146 with a large effect size.

1147 Suppose there are other non-reviewers involved in the
1148 discussion (*other_comment* exists). In that case, the pull re-
1149 quest is more likely to be merged, and contributors are
1150 advised not to give up and modify it according to the project
1151 requirements. As "developers need be more aware of the
1152 human-centric issues of their end-users," [62] one possible
1153 explanation for the influence of *other_comment* is that end-
1154 user feedback can help a lot in improving the quality of
1155 the software.¹⁶ The discussion may be closely related to
1156 the project requirements and development direction, which

1157 directly influences whether the contribution can be merged
1158 or not [63].

1159 *For pull request maintainers*, as the build outcome of CI
1160 tools significantly influences pull request decisions, we re-
1161 commend maintainers install related CI tools to help improve
1162 the merge rate of contributions.

1163 Contributions that remain unprocessed for a long time
1164 are likely not to be merged. On the one hand, maintainers
1165 purposely do not pick pull requests that are either not to
1166 their interest or do not need immediate attention. On the
1167 other hand, reviewers do not respond at the right time [64].
1168 The delay of response may lead to the loss of peripheral
1169 contributors [65] and produce many abandoned contribu-
1170 tions in the long run [66]. We think project managers can
1171 use the mention-bots to reduce the response time [67]. Or
1172 predict and alert on pull request remaining processing time
1173 to speed up the code review [3].

1174 *For both contributors and integrators*, we suggest they
1175 participate in the review process with a positive attitude and
1176 promote the merging of contributions encouragingly. Our
1177 study further solidifies the importance of positive emotion
1178 for pull request decisions by integrating multiple factors. A
1179 positive atmosphere is of great importance for intra-project
1180 communication and efficient collaboration [68].

1181 *For the improvement of the pull-based model*, as we find
1182 that self-integrated pull requests are likely to be rejected,
1183 and a previous study [55] found that self-approved contri-
1184 butions are bug-prone. Therefore, some adjustments can be
1185 made to self-integration. For self-integrated pull requests,
1186 the integrator's experience is a determinant factor for the
1187 decision of pull requests. We wonder if a warning flag
1188 could be added to pull requests integrated by inexperienced
1189 integrators to attract others for verification.

1190 **6 THREATS TO VALIDITY**

1191 Our work builds on a decade of research on pull-based de-
1192 velopment, extracting the features relevant for pull request
1193 decision-making. In this way, we stand on the shoulders
1194 of giants and hence benefit from it and inherit the limita-
1195 tions of the features they present. In addition, we face the
1196 following limitations and classify them into four categories,
1197 *i.e.*, construct validity, internal validity, external validity, and
1198 conclusion validity [69].

1199 **6.1 Construct Validity**

- 1200 • The measure of relative importance may change if we
1201 choose a different method, which may lead to a dif-
1202 ferent conclusion. There are different ways to calculate
1203 the importance of factors in a logistic regression model,
1204 *e.g.*, the percentage of variance explained by each fac-
1205 tor [45], which is similar to the percentage of total
1206 variance explained by least squares regression [39], the
1207 standardized coefficient [70], and the change in logistic
1208 pseudo partial correlation [71]. This is a research field
1209 in itself and relates to the choice of the algorithm [72],
1210 [73]. To compare the importance of factors in different
1211 models, in this paper, we choose the percentage of
1212 explained variance to represent factor importance. The
1213 choice of the metric may affect the consistency of the

16. <https://github.com/rails/rails/pull/20851>

TABLE 9: The recommended control factors for different contexts

	overall	other-integrated	self-integrated	has comment	no comment	use CI	no CI	early stage of projects
same_user	✓			✓	✓	✓	✓	✓
lifetime_minutes	✓	✓	✓	✓	✓	✓	✓	✓
prior_review_num	✓		✓	✓	✓	✓	✓	✓
has_comments	✓	✓	✓			✓	✓	✓
core_member	✓		✓	✓	✓	✓	✓	✓
num_commits	✓	✓	✓	✓		✓	✓	✓
other_comment	✓	✓	✓	✓		✓		✓
ci_exists	✓	✓			✓			
hash_tag	✓		✓	✓		✓	✓	
account_creation_days		✓			✓			
commits_on_files_touched		✓			✓		✓	✓
reopen_or_not			✓		✓			
open_pr_num			✓		✓			✓
prev_pullreqs			✓					
first_pr			✓					
files_added			✓					
contrib_open			✓					
perc_pos_emotion				✓				
description_length					✓			
ci_failed_perc						✓		
num_comments							✓	
files_changed								✓
followers								✓

Note: ✓ marks the recommended control factors when building logistic regression models for pull request decisions

conclusion to a certain extent. However, as this metric is widely used in many related works [10], [12], [74], our result can reflect the influence of factors on pull request decisions to a certain extent.

- The inconsistency between the GHTorrent dataset and the results returned by the GitHub API brought about errors in the time-related factors, which may influence the results. We checked 100 randomly selected records for each of the four factors *first_response_time*, *account_creation_days*, *project_age*, and *ci_latency*, and the precision was 98%, 97%, 96%, and 94%, respectively. Our dataset has inherited the problems, but from our investigation, the number of errors in our dataset is small compared to the size we have used for analysis.
- A developer may have multiple accounts in GitHub. We did not combine the accounts in our model. However, we analyzed this situation with a relevant tool [75] and found that 94% of the accounts in our dataset corresponds to only one developer. Due to the importance of the factor *same_user* in our model, we examined the reliability of the factor and found that the case of a user having multiple accounts does not affect its accuracy.
- For RQ2, we divided the data according to team size and the closing time of pull requests. This paper does not discuss the robustness of threshold selection, which may lead to less reliable conclusions. However, according to previous studies [52], [53], they split the data into three subsets for the trend analysis. Also, there are infinite ways to select the data division threshold, which can lead to differences in data size for different subsets. While optimizing the differences of data subsets, our result effectively reflects different contexts' influence on

pull request decisions.

6.2 Internal Validity

- The absence of factors may have an impact on the relative importance of factors in the conclusion. We consider factors that can be mined from archival data and exclude those factors, *e.g.*, eye track-related factors [21], that are difficult to quantify in a scalable manner. These factors also include factors that focus only on specific scenarios, *e.g.*, factors related to Microsoft [3] and npm ecosystems only [16]. Because these factors also influence pull request decisions, as mentioned in previous studies, removing them can impact factors' relative importance on affecting pull request decisions. We are not sure how these factors perform together with our collected factors. At least, we have collected as many relevant factors as possible, quantified them, and added them to our dataset. Also, during data preprocessing, we remove the factor *bug_fix* due to 99.3% missing values, and thus, we are not sure how this factor affects pull request decision-making. Although many tools can predict whether a pull request fixes a bug, we only use the manually added label to classify pull requests to ensure data's accuracy. Future studies that want to delve deeper into the impact of these deleted factors can use other tools to complement this data for further analysis.
- There lacks a careful consideration of different types of projects. It is undeniable that when building models, it's better to consider different kinds of projects separately. However, the heterogeneity of projects has many dimensions, not only limited to the code contribution and

1277 review process. Therefore it isn't easy to achieve accu-
1278 rate classification of projects. This paper has considered
1279 the issue of project heterogeneity to some extent, which
1280 includes many project related factors and treats team
1281 size as a project context.

1282 6.3 External Validity

- 1283 • Project selection introduces data bias when building
1284 models, resulting in our conclusions that may not apply
1285 to the complete set of GitHub data or some specific
1286 types of projects. *E.g.*, projects written in programming
1287 languages other than Java, Python, Ruby, JavaScript,
1288 Scala, and Go. Since it is impractical to model using
1289 data from the complete GitHub collection, the diversity
1290 of our data can help avoid this problem to a certain
1291 extent. Similarly, when selecting the projects, we se-
1292 lected the top 3% of projects in terms of the number
1293 of submitted pull requests and filtered out the projects
1294 in which the number of closed pull requests was less
1295 than 20. In Section 2.2, we mentioned that we specified
1296 these thresholds through discussion for the scalability
1297 and validity of the dataset. We cannot guarantee that
1298 our conclusions are available for other projects. We have
1299 at least tried to select the proper set of projects.
- 1300 • The generalizability of our study is not verified in other
1301 social coding platforms (other than GitHub) or other
1302 modern code review tools, *e.g.*, Gerrit. One major reason
1303 for differences can be the factors influencing pull re-
1304 quest decisions on different platforms. The comparison
1305 of factors' influence on contribution decisions under
1306 different platforms or tools belongs to another research
1307 in the future.

1308 6.4 Conclusion Validity

- 1309 • For logistic regression models, comparing the variance
1310 explained by the same factor in different models is
1311 not accurate. This may affect the correctness of the
1312 conclusions, as the variance explained by the factors
1313 in different regression models fluctuates when different
1314 models use different training sets. But there is not a
1315 good solution to the problem. However, in our study,
1316 we consider only the factors that change dramatically in
1317 different contexts. When building models with the same
1318 set of predictors, large changes in explained variance
1319 can be used to describe the change in factor importance.

1320 7 RELATED WORK

1321 The related work of this paper is mainly divided into four
1322 parts. The first subsection introduces modern code review.
1323 The second subsection introduces factors influencing pull
1324 request decisions. Third, we introduce papers that tried to
1325 integrate related factors and explain the relative importance
1326 of the factors influencing pull request decisions. Fourth, we
1327 discuss other studies that have introduced scientific research
1328 methods based on big data.

7.1 Modern Code Review

1329 Although Fagan et al. developed a structure of code in-
1330 spection in 1976 [76], it is very time-consuming and not
1331 applicable in practice [77]. Therefore, modern code review
1332 comes into being, which is informal, tool-based, and occurs
1333 regularly in practice [78].
1334

1335 Many tools or platforms support modern code review.
1336 Different companies and organizations use various tools
1337 and have their policy during the code review process [79].
1338 CRITICS [80], ReviewClique [81], and Mylyn Reviews [82]
1339 are code review tools integrated into IDE, combining the
1340 code review and development process. Another popular
1341 tool called Gerrit [83], which supported many projects in-
1342 cluding Android, OpenStack is a Git-based tool. CodeFlow,
1343 which is similar to Gerrit, is widely used by Microsoft [78].

1344 In recent years, the pull-based development model has
1345 become a new paradigm for distributed software develop-
1346 ment. Many code-hosting sites, notably GitHub, support the
1347 model by integrating it with code review systems [1]. Unlike
1348 Gerrit, pull requests on GitHub focus not only on a single
1349 commit but also on a whole branch [84]. In contrast, pull
1350 request is easy to participate in the contribution process
1351 without having to master many git operations [85]. Its well-
1352 designed user interface and support for social collaboration
1353 help improve the usability and code review process of
1354 GitHub [86]. These characteristics help GitHub get more
1355 than 79 million users and 238 million repositories. Therefore,
1356 we would like to start with GitHub's pull-based model to
1357 explain the factors associated with pull request decisions.

7.2 Factors influencing pull request decisions

1358 The factors influencing pull request decisions can be di-
1359 vided into three categories, namely, developer characteris-
1360 tics, project characteristics and pull request characteristics.
1361

7.2.1 Developer characteristics

1362 Developer characteristics are related to the contributor and
1363 the integrator. This category contains factors related to hu-
1364 man beings and interactions between two contributors or
1365 a contributor and a project. This category includes `basic`
1366 `information` on developers, including their *gender* [87],
1367 *country* information [12], and *affiliation* [88], [89]. Some stud-
1368 ies focus on `personal features`, including the *personality*
1369 and *emotion* of developers [2], [24], while others studied
1370 the *relationship* between the developer and the target
1371 project, including the *experience* of developers, which is
1372 conceptualized as the count of previous pull requests, ac-
1373 cepted commit count [90], days since account creation [91],
1374 whether it is the first pull request of the contributor [52],
1375 [53], the prior reviews of the integrator [89], the *coreness*
1376 of the contributor [10], [15], [52], [59], [92], [93], the *social*
1377 *distance* [15] and *social strength* [10] of contributor to the
1378 integrator, and the *response time* of the integrator to the pull
1379 request [10].
1380

7.2.2 Project characteristics

1381 Studies on project characteristics mainly talk about the
1382 `basic information` of target projects when submitting
1383 pull requests, which can be summarized into the follow-
1384 ing aspects: *programming language* [52], [58], [91], project
1385

1386 *popularity*, measured as watcher count [28], star count [28],
 1387 fork count [54], [91], *age* of the project [10], [15], *workload*
 1388 measured as the number of open pull requests [10], [89],
 1389 *activeness* measured as the time interval in seconds between
 1390 the opening time of the two latest pull requests [54], and
 1391 *openness* measured as the count of open issues [54].

1392 7.2.3 Pull request characteristics

1393 Related works focus on the `basic` information of pull
 1394 requests, which includes the *size of the change* measured
 1395 at the file level, commit level, and code level [10]; the
 1396 *complexity of a pull request* measured as the length of de-
 1397 scription [10]; the *nature of pull requests* measured as bug
 1398 fixes [58], [90], the *test inclusion* of pull requests [10], [15],
 1399 [92], and the *hotness* or relevance of a PR [1], [10], [15], [53],
 1400 [88], [90]. Additionally, some studies focus on the `process`
 1401 information of pull requests generated during the code
 1402 review process, including the *reference* of a contributor, issue
 1403 or pull request [10], [25]; the *conflict* of a pull request [1];
 1404 the *complexity of discussion* [28]; the *emotion in discussion* [24];
 1405 and *CI tool usage* during the review process [10], [11], [26],
 1406 [94], [95].

1407 7.3 Attempts at explaining pull request decisions

1408 Few studies have tried to integrate the factors related to
 1409 pull request decisions and have explored their relative
 1410 importance in predicting outcomes. Gousios et al. [1] first
 1411 collected a set of factors and performed a preliminary
 1412 exploration of relative importance based on the random
 1413 forest method. However, it was in the early stage of this
 1414 study area. Tsay et al. [15] used an explanatory method
 1415 to explore the importance of social and technical factors.
 1416 However, similar to Gousios et al.'s work [1], their work
 1417 also acted as groundbreaking research, leading to the emer-
 1418 gence of many other studies. Since then, a few follow-ups
 1419 have come into being, *e.g.*, personality-related factors [2],
 1420 geographical location [12], and CI-related factors [10]. In
 1421 2020, Dey et al. [16] collected 50 factors of 483,988 pull
 1422 requests based on 4,218 projects. They also used random
 1423 the forest method to determine the important factors in
 1424 predicting the decision. However, they focused only on the
 1425 npm community and gathered factors without conducting
 1426 a systematic literature review. As a result, factors related to
 1427 CI, personality, emotion, geographical, etc., were missing.
 1428 Furthermore, to the best of our knowledge, no study has
 1429 synthesized the existing body of knowledge to empirically
 1430 explain pull request decisions.

1431 7.4 Big-data-based scientific research methods

1432 Big data has provided many research opportunities, for
 1433 which there are mainly two research methods, *i.e.*, data-
 1434 driven and theory-driven methods. Maass et al. [96] dis-
 1435 cussed the difference between these two methods and found
 1436 that the data-driven method first focuses on the data and
 1437 then extracts patterns and forms into theory. However, the
 1438 theory-driven method first comes up with a theory and uses
 1439 data to prove it. Therefore, our study is data driven, finding
 1440 patterns in different subsets of data and forming them into
 1441 theory.

For the process of a data-driven study, Kar et al. [97] sug- 1442
 1443 gested that there are 6 main steps for building up a theory,
 1444 *i.e.*, data acquisition, data conversion, data analysis, factor
 1445 identification, theory development and model validation.

There are many studies in different research areas that 1446
 1447 have used data-driven research methods. For example,
 1448 Greenwood et al. [98] studied the influence of race, gender,
 1449 and socioeconomic status on the incidence rate of human
 1450 immunodeficiency virus (HIV) infection using data from
 1451 12 million patients. Likewise, other previous studies [1],
 1452 [10], [12], [15] on pull request decisions all used data-driven
 1453 methods.

However, for the data acquisition part, previous studies 1454
 1455 focused only on one specific type of factor or several self-
 1456 defined factors. Without including all the related factors, one
 1457 can hardly gain an overall grasp of the influence of all fac-
 1458 tors. Therefore, we conducted a systematic literature review
 1459 in this study. According to Kitchenham et al. [99], a system-
 1460 atic literature review is an important part of evidence-based
 1461 software engineering (EBSE), as it can aggregate all existing
 1462 evidence and provide guidelines for researchers.

1463 8 CONCLUSIONS

This study synthesizes the existing body of knowledge to 1464
 1465 empirically explain pull request decisions. Our mixed effects
 1466 logistic regression models built on large and diverse GitHub
 1467 project data show that a handful of factors (5 to 10) explain
 1468 pull request decisions the most. The most important factor
 1469 influencing pull request decisions is whether the contributor
 1470 and the integrator are the same user, explaining more than
 1471 30% of the variance. Surprisingly, this factor did not surface
 1472 in any of the prior works and is thus a contribution of this
 1473 study. In addition, positive emotions during discussion and
 1474 CI build results become relatively more important when a
 1475 pull request has comments and uses CI tools, respectively.
 1476 Furthermore, we noticed that the use of CI tools replaced the
 1477 function of comments, indicating changes in the influence
 1478 of these factors. We think that this study has empirically
 1479 synthesized an explanation for pull request decisions that is
 1480 useful for research and practice.

1481 ACKNOWLEDGMENT

This work is supported by National Key R&D Program of 1482
 1483 China (2020AAA0103504). Thank you Rahul N. Iyer, Frenk
 1484 van Mil, Celal Karakoc, Leroy Velzel, Daan Groenewegen,
 1485 and Sarah de Wolf for your help in implementation. Thanks
 1486 Mengluan Cai for validating the validity of factor extraction.

1487 REFERENCES

- 1488 [1] G. Gousios, M. Pinzger, and A. v. Deursen, "An exploratory
 1489 study of the pull-based software development model," in
 1490 *Proceedings of the 36th International Conference on Software
 1491 Engineering*, ser. ICSE 2014. New York, NY, USA: Association
 1492 for Computing Machinery, 2014, p. 345–355. [Online]. Available:
 1493 <https://doi.org/10.1145/2568225.2568260>
- 1494 [2] R. N. Iyer, S. A. Yun, M. Nagappan, and J. Hoey, "Effects of
 1495 personality traits on pull request acceptance," *IEEE Transactions
 1496 on Software Engineering*, pp. 1–1, 2019.

- [3] C. Maddila, C. Bansal, and N. Nagappan, "Predicting pull request completion time: a case study on large scale cloud services," in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2019, pp. 874–882.
- [4] J. Jiang, Y. Yang, J. He, X. Blanc, and L. Zhang, "Who should comment on this pull request? analyzing attributes for more accurate commenter recommendation in pull-based development," *Information and Software Technology*, vol. 84, pp. 48–62, 2017.
- [5] Y. Yu, H. Wang, G. Yin, and C. X. Ling, "Who should review this pull-request: Reviewer recommendation to expedite crowd collaboration," in *2014 21st Asia-Pacific Software Engineering Conference*, vol. 1, Dec 2014, pp. 335–342.
- [6] Y. Yu, Z. Li, G. Yin, T. Wang, and H. Wang, "A dataset of duplicate pull-requests in github," in *Proceedings of the 15th International Conference on Mining Software Repositories*, 2018, pp. 22–25.
- [7] Q. Wang, B. Xu, X. Xia, T. Wang, and S. Li, "Duplicate pull request detection: When time matters," in *Proceedings of the 11th Asia-Pacific Symposium on Internetware*, 2019, pp. 1–10.
- [8] Z. Liu, X. Xia, C. Treude, D. Lo, and S. Li, "Automatic generation of pull request descriptions," in *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2019, pp. 176–188.
- [9] E. v. d. Veen, G. Gousios, and A. Zaidman, "Automatically prioritizing pull requests," in *2015 IEEE/ACM 12th Working Conference on Mining Software Repositories*, 2015, pp. 357–361.
- [10] Y. Yu, G. Yin, T. Wang, C. Yang, and H. Wang, "Determinants of pull-based development in the context of continuous integration," *Science China Information Sciences*, vol. 59, no. 8, p. 080104, 2016. [Online]. Available: <https://doi.org/10.1007/s11432-016-5595-8>
- [11] B. Vasilescu, Y. Yu, H. Wang, P. Devanbu, and V. Filkov, "Quality and productivity outcomes relating to continuous integration in github," in *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*, ser. ESEC/FSE 2015. New York, NY, USA: Association for Computing Machinery, 2015, p. 805–816. [Online]. Available: <https://doi.org/10.1145/2786805.2786850>
- [12] A. Rastogi, N. Nagappan, G. Gousios, and A. van der Hoek, "Relationship between geographical location and evaluation of developer contributions in github," in *Proceedings of the 12th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, ser. ESEM '18. New York, NY, USA: Association for Computing Machinery, 2018. [Online]. Available: <https://doi.org/10.1145/3239235.3240504>
- [13] Z. Hu and E. Gehringer, "Use bots to improve github pull-request feedback," in *Proceedings of the 50th ACM Technical Symposium on Computer Science Education*, 2019, pp. 1262–1263.
- [14] Z. Peng and X. Ma, "Exploring how software developers work with mention bot in github," *CCF Transactions on Pervasive Computing and Interaction*, vol. 1, no. 3, pp. 190–203, 2019.
- [15] J. Tsay, L. Dabbish, and J. Herbsleb, "Influence of social and technical factors for evaluating contribution in github," in *Proceedings of the 36th International Conference on Software Engineering*, ser. ICSE 2014. New York, NY, USA: Association for Computing Machinery, 2014, p. 356–366. [Online]. Available: <https://doi.org/10.1145/2568225.2568315>
- [16] T. Dey and A. Mockus, "Which pull requests get accepted and why? a study of popular npm packages," *arXiv preprint arXiv:2003.01153*, 2020.
- [17] B. Kitchenham and S. Charters, "Guidelines for performing systematic literature reviews in software engineering," 2007.
- [18] A. W. Harzing, "The publish or perish book: Your guide to effective and responsible citation analysis," *International Review of Research in Open & Distance Learning*, vol. 13, no. 3, pp. 314–315, 2012.
- [19] M. Gusenbauer, "Google scholar to overshadow them all? comparing the sizes of 12 academic search engines and bibliographic databases," *Scientometrics*, vol. 118, 2019.
- [20] G. Jeong, S. Kim, T. Zimmermann, and K. Yi, "Improving code review by predicting reviewers and acceptance of patches," *Research on software analysis for error-free computing center Tech-Memo (ROSAEC MEMO 2009-006)*, pp. 1–18, 2009.
- [21] D. Ford, M. Behroozi, A. Serebrenik, and C. Parnin, "Beyond the code itself: How programmers really look at pull requests," in *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Society (ICSE-SEIS)*, 2019, pp. 51–60.
- [22] P. Pooput and P. Muenchaisri, "Finding impact factors for rejection of pull requests on github," in *Proceedings of the 2018 VII International Conference on Network, Communication and Computing*, ser. ICNCC 2018. New York, NY, USA: Association for Computing Machinery, 2018, p. 70–76. [Online]. Available: <https://doi.org/10.1145/3301326.3301380>
- [23] M. Ortu, M. Marchesi, and R. Tonelli, "Empirical analysis of affect of merged issues on github," in *2019 IEEE/ACM 4th International Workshop on Emotion Awareness in Software Engineering (SEmotion)*, 2019, pp. 46–48.
- [24] Iyer, Rahul, "Effects of personality traits and emotional factors in pull request acceptance." 2019. [Online]. Available: <http://hdl.handle.net/10012/14952>
- [25] F. Calefato, F. Lanubile, and N. Novielli, "A preliminary analysis on the effects of propensity to trust in distributed software development," in *2017 IEEE 12th International Conference on Global Software Engineering (ICGSE)*, 2017, pp. 56–60.
- [26] G. Gousios, A. Zaidman, M. Storey, and A. v. Deursen, "Work practices and challenges in pull-based development: The integrator's perspective," in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, 2015, pp. 358–368.
- [27] X. Zhang, A. Rastogi, and Y. Yu, "On the shoulders of giants: A new dataset for pull-based development research," in *Proceedings of the 17th International Conference on Mining Software Repositories*, 2020, pp. 543–547.
- [28] G. Gousios and A. Zaidman, "A dataset for pull-based development research," in *Proceedings of the 11th Working Conference on Mining Software Repositories*, ser. MSR 2014. New York, NY, USA: Association for Computing Machinery, 2014, p. 368–371. [Online]. Available: <https://doi.org/10.1145/2597073.2597122>
- [29] B. Vasilescu, A. Capiluppi, and A. Serebrenik, "Gender, representation and online participation: A quantitative study," *Interacting with Computers*, vol. 26, no. 5, pp. 488–511, 2014.
- [30] Q. Fan, Y. Yu, G. Yin, T. Wang, and H. Wang, "Where is the road for issue reports classification based on text mining?" in *2017 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. IEEE, 2017, pp. 121–130.
- [31] "Linking a pull request to an issue," <https://docs.github.com/en/issues/tracking-your-work-with-issues/linking-a-pull-request-to-an-issue>, [Online; accessed 4-November-2021].
- [32] X. Zhang, A. Rastogi, and Y. Yu, "Technical Report," https://github.com/zhangxunhui/new_pullreq_msr2020/blob/master/technical_report.pdf, 2020, [Online; accessed 3-March-2021].
- [33] StatsTest.com, "Cramer's V," <https://www.statstest.com/cramers-v-2/>, [Online; accessed 3-March-2021].
- [34] J. S. Jones, *Learn to Use the Eta Coefficient Test in SPSS With Data From the NIOSH Quality of Worklife Survey (2014)*, 2019.
- [35] J. Cohen, *Statistical power analysis for the behavioral sciences*. Academic press, 1969.
- [36] A. Galecki and T. Burzykowski, "Linear mixed-effects model," in *Linear Mixed-Effects Models Using R*. Springer, 2013, pp. 245–273.
- [37] D. M. Bates, "lme4: Mixed-effects modeling with r," 2010.
- [38] J. Frost, "Multicollinearity in Regression Analysis: Problems, Detection, and Solutions," <https://statisticsbyjim.com/regression/multicollinearity-in-regression-analysis/#:~:text=Multicollinearity%20occurs%20when%20independent%20variables%20in%20a%20regression,you%20fit%20the%20model%20and%20interpret%20the%20results.>, [Online; accessed 6-April-2021].
- [39] J. Cohen, P. Cohen, S. G. West, and L. S. Aiken, *Applied multiple regression/correlation analysis for the behavioral sciences*. Routledge, 2013.
- [40] P. C. Pndharkar and J. A. Rodger, "An empirical study of the impact of team size on software development effort," *Information Technology & Management*, vol. 8, no. 4, pp. 253–262, 2007.
- [41] S. W. Chou and M. Y. He, "The factors that affect the performance of open source software development – the perspective of social capital and expertise integration," *Information Systems Journal*, vol. 21, no. 2, pp. 195–219, 2011.
- [42] X. Zhang, Y. Yu, G. Gousios, and R. Ayushi, "Technical Report," https://github.com/zhangxunhui/TSE_pull-based-development/blob/main/technical_report.pdf, 2021, [Online; accessed 23-March-2022].
- [43] Y. Zhao, A. Serebrenik, Y. Zhou, V. Filkov, and B. Vasilescu, "The impact of continuous integration on other software development practices: A large-scale empirical study," in *2017 32nd IEEE/ACM*

- 1650 *International Conference on Automated Software Engineering (ASE)*,
 1651 2017, pp. 60–71.
- 1652 [44] Ø. Langsrud, “Anova for unbalanced data: Use type ii instead of
 1653 type iii sums of squares,” *Statistics and Computing*, vol. 13, no. 2,
 1654 pp. 163–167, 2003.
- 1655 [45] B. Ray, D. Posnett, V. Filkov, and P. Devanbu, “A large scale
 1656 study of programming languages and code quality in github,” in
 1657 *Proceedings of the 22nd ACM SIGSOFT International Symposium on*
 1658 *Foundations of Software Engineering*, ser. FSE 2014. New York, NY,
 1659 USA: Association for Computing Machinery, 2014, p. 155–165.
 1660 [Online]. Available: <https://doi.org/10.1145/2635868.2635922>
- 1661 [46] A. P. Bradley, “The use of the area under the roc curve in the
 1662 evaluation of machine learning algorithms,” *Pattern recognition*,
 1663 vol. 30, no. 7, pp. 1145–1159, 1997.
- 1664 [47] GitHub, “Approving a pull request with required reviews,” ht
 1665 tps://docs.github.com/en/pull-requests/collaborating-wit
 1666 h-pull-requests/reviewing-changes-in-pull-requests/appro
 1667 ving-a-pull-request-with-required-reviews, [Online; accessed
 1668 18-November-2021].
- 1669 [48] J. Corbet and G. Kroah-Hartman, “2017 linux kernel development
 1670 report,” *A Publication of The Linux Foundation*, 2017.
- 1671 [49] Y. Zhang, M. Zhou, A. Mockus, and Z. Jin, “Companies’ partici
 1672 pation in oss development-an empirical study of openstack,”
 1673 *IEEE Transactions on Software Engineering*, 2019.
- 1674 [50] M. Zhou, A. Mockus, X. Ma, L. Zhang, and H. Mei, “Inflow and
 1675 retention in oss communities with commercial involvement: A
 1676 case study of three hybrid projects,” *ACM Transactions on Software*
 1677 *Engineering and Methodology (TOSEM)*, vol. 25, no. 2, pp. 1–29,
 1678 2016.
- 1679 [51] J.-M. Hoc, *Psychology of programming*. Academic Press, 2014.
- 1680 [52] D. M. Soares, M. L. de Lima Júnior, L. Murta, and A. Plastino,
 1681 “Acceptance factors of pull requests in open-source projects,” in
 1682 *Proceedings of the 30th Annual ACM Symposium on Applied*
 1683 *Computing*, ser. SAC ’15. New York, NY, USA: Association for
 1684 Computing Machinery, 2015, p. 1541–1546. [Online]. Available:
 1685 <https://doi.org/10.1145/2695664.2695856>
- 1686 [53] D. M. Soares, M. L. d. L. Júnior, L. Murta, and A. Plastino,
 1687 “Rejection factors of pull requests filed by core team developers
 1688 in software projects with high acceptance rates,” in *2015 IEEE*
 1689 *14th International Conference on Machine Learning and Applications*
 1690 *(ICMLA)*, Dec 2015, pp. 960–965.
- 1691 [54] N. Khadke, M. H. Teh, and M. Shen, “Predicting acceptance of
 1692 github pull requests,” 2012.
- 1693 [55] S. McIntosh, Y. Kamei, B. Adams, and A. E. Hassan, “The im
 1694 pact of code review coverage and code review participation on
 1695 software quality: A case study of the qt, vtk, and itk projects,” in
 1696 *Proceedings of the 11th Working Conference on Mining Software*
 1697 *Repositories*, 2014, pp. 192–201.
- 1698 [56] D. Legay, A. Decan, and T. Mens, “On the impact of pull request
 1699 decisions on future contributions,” *CoRR*, vol. abs/1812.06269,
 1700 2018. [Online]. Available: <http://arxiv.org/abs/1812.06269>
- 1701 [57] O. Baysal, O. Kononenko, R. Holmes, and M. W. Godfrey,
 1702 “Investigating technical and non-technical factors influencing
 1703 modern code review,” *Empirical Software Engineering*, vol. 21,
 1704 no. 3, pp. 932–959, 2016. [Online]. Available: <https://doi.org/10.1007/s10664-015-9366-8>
- 1705 [58] R. Padhye, S. Mani, and V. S. Sinha, “A study of external
 1706 community contribution to open-source projects on github,” in
 1707 *Proceedings of the 11th Working Conference on Mining Software*
 1708 *Repositories*, ser. MSR 2014. New York, NY, USA: Association
 1709 for Computing Machinery, 2014, p. 332–335. [Online]. Available:
 1710 <https://doi.org/10.1145/2597073.2597113>
- 1711 [59] A. Bosu and J. C. Carver, “Impact of developer reputation
 1712 on code review outcomes in oss projects: An empirical
 1713 investigation,” in *Proceedings of the 8th ACM/IEEE International*
 1714 *Symposium on Empirical Software Engineering and Measurement*,
 1715 ser. ESEM ’14. New York, NY, USA: Association for Computing
 1716 Machinery, 2014. [Online]. Available: <https://doi.org/10.1145/2652524.2652544>
- 1717 [60] A. Lee and J. C. Carver, “Are one-time contributors different?
 1718 a comparison to core and periphery developers in floss reposi
 1719 tories,” in *2017 ACM/IEEE International Symposium on Empirical*
 1720 *Software Engineering and Measurement (ESEM)*, 2017, pp. 1–10.
- 1721 [61] M. Golzadeh, A. Decan, and T. Mens, “On the effect of discus
 1722 sions on pull request decisions,” 2019.
- 1723 [62] H. Khalajzadeh, M. Shahin, H. O. Obie, and J. Grundy, “How
 1724 are diverse end-user human-centric issues discussed on github?”
 1725 *arXiv preprint arXiv:2201.05927*, 2022.
- 1726 [63] J. Tsay, L. Dabbish, and J. Herbsleb, “Let’s talk about it: evaluat
 1727 ing contributions through discussion in github,” in *Proceedings of*
 1728 *the 22nd ACM SIGSOFT international symposium on foundations of*
 1729 *software engineering*, 2014, pp. 144–154.
- 1730 [64] X. Tan and M. Zhou, “How to communicate when submitting
 1731 patches: An empirical study of the linux kernel,” *Proceedings of*
 1732 *the ACM on Human-Computer Interaction*, vol. 3, no. CSCW, pp.
 1733 1–26, 2019.
- 1734 [65] I. Steinmacher, G. Pinto, I. S. Wiese, and M. A. Gerosa, “Almos
 1735 t there: A study on quasi-contributors in open-source software
 1736 projects,” in *2018 IEEE/ACM 40th International Conference on Soft*
 1737 *ware Engineering (ICSE)*, 2018, pp. 256–266.
- 1738 [66] Z. Li, Y. Yu, T. Wang, G. Yin, S. Li, and H. Wang, “Are you still
 1739 working on this an empirical study on pull request abandon
 1740 ment,” *IEEE Transactions on Software Engineering*, 2021.
- 1741 [67] Z. Peng and X. Ma, “Exploring how software developers work
 1742 with mention bot in github,” *CCF Transactions on Pervasive Com*
 1743 *puting and Interaction*, vol. 1, no. 3, pp. 190–203, 2019.
- 1744 [68] X. Lu, Y. Cao, Z. Chen, and X. Liu, “A first look at emoji usage
 1745 on github: An empirical study,” *arXiv preprint arXiv:1812.04863*,
 1746 2018.
- 1747 [69] L. Gren, “Standards of validity and the validity of standards
 1748 in behavioral software engineering research: the perspective of
 1749 psychological test theory,” in *Proceedings of the 12th ACM/IEEE*
 1750 *International Symposium on Empirical Software Engineering and*
 1751 *Measurement*, 2018, pp. 1–4.
- 1752 [70] S. Tonidandel and J. M. LeBreton, “Relative importance analysis:
 1753 A useful supplement to regression analysis,” *Journal of Business*
 1754 *and Psychology*, vol. 26, no. 1, pp. 1–9, 2011.
- 1755 [71] V. Agrawal, “Interpreting importance of features in logistic re
 1756 gression model [closed],” [https://stats.stackexchange.com/ques](https://stats.stackexchange.com/questions/233050/interpreting-importance-of-features-in-logistic-regression-model)
 1757 [tions/233050/interpreting-importance-of-features-in-logistic-re](https://stats.stackexchange.com/questions/233050/interpreting-importance-of-features-in-logistic-regression-model)
 1758 [gression-model](https://stats.stackexchange.com/questions/233050/interpreting-importance-of-features-in-logistic-regression-model), [Online; accessed 24-March-2021].
- 1759 [72] Aliweb, “How to choose the best algorithm for measuring attri
 1760 bute importance/relevance?” [https://stats.stackexchange.co](https://stats.stackexchange.com/questions/251248/how-to-choose-the-best-algorithm-for-measuring-attribute-importance-relevance)
 1761 [m/questions/251248/how-to-choose-the-best-algorithm-for](https://stats.stackexchange.com/questions/251248/how-to-choose-the-best-algorithm-for-measuring-attribute-importance-relevance)
 1762 [measuring-attribute-importance-relevance](https://stats.stackexchange.com/questions/251248/how-to-choose-the-best-algorithm-for-measuring-attribute-importance-relevance), [Online; accessed
 1763 24-March-2021].
- 1764 [73] M. Drury, “What are variable importance rankings useful for?”
 1765 [https://stats.stackexchange.com/questions/202277/what-ar](https://stats.stackexchange.com/questions/202277/what-are-variable-importance-rankings-useful-for)
 1766 [e-variable-importance-rankings-useful-for](https://stats.stackexchange.com/questions/202277/what-are-variable-importance-rankings-useful-for), [Online; accessed
 1767 24-March-2021].
- 1768 [74] C. Overney, J. Meinicke, C. Kastner, and B. Vasilescu, “How to
 1769 not get rich: An empirical study of donations in open ource,”
 1770 *Proceedings - International Conference on Software Engineering*, pp.
 1771 1209–1221, 2020.
- 1772 [75] B. Vasilescu, A. Serebrenik, and V. Filkov, “A data set for social
 1773 diversity studies of GitHub teams,” in *12th Working Conference on*
 1774 *Mining Software Repositories, Data Track*, ser. MSR. IEEE, 2015,
 1775 pp. 514–517.
- 1776 [76] M. Fagan, “Design and code inspections to reduce errors in
 1777 program development,” in *Software pioneers*. Springer, 2002, pp.
 1778 575–607.
- 1779 [77] F. Shull and C. Seaman, “Inspecting the history of inspections: An
 1780 example of evidence-based technology diffusion,” *IEEE software*,
 1781 vol. 25, no. 1, pp. 88–90, 2008.
- 1782 [78] A. Bacchelli and C. Bird, “Expectations, outcomes, and challenges
 1783 of modern code review,” in *2013 35th International Conference on*
 1784 *Software Engineering (ICSE)*. IEEE, 2013, pp. 712–721.
- 1785 [79] A. Bosu, “Modeling modern code review practices in open source
 1786 software development organizations,” in *Proceedings of the 11th*
 1787 *International Doctoral Symposium on Empirical Software Engineering*,
 1788 2013.
- 1789 [80] T. Zhang, M. Song, and M. Kim, “Critics: An interactive code
 1790 review tool for searching and inspecting systematic changes,” in
 1791 *Proceedings of the 22nd ACM SIGSOFT International Symposium on*
 1792 *Foundations of Software Engineering*, 2014, pp. 755–758.
- 1793 [81] M. Bernhart, A. Mauiczka, and T. Grechenig, “Adopting code
 1794 reviews for agile software development,” in *2010 Agile Conference*.
 1795 IEEE, 2010, pp. 44–47.
- 1796 [82] “Mylyn reviews,” [https://projects.eclipse.org/projects/mylyn](https://projects.eclipse.org/projects/mylyn-reviews)
 1797 [reviews](https://projects.eclipse.org/projects/mylyn-reviews), [Online; accessed 4-November-2021].
- 1798 [83] “Gerrit code review,” <https://www.gerritcodereview.com/>,
 1799 [Online; accessed 4-November-2021].

[84] L. Vogel, "Gerrit code review - tutorial," <https://www.vogella.com/tutorials/Gerrit/article.html>, 2020, [Online; accessed 4-November-2021].

[85] "Gerrit code review, or github's fork and pull model?" <https://softwareengineering.stackexchange.com/questions/173262/gerrit-code-review-or-githubs-fork-and-pull-model>, 2012, [Online; accessed 4-November-2021].

[86] "Gerritforge blog - git and gerrit code review supported and delivered to your enterprise," <https://gitenterprise.me/2013/10/17/gerrit-code-review-or-githubs-fork-and-pull-take-both/>, 2013, [Online; accessed 4-November-2021].

[87] J. Terrell, A. Kofink, J. Middleton, C. Rainear, E. Murphy-Hill, C. Parnin, and J. Stallings, "Gender differences and bias in open source: Pull request acceptance of women versus men," *PeerJ Computer Science*, vol. 3, p. e111, 2017.

[88] O. Kononenko, T. Rose, O. Baysal, M. Godfrey, D. Theisen, and B. de Water, "Studying pull request merges: A case study of shopify's active merchant," in *Proceedings of the 40th International Conference on Software Engineering: Software Engineering in Practice*, ser. ICSE-SEIP '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 124–133. [Online]. Available: <https://doi.org/10.1145/3183519.3183542>

[89] O. Baysal, O. Kononenko, R. Holmes, and M. W. Godfrey, "The influence of non-technical factors on code review," in *2013 20th Working Conference on Reverse Engineering (WCRE)*, Oct 2013, pp. 122–131.

[90] Y. Jiang, B. Adams, and D. M. German, "Will my patch make it? and how fast? case study on the linux kernel," in *2013 10th Working Conference on Mining Software Repositories (MSR)*, May 2013, pp. 101–110.

[91] M. M. Rahman and C. K. Roy, "An insight into the pull requests of github," in *Proceedings of the 11th Working Conference on Mining Software Repositories*, ser. MSR 2014. New York, NY, USA: Association for Computing Machinery, 2014, p. 364–367. [Online]. Available: <https://doi.org/10.1145/2597073.2597121>

[92] G. Pinto, L. F. Dias, and I. Steinmacher, "Who gets a patch accepted first? comparing the contributions of employees and volunteers," in *Proceedings of the 11th International Workshop on Cooperative and Human Aspects of Software Engineering*, ser. CHASE '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 110–113. [Online]. Available: <https://doi.org/10.1145/3195836.3195858>

[93] O. Baysal, O. Kononenko, R. Holmes, and M. W. Godfrey, "The secret life of patches: A firefox case study," in *2012 19th Working Conference on Reverse Engineering*, Oct 2012, pp. 447–455.

[94] F. Zampetti, G. Bavota, G. Canfora, and M. D. Penta, "A study on the interplay between pull request review and continuous integration builds," in *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, Feb 2019, pp. 38–48.

[95] Y. Tao, D. Han, and S. Kim, "Writing acceptable patches: An empirical study of open source project patches," in *2014 IEEE International Conference on Software Maintenance and Evolution*, Sep. 2014, pp. 271–280.

[96] W. Maass, J. Parsons, S. Puro, V. C. Storey, and C. Woo, "Data-driven meets theory-driven research in the era of big data: Opportunities and challenges for information systems research," *Journal of the Association for Information Systems*, 2018.

[97] A. K. Kar and Y. K. Dwivedi, "Theory building with big data-driven research – moving away from the "what" towards the "why"," *International Journal of Information Management*, vol. 54, 2020.

[98] B. N. Greenwood and R. Agarwal, "Matching platforms and hiv incidence: An empirical investigation of race, gender, and socioeconomic status," *Management Science*, vol. 62, no. 8, pp. pags. 2281–2303, 2016.

[99] B. Kitchenham, O. P. Brereton, D. Budgen, M. Turner, J. Bailey, and S. Linkman, "Systematic literature reviews in software engineering—a systematic literature review," *Information and software technology*, vol. 51, no. 1, pp. 7–15, 2009.

[100] P. Weißgerber, D. Neu, and S. Diehl, "Small patches get in!" in *Proceedings of the 2008 International Working Conference on Mining Software Repositories*, ser. MSR '08. New York, NY, USA: Association for Computing Machinery, 2008, p. 67–76. [Online]. Available: <https://doi.org/10.1145/1370750.1370767>

[101] J. Jiang, A. Mohamed, and L. Zhang, "What are the characteristics

of reopened pull requests? a case study on open source projects in github," *IEEE Access*, vol. 7, pp. 102751–102761, 2019.

[102] C. Hecht, "On the influence of developer coreness on patch acceptance: A survival analysis," 2020.

[103] V. Kovalenko and A. Bacchelli, "Code review for newcomers: Is it different?" in *Proceedings of the 11th International Workshop on Cooperative and Human Aspects of Software Engineering*, ser. CHASE '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 29–32. [Online]. Available: <https://doi.org/10.1145/3195836.3195842>

[104] M. Hilton, T. Tunnell, K. Huang, D. Marinov, and D. Dig, "Usage, costs, and benefits of continuous integration in open-source projects," in *2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2016, pp. 426–437.



Xunhui Zhang received his BS in Computer Science from Sichuan University in 2015. He received his MS in Software Engineering from National University of Defense Technology in 2017. He is now a PhD candidate in Software Engineering, National University of Defense Technology. His work interests include open source software engineering, data mining, recommendation system, cross community analysis and code clone.



Yue Yu is an associate professor in the College of Computer at National University of Defense Technology (NUDT). He received his Ph.D. degree in Computer Science from NUDT in 2016. He has won Outstanding Ph.D. Thesis Award from Hunan Province. His research findings have been published on ICSE, FSE, ASE, TSE, MSR, IST, ICSME, ICDM and ESEM. His current research interests include software engineering, data mining and computer-supported cooperative work.



Georgios Gousios is a research scientist at Facebook and an associate professor at the Delft University of Technology. His work is on applying techniques from the domains of static analysis, machine learning and software analytics to improve developer productivity and operational efficiency.



Ayushi Rastogi is an Assistant Professor in the Faculty of Science and Engineering at the University of Groningen, the Netherlands. Her research interests include software analytics, empirical software engineering, and mining software repositories. She studies human and social aspects of software engineering for improving developer productivity and promoting diversity and inclusion.

TABLE 10: Factors related to pull request decisions in related articles.

First column lists factors in alphabet ascending order in each class, the rest columns list related articles and the result of each factor.

Horizontal Line in the middle of shape (\ominus) means the factor is removed when building models because of multicollinearity.

Filling: Filled (\bullet) means *significance is reported* and unfilled (\circ) means *significance is not reported because of not using statistical model or inconsistent conclusions*.

Size of filled shape: Big shape (\bullet) shows *statistically significant relation* and small shape (\bullet) *statistically insignificant* with 95% confidence threshold.

Color: Blue \bullet means a *positive relation* (meaning increase in the chances of pull request acceptance), red \bullet means a *negative relation*, gray \bullet means *uncertain relation* because of not using statistical model or nonlinear conclusion.

	[1]	[15]	[10]	[57]	[2]	[24]	[12]	[54]	[100]	[88]	[94]	[52]	[53]	[61]	[58]	[56]	[101]	[91]	[93]	[92]	[59]	[60]	[102]	[103]	[87]	[104]	
Developer Characteristics																											
account_creation_days							\bullet											\circ									
agree_diff					\bullet	\bullet																					
cons_diff					\bullet	\bullet																					
contrib_affiliation				\circ						\bullet																	
contrib_agree					\bullet	\bullet																					
contrib_cons					\bullet	\bullet																					
contrib_country							\circ																				
contrib_extra					\bullet	\bullet																					
contrib_first_emo						\bullet																					
contrib_follow_integrator		\bullet			\bullet	\bullet	\bullet																				
contrib_gender																										\bullet	
contrib_neur					\bullet	\bullet																					
contrib_open					\bullet	\bullet																					
contrib_rate_author											\circ																
core_member		\bullet	\bullet		\bullet	\bullet	\bullet				\circ	\circ							\circ	\circ	\bullet	\bullet	\circ				
extra_diff					\bullet	\bullet														\circ	\circ	\bullet	\bullet	\circ			
first_pr													\circ	\circ								\bullet		\circ			
first_response_time				\bullet																		\bullet			\circ		
followers		\bullet	\bullet		\bullet	\bullet	\bullet																				
inte_affiliation				\circ																							
inte_agree					\bullet	\bullet																					
inte_cons					\bullet	\bullet																					
inte_extra					\bullet	\bullet																					
inte_first_emo						\bullet																					
inte_neur					\bullet	\bullet																					
inte_open					\bullet	\bullet																					
neur_diff					\bullet	\bullet																					
open_diff					\bullet	\bullet																					
perc_contrib_neg_emo						\bullet																					
perc_contrib_pos_emo						\bullet																					
perc_inte_neg_emo						\bullet																					
perc_inte_pos_emo						\bullet																					

TABLE 10: Factors related to pull request decisions in related articles.

First column lists factors in alphabet ascending order in each class, the rest columns list related articles and the result of each factor.

Horizontal Line in the middle of shape (\ominus) means the factor is removed when building models because of multicollinearity.

Filling: Filled (\bullet) means *significance is reported* and unfilled (\circ) means *significance is not reported because of not using statistical model or inconsistent conclusions*.

Size of filled shape: Big shape (\bullet) shows *statistically significant* relation and small shape (\bullet) *statistically insignificant* with 95% confidence threshold.

Color: Blue \bullet means a *positive relation* (meaning increase in the chances of pull request acceptance), red \bullet means a *negative relation*, gray \bullet means *uncertain relation* because of not using statistical model or nonlinear conclusion.

	[1]	[15]	[10]	[57]	[2]	[24]	[12]	[54]	[100]	[88]	[94]	[52]	[53]	[61]	[58]	[56]	[101]	[91]	[93]	[92]	[59]	[60]	[102]	[103]	[87]	[104]	
prev_pullreqs	\circ			\bullet			\bullet			\bullet						\circ											
prior_interaction		\bullet			\bullet	\bullet																					
prior_review_num				\bullet																							
requester_succ_rate	\circ						\bullet	\circ			\circ																
same_affiliation				\circ																							
same_country							\bullet																				
social_strength			\bullet																								
Project Characteristics																											
asserts_per_kloc	\circ						\ominus																				
fork_num		\ominus								\circ										\circ							
integrator_availability				\bullet																							
language												\circ															
open_issue_num																\circ											
open_pr_num				\bullet	\circ																						
perc_external_contribs	\circ										\bullet																
project_age		\bullet	\bullet		\bullet	\bullet	\bullet																				
pr_succ_rate																\circ											
pushed_delta											\circ																
sloc	\circ										\bullet																
stars		\bullet			\bullet	\bullet	\bullet				\bullet					\circ											
team_size	\circ	\bullet	\bullet		\bullet	\bullet	\ominus	\circ																			
test_cases_per_kloc	\ominus																										
test_lines_per_kloc	\circ						\bullet																				
Pull Request Characteristics																											
at_tag			\bullet																								
bug_fix											\circ					\circ											
churn_addition				\bullet																							
churn_deletion				\bullet																							
ci_build_num																											
ci_exists																										\bullet	
ci_failed_perc											\circ																
ci_first_build_status											\bullet																
ci_last_build_status											\bullet																

TABLE 10: Factors related to pull request decisions in related articles.

First column lists factors in alphabet ascending order in each class, the rest columns list related articles and the result of each factor.

Horizontal Line in the middle of shape (\ominus) means the factor is removed when building models because of multicollinearity.

Filling: Filled (\bullet) means *significance is reported* and unfilled (\circ) means *significance is not reported because of not using statistical model or inconsistent conclusions*.

Size of filled shape: Big shape (\bullet) shows *statistically significant* relation and small shape (\bullet) *statistically insignificant* with 95% confidence threshold.

Color: Blue \bullet means a *positive relation* (meaning increase in the chances of pull request acceptance), red \bullet means a *negative relation*, gray \bullet means *uncertain relation* because of not using statistical model or nonlinear conclusion.

	[1]	[15]	[10]	[57]	[2]	[24]	[12]	[54]	[100]	[88]	[94]	[52]	[53]	[61]	[58]	[56]	[101]	[91]	[93]	[92]	[59]	[60]	[102]	[103]	[87]	[104]	
ci_latency			\bullet																								
ci_test_passed			\bullet																								
comment_conflict	\circ																										
commits_on_files_touched	\circ		\bullet																								
contrib_comment														\circ													
description_length			\bullet																								
files_added												\circ															
files_changed	\circ	\bullet			\ominus	\ominus	\bullet	\circ		\ominus	\circ	\circ	\circ														
files_deleted												\circ															
friday_effect			\bullet																								
has_comments													\circ	\circ													
has_exchange														\circ													
hash_tag	\circ		\bullet																								
has_participants														\circ													
inte_comment														\circ													
lifetime_minutes											\circ	\circ				\circ											
num_code_comments										\ominus	\circ																
num_code_comments_con										\ominus																	
num_comments	\circ	\bullet	\bullet		\bullet	\bullet	\bullet			\ominus	\circ																
num_comments_con										\bullet																	
num_commits	\circ		\bullet					\circ		\ominus	\circ	\circ	\circ														
num_participants	\circ									\bullet																	
other_comment														\circ													
part_num_code										\bullet																	
perc_neg_emotion						\bullet																					
perc_pos_emotion						\bullet																					
reopen_or_not																							\circ				
core_comment														\circ													
src_churn	\circ	\bullet		\bullet	\bullet	\bullet	\bullet	\circ	\circ	\bullet																	
test_churn	\circ																										
test_inclusion		\bullet	\bullet		\bullet	\bullet	\bullet					\circ															

1932 APPENDIX A
1933 RESULTS OF DIFFERENT CONTEXTS

TABLE 11: Results in different contexts

	Dependent variable: merged_or_not=1											
	same user or not		has comments or not		ci exists or not		different team sizes			different periods		
	yes	no	yes	no	yes	no	small	mid	large	before 2016.6	2016.6-2018.6	after 2018.6
(Intercept)	10.4***	34.4***	13.1***	42.4***	20.4***	13.3***	24.9***	20.7***	15.9***	6.9***	16.6***	7.1***
prior_review_num	2.86*** [31.17]	0.98*** [0.04]	1.51*** [14.40]	1.91*** [22.12]	1.53*** [13.76]	1.53*** [11.95]	1.59*** [11.27]	1.41*** [8.80]	1.57*** [18.89]	1.30*** [6.25]	1.63*** [14.12]	1.72*** [17.00]
lifetime_minutes	0.66*** [19.09]	0.52*** [43.67]	0.61*** [29.79]	0.70*** [12.97]	0.60*** [21.52]	0.61*** [20.78]	0.54*** [24.47]	0.61*** [20.16]	0.67*** [16.55]	0.65*** [20.26]	0.57*** [20.83]	0.62*** [12.69]
core_member	1.26*** [9.36]	1.13*** [1.31]	1.29*** [5.55]	1.33*** [5.85]	1.30*** [5.28]	1.26*** [4.66]	1.42*** [5.90]	1.28*** [5.24]	1.19*** [3.24]	1.27*** [5.99]	1.34*** [4.94]	1.29*** [3.14]
prev_pullreqs	0.61*** [5.85]	1.17*** [1.24]	1.13*** [0.69]	0.95*** [0.09]	1.14*** [0.64]	1.12*** [0.56]	1.21*** [0.79]	1.21*** [1.37]	1.13*** [0.84]	1.08*** [0.29]	1.26*** [1.59]	1.24*** [1.32]
num_commits	1.23*** [3.78]	1.46*** [10.43]	1.49*** [11.33]	0.98** [0.03]	1.32*** [4.85]	1.25*** [3.57]	1.36*** [4.64]	1.31*** [4.53]	1.26*** [4.36]	1.18*** [2.31]	1.32*** [4.13]	1.36*** [4.84]
hash_tag	1.14*** [2.52]	1.10*** [1.34]	1.14*** [2.27]	1.09*** [0.65]	1.12*** [1.46]	1.10*** [1.14]	1.13*** [1.37]	1.11*** [1.16]	1.11*** [1.38]	1.04*** [0.22]	1.13*** [1.44]	1.13*** [1.30]
first_pr	0.91*** [1.82]	0.96*** [0.30]	0.95*** [0.32]	0.95*** [0.31]	0.94*** [0.45]	0.96*** [0.27]	0.95*** [0.23]	0.97*** [0.15]	0.96*** [0.31]	0.95*** [0.50]	0.96*** [0.20]	0.94*** [0.32]
files_added	0.88*** [1.57]	0.95*** [0.30]	0.90*** [0.83]	0.92*** [0.57]	0.90*** [0.90]	0.92*** [0.53]	0.90*** [0.52]	0.89*** [1.04]	0.90*** [0.94]	0.97*** [0.11]	0.88*** [0.95]	0.86*** [1.34]
reopen_or_not	0.93*** [1.52]	0.99*** [0.01]	0.97*** [0.17]	0.92*** [2.39]	0.96*** [0.40]	0.97*** [0.24]	0.97*** [0.23]	0.96*** [0.48]	0.96*** [0.43]	0.99*** [0.03]	0.97*** [0.18]	0.93*** [1.06]
open_pr_num	0.73 [1.37]	1.06*** [0.05]	0.92*** [0.09]	0.60 [3.07]	0.83*** [0.33]	0.76*** [0.93]	0.81 [1.26]	0.83*** [0.81]	0.98 [0.01]	0.77 [1.34]	0.72*** [0.53]	0.75*** [0.28]
contrib_open	1.13*** [1.28]	1.06*** [0.38]	1.05*** [0.25]	1.09*** [0.24]	1.05*** [0.15]	1.07*** [0.34]	1.05*** [0.12]	1.12*** [0.70]	1.05*** [0.25]	1.07*** [0.46]	1.07*** [0.27]	1.02*** [0.02]
description_length	1.06*** [0.45]	1.02*** [0.05]	1.01*** [0.03]	1.12*** [1.11]	1.04*** [0.17]	1.04*** [0.13]	1.03*** [0.08]	1.03*** [0.07]	1.05*** [0.29]	0.99* [0.01]	1.06*** [0.26]	1.07*** [0.36]
commits_on_files_touched	1.06*** [0.41]	1.11*** [1.16]	1.05*** [0.23]	1.18*** [1.86]	1.06*** [0.24]	1.13*** [1.39]	1.10*** [0.61]	1.12*** [0.90]	1.05*** [0.20]	1.30*** [7.25]	0.99 [0.00]	0.99 [0.00]
stars	0.83*** [0.40]	0.94*** [0.05]	0.83*** [0.39]	0.83*** [0.37]	0.85*** [0.24]	0.75*** [0.81]	0.81*** [0.69]	0.89*** [0.17]	1.05** [0.01]	0.83*** [0.45]	0.72*** [0.44]	0.71*** [0.50]
project_age	1.08*** [0.19]	1.24*** [1.37]	1.08*** [0.16]	1.16*** [0.53]	1.10*** [0.20]	1.15*** [0.50]	1.04*** [0.05]	1.08*** [0.15]	0.98** [0.01]	0.89*** [0.40]	1.69*** [2.27]	3.80*** [4.81]
files_changed	0.94*** [0.18]	0.90*** [0.54]	0.95*** [0.11]	0.90*** [0.43]	0.94*** [0.15]	0.90*** [0.47]	0.93*** [0.16]	0.91*** [0.33]	0.93*** [0.21]	0.86*** [1.13]	0.95*** [0.09]	0.94*** [0.13]
test_churn	1.05*** [0.15]	1.09*** [0.52]	1.08*** [0.39]	0.99 [0.99]	1.07*** [0.27]	1.05*** [0.15]	1.11*** [0.45]	1.07*** [0.25]	1.04*** [0.11]	1.07*** [0.34]	1.10*** [0.40]	1.08*** [0.23]
account_creation_days	1.03*** [0.13]	1.11*** [1.70]	1.05*** [0.28]	1.17*** [2.26]	1.08*** [0.58]	1.04*** [0.22]	1.06*** [0.33]	1.11*** [1.04]	1.02*** [0.06]	0.99* [0.01]	1.02*** [0.02]	1.03*** [0.06]
team_size	0.94*** [0.08]	1.09*** [0.19]	1.06*** [0.07]	0.92*** [0.14]	1.02* [0.00]	0.96* [0.02]	1.06*** [0.30]	1.00 [0.00]	0.96*** [0.04]	0.88*** [0.37]	1.09*** [0.07]	0.85*** [0.16]
pushed_delta	1.02*** [0.07]	1.06*** [0.46]	1.03*** [0.15]	1.06*** [0.38]	1.04*** [0.17]	1.04*** [0.18]	1.06*** [0.31]	1.05*** [0.25]	1.02*** [0.08]	1.04*** [0.26]	1.03*** [0.05]	1.04*** [0.12]
integrator_availability	0.98*** [0.07]	1.03*** [0.15]	1.00 [0.00]	0.99 [0.99]	0.99*** [0.03]	1.01* [0.03]	1.03*** [0.07]	1.01 [0.00]	0.97*** [0.14]	1.00 [0.00]	0.97*** [0.06]	0.98** [0.03]
test_inclusion	1.03*** [0.06]	1.00 [0.00]	1.02*** [0.02]	1.02* [0.02]	1.03*** [0.05]	0.97*** [0.07]	1.00 [0.00]	1.03*** [0.05]	1.01** [0.02]	1.00 [0.00]	1.01* [0.01]	1.01 [0.00]
contrib_neur	1.02*** [0.04]	1.05*** [0.27]	1.01* [0.01]	1.04*** [0.05]	1.01 [0.00]	1.03*** [0.05]	1.00 [0.00]	1.07*** [0.25]	0.99** [0.02]	1.02*** [0.03]	1.02*** [0.03]	0.99 [0.00]
contrib_cons	1.02*** [0.04]	1.04*** [0.17]	1.05*** [0.20]	0.94*** [0.12]	1.03*** [0.05]	1.02** [0.04]	1.05*** [0.12]	1.03*** [0.04]	1.03*** [0.07]	1.01** [0.02]	1.03*** [0.05]	1.12*** [0.54]
contrib_gender	0.98*** [0.04]	0.97*** [0.13]	0.97*** [0.10]	0.99 [0.99]	0.98*** [0.06]	0.98** [0.04]	0.97*** [0.08]	0.97*** [0.09]	0.99** [0.02]	0.99 [0.01]	0.97*** [0.10]	1.01 [0.00]
pr_succ_rate	0.98*** [0.04]	0.99* [0.01]	0.98*** [0.03]	0.97*** [0.05]	0.97*** [0.06]	1.02** [0.04]	0.94*** [0.25]	0.99 [0.01]	0.96*** [0.10]	0.97*** [0.14]	1.00 [0.00]	1.11*** [0.13]
contrib_agree	0.98*** [0.04]	0.99* [0.01]	0.98*** [0.03]	0.96*** [0.04]	0.99*** [0.01]	0.97*** [0.07]	0.96*** [0.09]	0.99 [0.00]	0.98*** [0.02]	0.97*** [0.05]	0.96*** [0.08]	1.00 [0.00]
friday_effect	1.01*** [0.03]	1.01* [0.01]	1.01*** [0.03]	1.01 [0.01]	1.01** [0.01]	1.02** [0.06]	1.01 [0.00]	1.01*** [0.02]	1.01*** [0.02]	1.02*** [0.05]	1.01* [0.01]	1.01 [0.00]
contrib_extra	0.98*** [0.03]	0.99* [0.01]	0.98*** [0.05]	1.08*** [0.18]	0.99* [0.00]	1.00 [0.00]	0.98*** [0.02]	0.99 [0.00]	1.00 [0.00]	0.99** [0.02]	0.97*** [0.06]	0.95*** [0.11]
src_churn	0.99** [0.02]	1.01 [0.01]	1.05*** [0.15]	0.90*** [0.65]	1.00 [0.00]	1.00 [0.00]	1.05*** [0.10]	1.00 [0.00]	0.96*** [0.10]	0.98*** [0.05]	1.01 [0.00]	1.00 [0.00]
open_issue_num	0.96** [0.02]	1.15*** [0.22]	0.99 [0.99]	1.12*** [0.13]	1.07*** [0.04]	0.90*** [0.15]	1.02 [0.01]	1.03** [0.01]	0.96 [0.01]	0.94*** [0.04]	1.07** [0.02]	1.05 [0.01]
sloc	1.02* [0.01]	1.04*** [0.03]	1.04*** [0.04]	0.96** [0.04]	0.99 [0.00]	1.00 [0.00]	1.00 [0.00]	0.98** [0.01]	1.07*** [0.08]	1.13*** [0.33]	0.92*** [0.07]	0.94*** [0.05]
files_deleted	0.99* [0.01]	0.98*** [0.08]	0.96*** [0.18]	1.05*** [0.28]	0.98*** [0.06]	1.00 [0.00]	0.97*** [0.06]	0.98*** [0.04]	0.99** [0.01]	1.02*** [0.05]	0.98*** [0.04]	0.97*** [0.07]
test_lines_per_kloc	0.98* [0.01]	0.99 [0.00]	1.03*** [0.02]	0.93*** [0.14]	1.01 [0.00]	0.92*** [0.17]	0.98** [0.01]	1.02* [0.01]	0.98 [0.01]	1.09*** [0.22]	0.95*** [0.05]	0.94*** [0.06]
followers	1.00 [0.00]	0.96*** [0.18]	1.03*** [0.06]	1.05*** [0.12]	1.04*** [0.12]	1.02** [0.04]	1.04*** [0.07]	1.01 [0.00]	1.07*** [0.39]	1.14*** [1.40]	1.06*** [0.16]	1.04*** [0.07]
has_comments	0.68*** [10.43]	0.50*** [27.35]	-	-	0.65*** [10.11]	0.52*** [24.50]	0.57*** [13.21]	0.64*** [10.27]	0.66*** [11.93]	0.63*** [14.94]	0.62*** [9.60]	0.55*** [13.78]
other_comment	1.24*** [5.84]	1.18*** [3.60]	-	-	1.23*** [4.18]	1.08*** [0.71]	1.31*** [6.33]	1.23*** [4.15]	1.14*** [1.92]	1.14*** [2.25]	1.22*** [3.07]	1.25*** [2.84]
ci_exists	1.11*** [0.95]	1.19*** [2.52]	1.14*** [1.64]	1.16*** [1.16]	-	-	1.13*** [0.76]	1.12*** [0.79]	1.09*** [0.66]	1.08*** [0.64]	1.12*** [0.44]	1.17*** [0.66]
num_comments	1.12*** [0.88]	0.96*** [0.11]	-	-	1.01* [0.00]	1.18*** [1.38]	0.91*** [0.37]	1.01* [0.01]	1.13*** [0.79]	1.08*** [0.31]	1.02*** [0.01]	1.07*** [0.16]
comment_conflict	1.01*** [0.03]	1.01*** [0.04]	-	-	1.01* [0.01]	1.02** [0.06]	1.01** [0.02]	1.00 [0.00]	1.01*** [0.02]	1.01*** [0.04]	1.00 [0.00]	1.02*** [0.04]
same_user	-	-	0.56*** [29.27]	0.42*** [41.50]	0.51*** [32.75]	0.59*** [23.27]	0.49*** [24.47]	0.49*** [36.03]	0.55*** [32.58]	0.57*** [31.20]	0.46*** [33.20]	0.46*** [28.79]
inte_open	-	-	1.10*** [0.61]	1.01 [0.01]	1.10*** [0.51]	1.04*** [0.07]	0.98* [0.01]	0.92*** [0.25]	1.18*** [2.12]	0.97*** [0.05]	1.03*** [0.04]	1.25*** [2.21]
inte_neur	-	-	1.03*** [0.05]	0.99 [0.01]	1.06*** [0.15]	0.93*** [0.24]	0.96*** [0.05]	0.93*** [0.18]	1.10*** [0.58]	0.96*** [0.11]	0.98* [0.01]	1.13*** [0.57]
inte_agree	-	-	1.00 [0.00]	1.05*** [0.06]	1.00 [0.00]	1.07*** [0.14]	1.03*** [0.02]	1.01 [0.00]	0.98*** [0.03]	1.02** [0.02]	1.02** [0.01]	0.92*** [0.16]
inte_extra	-	-	1.01 [0.00]	0.97* [0.02]	1.02*** [0.01]	0.97** [0.03]	1.06*** [0.11]	1.07*** [0.21]	0.95*** [0.14]	1.02*** [0.03]	1.04*** [0.06]	1.02 [0.02]
inte_cons	-	-	1.00 [0.00]	1.05*** [0.06]	1.01 [0.00]	0.98* [0.02]	1.01 [0.00]	1.00 [0.00]	0.99** [0.01]	1.01 [0.01]	1.03*** [0.02]	0.97** [0.03]
Observations	950,985	1,010,937	1,152,714	809,208	1,611,277	350,645	601,460	703,396	701,900	512,707	585,401	274,121
AUC_train	0.862	0.874	0.837	0.872	0.843	0.884	0.877	0.843	0.837	0.850	0.867	0.879

© 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.