

An Empirical Study of Sentiments in Code Reviews

Ikram El Asri, Nouredine Kerzazi, Gias Uddin, Foutse Khomh, M. A.
Janati Idrissi

*Mohammed V University in Rabat, Morocco, ENSIAS.
Polytechnique Montreal, Canada.*

Abstract

Context: Modern code reviews are supported by tools to enhance developers' interactions allowing contributors to submit their opinions for each committed change in form of comments. Although the comments are aimed at discussing potential technical issues, the text might enclose harmful sentiments that could erode the benefits of suggested changes.

Objective: In this paper, we study empirically the impact of sentiment embodied within developers' comments on the time and outcome of the code review process.

Method: Based on historical data of four long-lived Open Source Software (OSS) projects from a code review system we investigate whether perceived sentiments have any impact on the interval time of code changes acceptance.

Results: We found that (1) contributors frequently express positive and negative sentiments during code review activities; (2) the expressed sentiments differ among the contributors depending on their position within the social network of the reviewers (*e.g.*, core vs peripheral contributors); (3) the sentiments expressed by contributors tend to be neutral as they progress from the status of newcomer in an OSS project to the status of core team contributors; (4) the reviews with negative comments on average took more time to complete than the reviews with positive/neutral comments, and (5) the reviews with controversial comments took significantly longer time in one project.

Conclusion: Through this work, we provide evidences that text-based sentiments have an impact on the duration of the code review process as well as the acceptance or rejection of the suggested changes.

Keywords: Empirical Software Engineering, Code review, Sentiment Analysis, Opinion Mining, Affective Analysis, Propensity Score Matching.

1. Introduction

Peer code review is the practice where a developer submits a piece of code (*i.e.*, code changes) to peers to judge its eligibility to be integrated into the main project code-base [1]. It aims to assess the quality of source code changes made by contributors before they are integrated into the mainstream. Beyond technical information, the textual comments of reviews could contain either positive or negative sentiments, which might alter the perception of their benefits. Past studies have shown that mailing lists of virtual communities include not only useful information such as ideas for improvements, but also contributor opinions, and feelings about the introduced changes [2]. There are also evidences that developers’ opinions play a key role in the decision-making process of source code reviews [3, 4, 5]. However, little is known about the impact of the expressed sentiments on the effectiveness of the review process.

Previously, Baysal *et al.* [6] have explored the impact of technical and non-technical factors on the duration of source code reviews. They observed that non-technical factors, such as reviewer experience can significantly impact code review outcomes. An empirical understanding of the impact of sentiments in code review process can add a novel dimension to the findings of Baysal *et al.* [6] – notably to guide the design of better code review approaches and tools to facilitate improved productivity.

With a view to understand the prevalence and impact of sentiments in modern code reviews, we empirically studied the code reviews of four long-lived software projects. In particular, we answer four research questions:

RQ1: What is the performance of Sentiment Detectors When Applied on Code Reviews?

Recent studies [7, 8, 9] have raised uncertainties related to the unsuccessful application of sentiment analysis tools for software engineering. Indeed, existing tools might require customization to satisfy needs of a specific usage context such as technical software engineering. Following Novielli *et al.* [10], we carried out a benchmark-based study of three sentiment detection tools that are widely used in software engineering research (*Senti4SD* [11], *SentiCR* [12], and *Sentistrength_SE* [13]). We found that Senti4SD tool provides the best performance (F1 79) when

35 applied to our code review samples datasets. We used Senti4SD [11] in
36 our subsequent analysis.

37 **RQ2: How Prevalent are Sentiments in Code Reviews?**

38 We found that contributors express sentiments in their review com-
39 ments (**13.94% of comments were positive, 2.24% negative, and**
40 **83.81% were identified as neutral**). We observed that both core
41 and peripheral contributors do express sentiments in the code reviews.
42 Core members are those developers that contribute intensively and con-
43 sistently to the OSS project, and thus, lead the community, while pe-
44 ripheral ones are occasional contributors with less frequent commits.
45 We built Social Network Graphs of reviewers to segregate Core and
46 Peripheral contributors. Our analysis reveals that the sentiments of
47 Core contributors tend to become more neutral over time.

48 **RQ3: How do the presence of sentiments in code reviews correlate** 49 **with the outcome of the reviews?**

50 We examined the effect of sentiments on the outcome of code reviews.
51 We observed that reviews with negative comments on average take
52 longer time to complete. In contrast, the reviews with positive senti-
53 ments had a lower duration. Reviews that contain positive sentiments
54 required, on average, **1.32 day less time** to be closed than those with
55 negative sentiments. Moreover, we found that **91.81%** of successful re-
56 views were identified with positive sentiments, and **64.44%** of aborted
57 reviews contained negative sentiments.

58 **Contributions.** This paper makes the following contributions:

- 59 1. We provide empirical evidence on the effect of expressed sentiments
60 on the outcome of code reviews. Providing stakeholders with a bet-
61 ter understanding of the impact of contributors' sentiments on team
62 dynamics and their productivity;
- 63 2. We investigate whether the core (*i.e.*, experienced) developers and the
64 peripheral (*i.e.*, newcomers) developers express different types of sen-
65 timents and the effect of these sentiments on the efficiency of code
66 reviews;

67 3. We monitor the evolvement of sentiments of the top 5% contributors
68 across time, for four OSS projects, as they progress and gain more
69 experience, aiming at understanding the correlation between notoriety
70 (*i.e.*, experience) and the trend of sentiments expressed in text-based
71 interactions.

72 **Paper organization.** Section 2 provides background information on sen-
73 timent analysis, the code review process, and the social network analyses
74 conducted in this paper. Section 3 discusses the related literature. Section 4
75 describes the methodology of our case study. Section 5 reports our findings.
76 Section 6 discusses our results. Section 7 highlights threats to the validity
77 of our study and Section 8 concludes the paper and outlines directions for
78 future work.

79 2. Background

80 This section provides background information about sentiment analysis,
81 code review, and social network analysis.

82 2.1. What Does Sentiment Analysis Stand for?

83 Emotion and sentiment are terms relating to human subjectivity [14] un-
84 derstood in the same way and used interchangeably in different domains even
85 if they are not synonymous. Sentiment detection focuses on the detection
86 of subjectivity in a given input (*e.g.*, a sentence). A subjectivity can be of
87 three types: (1) Positive, (2) Negative, and (3) Neutral. Emotion detection
88 focuses on a finer-grained detection of the underlying expressions carried over
89 by the sentiments, such as, anger, frustration. Gerrod Parrott identified six
90 prominent emotions in social psychology [15]: (1) Joy, (2) Love, (3) Surprise,
91 (4) Anger, (5) Sadness, and (6) Fear. This paper focuses on the analysis of
92 sentiments in code reviews, because sentiment detection is predominantly
93 used in other domains (*e.g.*, cars, movies) to mine and summarize opinions
94 about entities [16]. Although, analyzing sentiments and emotions in text
95 data similarly related to one another, actually the granularity is quite dif-
96 ferent. For example, "*this new feature wasn't what I expected*" and "*I hate*
97 *using this API with buggy source code*" are both negative sentiments. While
98 a Sentiment Analysis seeks to catch the general feel or impression people get
99 from consuming a piece of content, Emotion Analysis stresses the specific
100 articulate emotions such as happy, angry, sad, etc.

101 Sentiment analysis can be performed typically at one of the three levels:
102 document level, sentence level, feature level [17]. In this study, we perform
103 a document level analysis.

104 *2.2. Modern Code Review Practice*

105 Code change review is a well-established practice to improve code quality
106 in software engineering. Developers read and assess each other’s code change
107 before it is integrated into the mainstream line of code towards a release.
108 Gerrit¹ is one of the tools providing infrastructure for online reviews as a
109 substitute to face-to-face meetings or mailing lists. It is an online tool that
110 supports the traceability of the code review process by explicitly linking
111 changes to a software system recorded in a Version Control System (VCS)
112 to their respective code review discussions.

113 Figure 1 illustrates the overall process underpinning the code review flow
114 into Gerrit tool. There are three roles into Gerrit: Author, Reviewer, and
115 Verifier as shown in Figure 1. Authors commit code changes into VCS and re-
116 quest a review. Reviewers are responsible for passing throughout the changes
117 and then proposing and discussing adjustments within comments. In other
118 words, reviewers might spot potential defects that authors are not consciously
119 aware of. Then, the author addresses the comments and produces a new code
120 revision. Verifiers are responsible for executing tests to ensure that proposed
121 changes are bug-free and do not cause any regression of the system. They
122 can also leave comments to describe verification issues that they might en-
123 counter during testing. Once the criteria for a review are satisfied, changes
124 are integrated into the mainstream repository and flagged as “Merged”. This
125 lifecycle may have another different transition “Abandoned” when the review
126 has not passed the evaluation and is no longer active.

127 *2.3. Code Review Factors*

128 One of the main concern of developers when submitting patches for code
129 review is maximizing the chances of their patches being examined in the
130 shortest possible time. However, the outcome and duration of the code re-
131 view process can be affected by a variety of technical factors. These influ-
132 encing factors might introduce some bias when analyzing the real effect of
133 contributor’s sentiments on review fixing time and review outcome.

¹<https://www.gerritcodereview.com/>

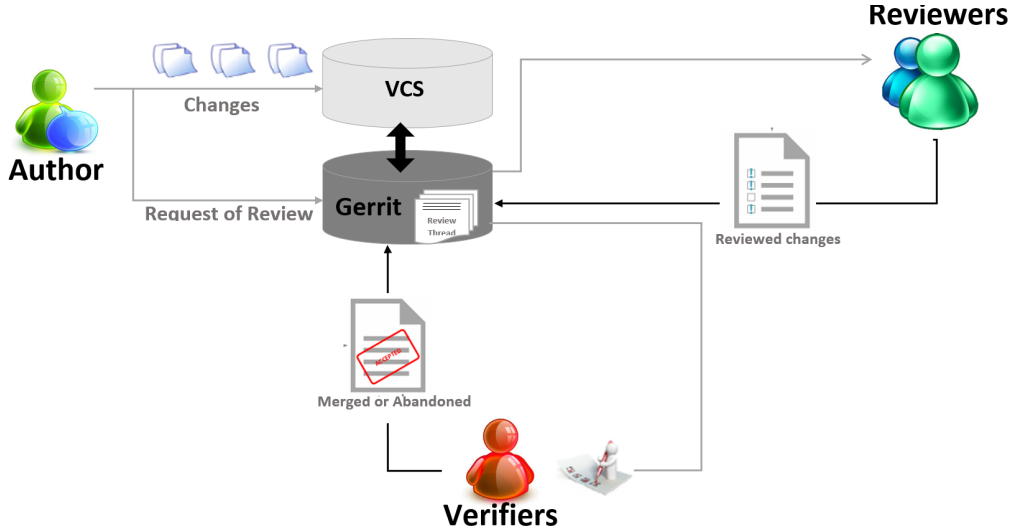


Figure 1: Code Review Flow

134 The most intuitive factor is patch size (Churn); previous studies have
 135 found that smaller patches are more likely to receive faster responses [18, 19]
 136 since larger patches would be more difficult to review, and hence require
 137 more time. Another important factor is how many times a developer had
 138 to resubmit his patch for an additional review (Count Patches); a patch
 139 requiring multiple revisions and re-submission(s) before being accepted con-
 140 sumes more time. Moreover, the more wide-spread a change is across files
 141 (Edited Files), the more concepts it touches in a system which often results
 142 in more rework [20]. Based on a survey of 88 open source core developers,
 143 Kononenko et al. [21] confirmed the important influence of these technical
 144 factors on reviews process and outcome. Authors report that the length of
 145 the discussion (count comments) and the amount of people involved in the
 146 discussion (Distinct involved Contributors) were judged as influencing factors
 147 by the interviewed contributors.

148 For our study, we select these widely used technical metrics to charac-
 149 terize reviewed patches. Then, we use the propensity score matching (PSM)
 150 technique (see Section 4.3) to ensure that our analysis is not biased by dif-
 151 ferent technical characteristics. PSM [22] is a statistical matching technique
 152 that allows us to create groups of reviews that share similar characteristics.
 153 The technical characteristics considered in our study are:

- 154 • **Count Comments** : The number of comments posted on each code

- 155 review request (*i.e.*, about the proposed code change).
- 156 • Count Patches : The number of patches submitted before the proposed
157 code change is accepted or rejected.
 - 158 • Edited Files (discrete count) : The number of files modified by the
159 proposed code change.
 - 160 • Distinct involved Contributors : The number of developers that par-
161 ticipated in the review of the proposed code change.
 - 162 • Code Churn (Cumulative count) : The number of added and deleted
163 lines that are performed in the reviewed code changes.

164 2.4. Social Network Analysis

165 Social Network Analysis (SNA) is the process of investigating social struc-
166 tures through the use of networks and graph theory [23]. A Network is
167 typically modeled using a graph structure consisting of vertices and edges.
168 Vertices represent individuals or organizations. An edge connecting two ver-
169 tices represents some type of relationships between the two individuals or
170 organizations. Social network analysis focuses on studying social network
171 graphs to understand the patterns of interactions and the relative positions
172 of individuals in a social setting [24]. SNA provides various global or node-
173 specific computed metrics for a network, that are useful for making general
174 statements about specific nodes or classes of nodes. Examples of such metrics
175 are betweenness, diameter, distance, density, betweenness centrality, degree
176 centrality, or eigenvector centrality [23].

177 SNA is being widely used by researchers to model the social structure of
178 OSS communities and barely used in analyzing Open Source Software Peer
179 Review [25]. Previous studies using SNA in OSS generally indicated a few
180 central persons being responsible for most of the interactions in the network
181 (*Core*) and a less connected large group of contributors (*Peripheral*) [26].
182 Through sentiment analysis, we aim to get insights about contributor’s pos-
183 itivity/negativity in relation to their position in code review interactions
184 networks.

185 3. Related Works

186 Several works have focused the attention of the research community on
187 sentiments analysis. These works span many fields ranging from happiness at

workplaces [27] to emotions in social networks' messages such as Yahoo and Twitter [4, 28] and online Q&A such as Stack Overflow posts [29]. Guillory *et al.* [30] went a step further and examined the spread of negative emotions into online communities. Their analyses suggest that contagion of negative emotions can occur in groups of people and impact their performance.

Guzman and Bruegge [31] presented a position paper that describes emotional awareness in software development teams. The paper was motivated by the same concerns that have motivated our approach. Their approach investigates the collective emotional awareness of developers in distributed teams. It extracts emotional state from a 1,000 of collaboration artifacts aiming to summarize emotions expressed in those artifacts by extracting topics and assigning them an average emotion score. Authors presented the emotion average fluctuation to the project leaders, whom confirmed the correlation of positive and negative emotion peaks with team performance, motivation and important deadlines. Our work improves and expands their idea by using propensity score to allow for more accurate comparisons, and apply them on comments related to code reviews instead of comments from commits.

Sinha *et al.* [32] analyzed developers commits logs for a large set of Github projects and found that the majority of the sentiment expressed by developers is neutral. They also found that negative comments are more present than positive ones (respectively 18.05% vs. 7.17%). Similarly, Guzman *et al.* [33] examined the sentiments expressed by developers in comments related to commits from 29 open source projects and found an approximately equal distribution of positive, negative and neutral sentiments. Paul *et al.* [34] explored the difference of expressed sentiments between men and women during various software engineering tasks including the code review practice. The authors report that women are less likely to express their sentiment than men and that sentiment words, emoticons, and expletives vary cross-gender. However, their study did not investigate the effect of expressed sentiment on the productivity of the code review activity according to the duration and results.

Khan *et al.* [35] conducted two studies to explore the impact of sentiments on developer's performance. They found that programmers' moods influence positively some programming tasks such as debugging. Similarly, Ortu *et al.* [36] studied the impact of developers' affectiveness on productivity focusing on the correlation between emotional states and productivity in terms of issues fixing time. They report that the happier developers are, *i.e.*, expressing emotions such as joy and love in their comments, the shorter

Table 1: Existing Sentiment Analysis Tools.

Tool	Purpose	Technique	Trained on	Ref.
Sentistrength	General	Rule-based	Twitter	[28]
Sentistrength_SE	Focused	Rule-based	Jira	[13]
Senti4SD	Focused	Lexical Features	Stack Overflow	[11]
SentiCR	Focused	Lexical Features	Code Reviews	[12]

the issue fixing time is likely to be. They also report that emotions such as sadness are linked to longer issue fixing time. Also, Destefanis *et al.* [37] investigated social aspects among developers working on software projects and explored whether the politeness of comments affected the time required to fix any given issue. Their results showed that the level of politeness in the communication process among developers does have an effect on the time required to fix issues and, more specifically the more polite the developers were, the less time it took to fix an issue. We complement existing work on the impact of sentiment on productivity by studying the influence of text-based expressed sentiment on the duration and outcome of code reviews.

Recent studies have investigated factors affecting the effectiveness of code review comments. Rahman *et al.* [38] extracted a number of features from the text of the review comments attempting to predict the usefulness of code review comments using textual features. However, their empirical study was limited to structural characteristics of the text without considering emotions/sentiments expressed in them. Efstathiou and Spinellis [7] studied the language of code review comments and report that language does matters. In this paper, we continue this line of work by investigating the role of sentiments expressed in code review comments on the outcome of code review. Since Lin *et al.* [8] recently highlighted issues with the accuracy of existing sentiment analysis tools from the literature, we have choose the most powerful sentiment analysis tool based on a benchmarking of several sentiment analysis tools. In Table 1, we present a summary of existing sentiment analysis tools that are designed and tested using data from software artifacts.

4. Empirical Study Design

Our overall goal is to understand the influence of expressed sentiment, throughout comments, on time and outcomes of code reviews. Figure 2

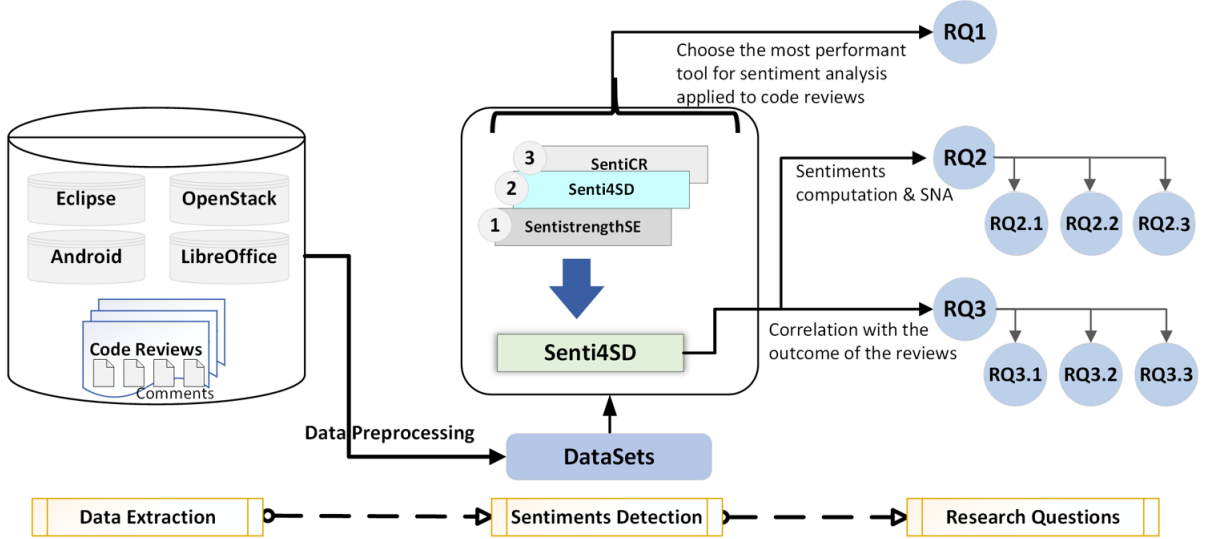


Figure 2: Overview of Our Empirical Study.

253 presents an overview of the steps of our study and how they relate to our
 254 research questions. In the remainder of this section, we describe each step in
 255 details.

256 4.1. Data Collection

257 We conduct our empirical study based on publicly available code re-
 258 view data, mined from Gerrit system and organized in a portable database
 259 dump [39]. We selected this data set because it contains a substantial
 260 volume of data from well-known open source projects organized in a rela-
 261 tional database² as depicted in Figure 3. In our study we used data of four
 262 well-known open-source systems, OpenStack³, Eclipse⁴, Android⁵ and Libre-
 263 Office⁶. OpenStack is a software platform for cloud computing, controlling
 264 large pools of computing, storage, and networking resources throughout a
 265 data center. Eclipse is an integrated development environment (IDE) used
 266 in computer programming. Android is a free software stack for a wide range

² <http://kin-y.github.io/miningReviewRepo/>

³ <http://openstack.org>

⁴ <https://eclipse.org/>

⁵ <https://source.android.com/>

⁶ <https://www.libreoffice.org/>

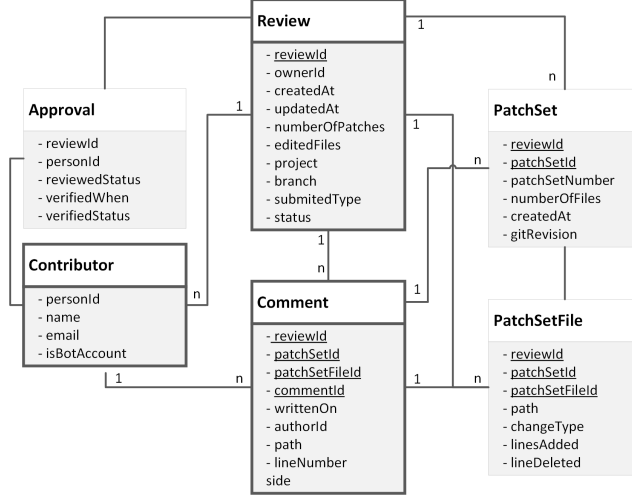


Figure 3: Simplified database schema of Gerrit data.

267 of mobile devices led by Google. LibreOffice is a fork from the OpenOffice.org
 268 project. We selected these projects because they have been actively devel-
 269 oped for more than five years and hence provide a rich data set of reviews.
 270 Also, they are from different domains, are written in different programming
 271 languages, and have been quite studied in other research domains.

272 The original dataset⁷ is stored in a relational database (335,626 reviews
 273 and contains over 5 million comments) under the schema depicted in Figure 3.
 274 In general, a contributor (*i.e.*, personId, name, email) requests a review
 275 characterized by a reviewId, the creation time (createdAt), the last time
 276 modified (updatedAt), related project and the source code branch. A review
 277 includes a set of patches when the author repeatedly update the change by
 278 committing new resubmissions with the same review request ID and a list
 279 of edited files. The history of launched discussion over proposed changes is
 280 recorded in the table ‘*Comment*’.

281 We retrieved and exported required data into separate csv files to ease
 282 our data pre-processing. Table 2 shows descriptive statistics regarding the
 283 studied projects.

⁷<http://kin-y.github.io/miningReviewRepo/>

Table 2: A statistical summary for each studied system

Projects	#Reviews	#Comments	#Contributors
Openstack	228,099	5,021,264	8,088
Eclipse	15,887	153,176	1,082
Android	63,610	355,765	3,334
LibreOffice	28,030	174,181	634

284 4.2. Data Preprocessing

285 In order to improve the quality of our dataset with respect to our main
 286 goal which is studying the human sentiments expressed in code review com-
 287 ments, we performed three pre-processing steps on the raw data:

- 288 1. We discarded comments generated automatically such as those gener-
 289 ated by build automation and continuous integration services. Those
 290 comments contain key-words such as: ‘Jenkins’, ‘Hudson’, ‘Bot’, or
 291 ‘CI Servers’ (about 29% of the total comments were excluded). Using
 292 regular expressions, we excluded automatic expressions (*e.g.*, Build suc-
 293 ceed, Build failed, etc.). In addition, we removed reviews with status
 294 = “New” since their final status remains unknown ($\sim 5\%$ of total re-
 295 views). We limited our analysis to closed reviews (i.e., reviews marked
 296 as “Merged” or “Abandoned”) that contain at least one comment. The
 297 remaining data set contains 4,426,451 comments belonging to 317,373
 298 reviews.
- 299 2. For each review, we gathered key information such as timestamp of
 300 opening and closing of the review, count of edited files, count of patches,
 301 and number of added and deleted lines. Clustering reviews based on
 302 these metrics help us to make unbiased comparisons later on.

303 4.3. Data Analysis

304 We use the propensity score matching (PSM)⁸ [22] method to regroup
 305 reviews homogeneously according to some characteristics (*e.g.*, size of the
 306 review, code churn, number of comments, etc.). Previous works report that
 307 code reviews are affected by a variety of technical factors such as the size

⁸https://en.wikipedia.org/wiki/Propensity_score_matching

308 of the source code [6]. Using PSM allows us to be able to compare reviews
 309 that are logically comparable in terms of these known affecting factors. PSM
 310 is a statistical matching technique widely used to compress covariates to a
 311 variable (*i.e.*, compare technical factors and generate a propensity score).
 312 PSM is proposed to treat the effects of confounding factors. [40] presents an
 313 evaluation of the efficiency of PSM in mitigating confounding factors.
 314 In this work, we used the R package called *Matchit* to carry out the first two
 315 steps enumerated bellow, while step 3 required a manual verification. The
 316 three steps are described as follows:

- 317 1. A logistic regression model is built based on a high-dimensional set
 318 of characteristics. The revision sentiment (Positive or Negative) is set
 319 as the dependent variable, and reviews technical characteristics: (*i*)
 320 amount of comments for the review; (*ii*) count of patchSets, (*iii*) num-
 321 ber of edited files, (*iv*) distinct involved Contributors, and (*v*) churn are
 322 set as its independent variables. The output of the logistic regression
 323 model is a fitted value (a probability value) called propensity score.
- 324 2. The propensity scores are used to match pairs of data points. Each pair
 325 has different values of the dependent variable. Similar values of the
 326 propensity score imply a similarity of reviews technical characteristics.
 327 For our purpose we used the Genetic matching algorithm to match
 328 appropriate pairs of reviews. Then matched pairs are combined into a
 329 new dataset.
- 330 3. The final step is to verify the balance of covariate characteristics. To do
 331 that, we manually compared the means differences for each covariate
 332 variable across matched reviews.

333 The output of PSM is two groups : one for positive reviews and the other
 334 one for negative reviews. Although the final step of PSM is to verify the
 335 balance of covariate characteristics, we carried out a manual validation of
 336 confounding bias on PSM outputs as shown in Table 3. For instance, the
 337 mean difference of total comments for positive and negative reviews shift
 338 respectively from (9.96, 12.6) to (6.2, 6.2), which means more homogeneous
 339 groups. One can also notice greater p -values⁹ with matched reviews; meaning

⁹The p -values are calculated using Mann-Whitney U test [41]

340 an equivalent distribution regarding technical characteristics.

Table 3: Mean differences of technical characteristics before and after Propensity Score Matching (Eclipse Project).

	Befor PSM			After PSM		
	Positive reviews	Negative reviews	<i>p</i> -value	Positive reviews	Negative reviews	<i>p</i> -value
Comments	9.96	12.60	9.07e-06	6.20	6.20	0.09
Patchesets	2.60	2.63	2e-03	1.51	1.51	0.06
Edited_files	13.25	14.22	0.33	1.66	1.66	0.37
Churn	4993.70	8431.20	0.06	51.49	42.29	0.34
Distinct contrib	2.94	3.27	1.60e-7	2.66	2.66	2e-4

341
 342 The new balanced dataset, used later in answering RQ3.1, contains (2,393
 343 positive vs 876 negative reviews) for Openstack, (811 positive vs 155 negative
 344 reviews) for Eclipse, (11,831 positive vs 2,373 negative reviews) for Android,
 345 and (9,002 positive vs 498 negative reviews) for LibreOffice.

346 5. Findings

347 We now present the findings of the sentiment analysis conducted on four
 348 OSS projects. For each of our four research questions, we present our moti-
 349 vation, the approach, and results.

350 RQ1. What is the performance of Sentiment Detectors When Ap- 351 plied on Code Reviews?

352 • **Motivation.** To investigate whether the negative and positive sentiments
 353 expressed in developers’ text-based review interactions affect the code re-
 354 view process, we need a tool capable of detecting sentiments in code review
 355 comments accurately. So far three different sentiment detection engines have
 356 been proposed in the software engineering literature [13, 11, 12], but not
 357 trained specifically on comments of code reviews. However, these previous
 358 attempts to analyze text-based sentiments for software engineering have been
 359 either incomplete nor reusable for other domains [8, 7]. A major reason of
 360 this inadequacy is that software engineering encompasses vocabulary from
 361 diverse sub-domains [7]. Therefore, a tool trained and successfully tested in

one sub-domain (e.g., Q&R Stack Overflow) may not be useful enough for another sub-domain (i.e., comments within Jira issues system). Consequently, we were more cautious on how to choose our tools.

• **Approach.** We compared the performance of three sentiment detection tools: SentistrengthSE [13], Senti4SD [11], and SentiCR [12] aiming to choose the most adequate tool for the domain of source code reviews. These three tools have been trained previously to detect sentiments in software engineering using specific datasets (see Table 1). In order to compare cautiously the performance of the three tools, we carried out a manual annotation (by four raters) on a subset of comments. To strengthen our sampling, we built an over-sampling approach in which the minority class (i.e., negative sentiments) is equally represented. Concretely, following the approach used by Novieli et al. [42], we built four sub-datasets by performing opportunistic sampling. The first sample is created based on the output of SentiStrengthSE, it contains 1,200 comments equally distributed (for each project we have 100 Positive, 100 Negative and 100 Neutral, making $300 \times 4 = 1,200$). The second and the third samples retrieved respectively from Senti4SD and SentiCR contain 360 comments each. The final sample contains 300 random comments. Then each review comment was manually annotated (positive, negative, neutral) by the first author and one of the other authors to ensure a stable annotation. The same approach was previously used by Lin et al. [8] to produce their sentiment benchmark. The agreement between the two coders, measured using Cohen’s kappa, ranged from 81% to 95% (83% for Senti4SD sample, 95% for SentiCR sample, 81% for SentiStrengthSE and 91% for Random sample). To resolve the disagreements between raters, the annotations were discussed and the guideline of annotation was updated by the first author. For instance, an example of a disagreement between two annotators happened for the following sentence: *”Patch Set 2: Fails Merges in public tree, but does not build. Please fix and reupload. Thanks!”*. The first annotator classified this comment as Positive, while the second one classified it as Negative. Through mutual consent we decided to tag this typical comment as Positive since the commenter was very polite and used *”Please”* and *”Thanks”* in his text.

Our sub-datasets as well as the original dataset (5 millions comments) are available in the companion on line appendix [43] for the purpose of replication.

• **Results.** Table 4 reports the performance obtained in terms of recall,

precision, and F1-measure, for each polarity classes (Positive, Negative, and Neutral) as well as the overall performance, for the three tools when applied to our data samples. We highlight the best values for each metric. Surprisingly, when expecting a good performance from SentiCR tool, Senti4SD shows a slightly better overall performance than the other tools ($F1 = 0.79^{10}$). Again, Lin et al [8] pointed out that sentiment analysis tools should always be carefully evaluated in the specific context of usage. We double check our results by performing McNemar¹¹ statistical test [45] in order to compare the classification results of the three tools. The performance differences between Senti4SD and other classifiers were found to be statistically significant ($p_value < 0.05$ and $z\ scores = 11.49 > 0$) indicating that Senti4SD performs better than SentiCR. Moreover, when comparing this result with our manual tagging, we noticed that 472 comments were correctly classified by Senti4SD and misclassified by SentiCR while only 178 comments correctly classified by SentiCR and misclassified by Senti4SD.

415

RQ1 What is the performance of Sentiment Detectors When Applied on Code Reviews?

On average Senti4SD led to the best performance (Precision **79%**, F1 **79%**) when applied to our code review data samples.

416

417 RQ2. How Prevalent are Sentiments in Code Reviews?

418 • **Motivation.** Sentiments are ubiquitous in human activity: There is an
 419 old saying “*Feeling Good-Doing Good*” [46]. OSS contributors may under-
 420 perform if they do not feel safe and happy [35]. Negative emotions like anger
 421 can make people less motivated and thus less creative [36]; two key fac-
 422 tors to ensure productivity within modern software organizations [27]. For
 423 instance, Linus Torvalds sent out an email¹² to the Linux developers’ com-
 424 munity admitting his verbal abuse in communications [“*My flippant attacks*”

¹⁰We used the F1-measure to determine the best performing classifiers, following standard practices in Information Retrieval [44]

¹¹<https://stat.ethz.ch/R-manual/R-devel/library/stats/html/mcnemar.test.html>

¹²<https://gizmodo.com/linux-founder-takes-some-time-off-to-learn-how-to-stop-1829105667>

Table 4: Performance of Sentiment Detectors in Code Review samples (P = Precision, R = Recall, F1 = F1-Measure)

Dataset	Class	Sentistrength_SE			Senti4SD			SentiCR		
		P	R	F1	P	R	F1	P	R	F1
SentistrengthSE_based	Positive	0.86	0.85	0.83	0.81	0.84	0.82	0.59	0.89	0.71
	Negative	0.91	0.61	0.73	0.66	0.68	0.67	0.59	0.66	0.62
	Neutral	0.7	0.98	0.81	0.83	0.81	0.82	0.88	0.69	0.77
	Micro-avg.	0.8	0.8	0.8	0.79	0.79	0.79	0.72	0.72	0.72
	Macro-avg.	0.82	0.8	0.79	0.77	0.77	0.77	0.69	0.75	0.7
Senti4SD_based	Positive	0.83	0.92	0.87	0.91	0.91	0.91	0.3	0.81	0.43
	Negative	0.57	0.8	0.67	1	0.78	0.87	0.42	0.8	0.55
	Neutral	0.89	0.7	0.78	0.79	0.96	0.87	0.93	0.51	0.66
	Micro-avg.	0.78	0.78	0.78	0.88	0.88	0.88	0.59	0.59	0.59
	Macro-avg.	0.76	0.81	0.77	0.9	0.88	0.88	0.55	0.71	0.55
SentiCR_based	Positive	0.6	0.82	0.7	0.74	0.72	0.73	0.73	0.7	0.72
	Negative	0.52	0.4	0.45	0.67	0.46	0.55	0.91	0.25	0.4
	Neutral	0.82	0.76	0.79	0.76	0.83	0.79	0.53	0.93	0.67
	Micro-avg.	0.73	0.73	0.73	0.75	0.75	0.75	0.63	0.63	0.63
	Macro-avg.	0.65	0.66	0.64	0.72	0.67	0.69	0.72	0.63	0.6
Random	Positive	0.27	0.95	0.42	0.83	0.89	0.86	0.05	0.44	0.09
	Negative	0.11	0.15	0.13	0.58	0.76	0.66	0.05	0.14	0.08
	Neutral	0.95	0.75	0.84	0.97	0.93	0.95	0.96	0.71	0.81
	Micro-avg.	0.74	0.74	0.74	0.92	0.92	0.92	0.69	0.69	0.69
	Macro-avg.	0.44	0.62	0.46	0.8	0.86	0.82	0.35	0.43	0.33
Overall Micro-avg.		0.76	0.76	0.76	0.83	0.83	0.83	0.65	0.65	0.65
Overall Macro-avg.		0.66	0.72	0.66	0.79	0.79	0.79	0.57	0.63	0.54

in emails have been both unprofessional and uncalled for,”]. Torvalds stepped down because people were complaining about his lack of care sentiments in his communications which has hurt some contributors and may have driven some away from working in kernel development altogether [“I’m going to take time off and get some assistance on how to understand people’s emotions and respond appropriately”]. Empirical evidence of the effect of expressed sentiments contained into comments on code reviews could help developers pay more attention to the way they comment on other’s work, especially in a context of virtual communities such as Github characterized by multicultural contributors.

We are also interested in understanding how the expressed sentiments of contributors evolve over time as they gain in seniority within a project. There has been research examining OSS contributors’ involvement over time [26], in particular, researchers pointed out that empirical analyses that mix the two groups will likely yield invalid results. Surprisingly little research has examined the evolution of text-based sentiments when contributors gain reputation

441 (*i.e.*, belong to the core team leading the project). Reviewer’s sentiment may
442 wax and wane as project progresses. We thus derive the following research
443 question.

444

- 445 • **RQ2.1: How are positive and negative sentiments expressed**
446 **in code reviews?**

447 Previous research [2, 33, 36] that observed significant presence of senti-
448 ments and emotions in code reviews and issue comments . Therefore, before
449 analyzing the relationship between sentiments expressed in code reviews and
450 code review outcomes, it is important to learn whether sentiments are also
451 prevalent in our dataset of code review comments. Given that developers
452 can express as well as seek opinions in diverse development scenarios [47, 48],
453 their expression of opinions in code review comments may be influenced by
454 diverse development needs and situations. Therefore, it is necessary to learn
455 how developers expressed those sentiments and what could have triggered
456 the developers to express those opinions.

- 457 • **RQ2.2: How do the prevalence of Expressed Sentiments of**
458 **Reviewers Evolve Over Time?**

459 We are interested in analyzing potential differences in expressed senti-
460 ments between core and peripheral contributors. Core members are those
461 developers that contribute intensively and sustainably to the OSS project,
462 and thus, lead the community, while peripheral ones are occasional contribu-
463 tors with less frequent commits. Our main purpose is to study the correlation
464 between a gain of contributors reputation and the nature of sentiments they
465 express within reviews comments. We hypothesize that newcomers try to
466 imitate contributors with a certain reputation, which might affect the cul-
467 ture of commenting. Hence, we formulate the following research questions:

468

- 469 • **RQ2.3: Do Core and Peripheral Contributors Express Dif-**
470 **ferent types of Sentiment According to their Position in a**
471 **Collaborative Social Network Graph?**

472 *RQ2.1: How are positive and negative sentiments expressed in code reviews?*

- 473 • **Approach:** Our dataset contains more than 4.4 million comments on code
474 reviews regarding four long-lived and well known OSS projects: *Openstack*,

475 *Eclipse*, *Android*, and *LibreOffice*. The distribution of comments is shown in
476 Table 5. Next, we conducted a sentiment analysis on comments using natural
477 language processing techniques. We used *Senti4SD* tool, a fully-automated
478 algorithm, to compute the sentiment score for each comment on each review.
479 To determine whether sentiment scores are consistent in the projects, we cal-
480 culate the skewness and kurtosis of the sentiment scores. The skewness of
481 a distribution captures the level of symmetry in terms of mean and median.
482 For instance, a negative skew means that the overall reviews are towards
483 negativity, while a positive skew means that the reviewers overall express
484 more positivity. Kurtosis explains the shape of the distribution (univariate
485 normal distribution is 3). A kurtosis lower than 3 means that the reviewers
486 have a strong consensus, while a kurtosis greater than 3 means a divergence.
487 In order to investigate further the type of sentiments expressed within code
488 review comments, we manually tagged positive (666) and negative (443) com-
489 ments from the dataset used to answer *RQ1*. To do so, we leveraged on
490 categorization provided by Tourani et al. [49], which categorized positive
491 sentiments into six categories and negative ones into four categories as de-
492 scribed in Table 6. Each comment was manually categorized by two raters,
493 the agreement between the two coders, measured using Cohen’s kappa, was
494 61% for positive comments and 65% for negative comments.

495 • **Results: 8.31% of comments were reported as positive** (score =
496 1) in the Eclipse project. **While 89.92% of comments were neutral**
497 (score = 0) and **1.77% were negative sentiments** (score = -1). Table 5
498 summarizes the results of sentiment computation for the studied projects and
499 provide some descriptive statistics, *i.e.*, the mean, standard deviation, kurto-
500 sis, and skewness. We noticed relatively similar distributions concerning the
501 proportion of expressed sentiment within comments across the four projects.
502 Neutral comments are the most present (**83.81%**) which confirms the re-
503 sults of previous studies [32]. The large amount of neutral sentiments can be
504 mainly explained by the presence of technical vocabulary within comments.
505 For instance, in Eclipse project, over 153 thousand comments, 89.92% was
506 reported as neutral.

507 We found that **13.94%** of comments related to all projects were positive
508 (*e.g.*, “Thanks for the most excellent review. :)”), while around **2.24%** of
509 comments were identified as negative (*e.g.*, “Horrible :”).

510 Eclipse is the only project with a Kurtosis value greater than 3 which
511 suggests that sentiment are diverse among the contributors while Openstack,
512 Android, and LibreOffice have a Kurtosis slightly less than 3, meaning that

Table 5: Distributions of Sentiments in Reviews

Project	Positive	Neutral	Negative	Mean	SD	Kurtosis	Skewness
Openstack	16.27 %	81.04 %	2.68 %	0.13	0.41	1.63	0.89
Eclipse	8.31 %	89.92 %	1.77 %	0.06	0.31	6.25	1.53
Android	14.06 %	84.09 %	1.86 %	0.12	0.37	2.43	1.22
LibreOffice	17.14 %	80.21 %	2.65 %	0.14	0.42	1.39	0.87

513 reviewers have a strong consensus on sentiment expression in source code
514 reviews.

515 Overall, results reveal that the distribution is highly positively skewed for
516 Eclipse and Android while moderately skewed for Openstack and LibreOffice.

517 Manual annotation revealed that ‘Friendly Interaction’ is the most preva-
518 lent category of positive sentiments with a proportion of 37.1% as reported
519 in Table 6. This means that, to a large extent, 37.1% of positive interactions
520 between the community’ members are guided with respect and positive at-
521 titudes. ‘Satisfactory Opinion’ and ‘Announcement’ count respectively for
522 19.6% and 15.5%. While the most common category of negative sentiment
523 is ‘Uncomfortable Situation’ with a percentage of 62.30%. This means that
524 62.30% of negative comments express strong pressures such as time con-
525 straints that could overwhelm them, confusion about inexplicable behavior
526 of the software system, or concerns about risks and fears. ‘Unsatisfied Opin-
527 ion’, ‘Aggression’, and ‘Sadness’ count respectively for 16.03%, 15.35% and
528 6.32%.

529 As stated by [49], well-mannered interactions with a positive undertone
530 might lead to a higher productivity. Our RQ3 will investigate the impact of
531 expressed sentiments into comments on the duration and the outcome of a
532 source code reviews.

Table 6: Categorization of sentiments within code review comments.

Sentiment	Category	Example	Total	ratio
Positive Sentiment	Satisfactory Opinion	Thank you Andrey! I really appreciate that you took the time to check that, and I'm glad to hear that performance is now okay :)	130	19.6%
	Friendly Interaction	Works fine no issues	247	37.1%
	Explicit Signals	Restored I'll revive this, it makes the debug info analysis much more pleasant.	82	12.3%
	Announcement	I'm sure you see how having all your patches in one chain helps for sanity	103	15.5%
	Socializing	My pleasure :).	59	8.9%
Negative Sentiment	Curiosity	Before the change, the tests were working fine on the command line.	44	6.6%
	Unsatisfied Opinion	Forgot to publish these. Sorry!	71	16.03%
	Aggression	PS, I hate this change and this API.	68	15.35%
	Uncomfortable Situation	I'm sorry but your approach looks like overkill to me.	276	62.30%
Neutral Sentiment	Sadness	I really dislike this patch and "I would prefer that you didn't submit this" but I don't know if it's a valid reason to -1 it.	28	6.32%
	-	-	1111	

RQ2.1 How are positive and negative sentiments expressed in code reviews?

Open source software developers do express sentiments when they are reviewing each other source code. A percentage of 13.94% of comments related to all projects were positive, while around 2.24% of comments were identified as negative. Also, 'Friendly Interaction' is the most prevalent category of positive sentiments with a proportion of 37.1%, 'Uncomfortable Situation' is the most common category of negative sentiment with a percentage of 62.30%.

533

534 RQ2.2: How do the prevalence of Expressed Sentiments of Reviewers Evolve 535 Over Time?

536 • **Approach:** To answer this research question, we proceeded as follows.
537 First, we examined the sentiment evolution of top 5% contributors for each
538 project during the complete time period under study. We pick the top 5% to
539 ensure that we have the most active contributors without any discontinuity
540 in the review activity. In total, over the four studied projects, we analyzed
541 the evolution of expressed sentiments of the top 5% (484 out of 9680 contrib-
542 utors) who have created 1,493,224 comments (33.73% of total comments).
543 After zooming on this group of contributors, we explored manually in details
544 the time series of the top five contributors for one project, which produced
545 7,184 comments (0.16% of the total comments) to ground sentiment evolu-
546 tion patterns. We focused only on 5 members because of the high cost of the
547 analysis.

548 • **Results:** We observed a trend toward neutral sentiments correlated with
 549 the progression of contributors toward the core team. **The more a con-**
 550 **tributor gains reputation, the more he is likely to express neutral**
 551 **sentiments.** Figure 4 shows the average of sentiment evolution per month
 552 of top 5% contributors (*i.e.*, reviewers). However, we cannot conjecture that
 553 this trend towards neutral sentiments is due to a gain of reputation by con-
 554 tributors. It could also be simply due to cultural changes in the studied
 555 projects. Further analysis are necessary to better understand the evolution
 556 of developers' sentiments in OSS projects.

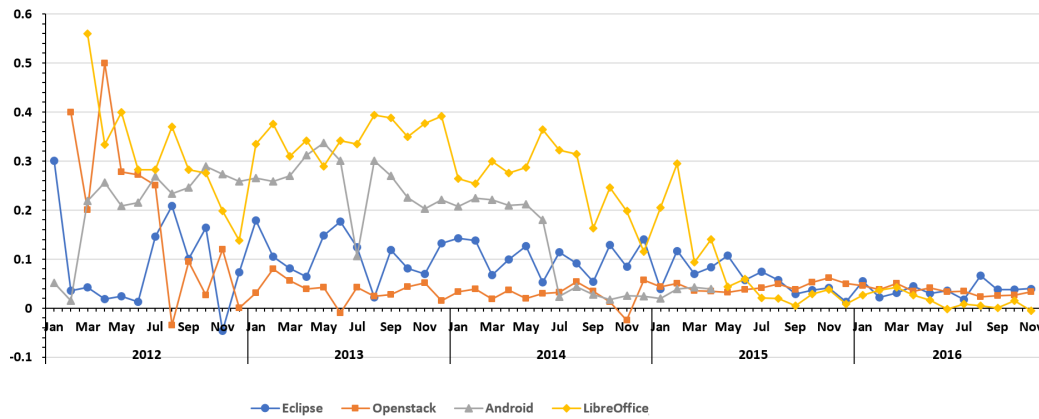


Figure 4: Average Sentiment Evolution per Month of the Top 5% Contributors.

557 In order to get more insights on sentiment average evolution, we monitor
 558 the top 5 core contributors for the Eclipse project. We choose the Eclipse
 559 project because it is intensively studied in the literature. As one can see in
 560 Figure 5, the sentiment averages vary significantly over years, decreasing from
 561 positive towards neutral. The sentiments of the top 5 reviewers in Eclipse
 562 decreased to neutral over time. As mentioned earlier, an interesting future
 563 qualitative research would be surveying the behavior of the most productive
 564 contributors.

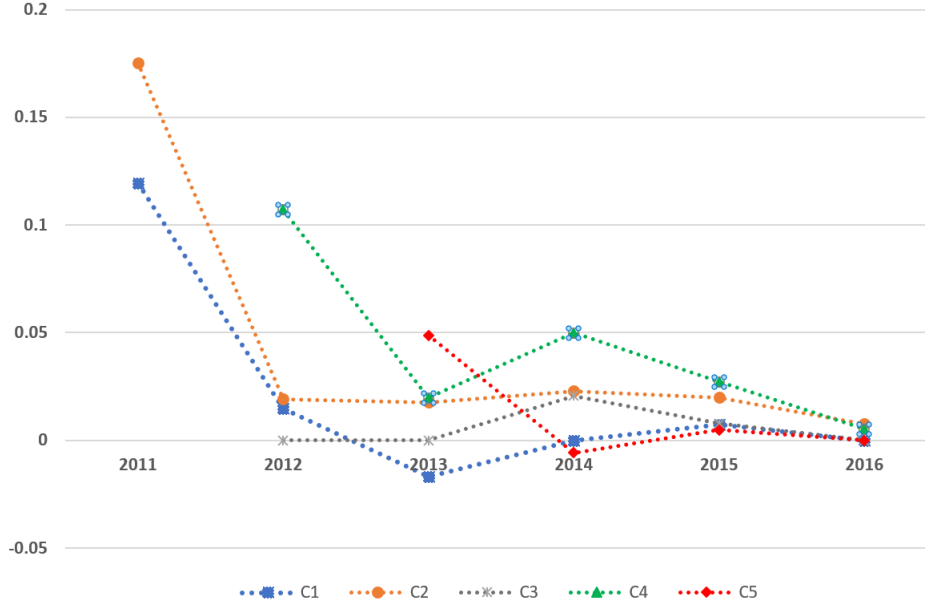


Figure 5: Evolution of the sentiment average (per year) for the top 5 core contributors of the Eclipse project.

RQ2.2 How do the prevalence of Expressed Sentiments of Reviewers Evolve Over Time?

Sentiments expressed by code review contributors tend to be neutral as they progress from the status of newcomer in an OSS project to the status of core team contributor.

566

567 *RQ2.3: Do Core and Peripheral Contributors Express Different Types of Sen-*
 568 *timent According to their Position in a Collaborative Social Network Graph?*

569 • **Approach.** To answer **RQ2** and the sub research questions we need to
 570 build Social Networks for each project in order to detect Core and Periph-
 571 eral contributors. To do so, we calculated the number of interactions between
 572 each pair of developers in each project. Then, we generated our social net-
 573 work graphs as undirected, weighted graphs where nodes represent developers
 574 and edge weights represent the amount of co-edited files by those contribu-
 575 tors. Finally, to locate core and peripheral members we followed the same
 576 approach described in [50]. We used the Kmeans clustering method based

on SNA centrality measures. Centrality measures used for this approach are : Degree centrality, Betweenness centrality, Closeness centrality, Eigenvector centrality, Eccentricity and PageRank. Each metric calculates centrality in a different way and has a different interpretation of a central node [24].

Concretely, we used the Python package *NetworkX* to calculate the six centrality measures for each node in each graph. Then, we used the R implementation of the Kmeans clustering algorithm to partition the nodes into core and peripheral groups based on the six centrality scores. K-means groups the project contributors into two mutually exclusive clusters in which each contributor belongs to the cluster with the nearest mean (measured using different centrality measures). K-means treats each contributor as an object having a location in space. It finds a partition in which objects within each cluster are as close to each other as possible and as far from objects in other clusters as possible. We used the `kmeans()`¹³ function within R with default configuration options to identify core and peripheral contributors. Table 7 provides a description of the core-periphery partitions obtained for the four projects in this study, alongside with the goodness which is the `between_SS / total_SS` values provided as a result when using `kmeans()` of the classification.

Table 7: Core-Periphery distribution in studied Projects

Projects	Size of the Core	Size of the Periphery	Goodness (%)
Openstack	1,081	4,921	82.6
Eclipse	121	628	82.3
Android	189	2,295	84.2
LibreOffice	45	437	85.3

Figure 6 shows the social network structure of the Openstack project as generated by Cytoscape¹⁴, a tool for networks’ visualization. Core developers (shown in yellow) represents a small set of contributors (between 4.55% and 12.14%, for the studied projects) who have generally been involved with the OSS project for a relatively long time and are making significant contributions to guide the development and evolution of the project. Peripheral developers (shown in red) are a larger set of contributors whom occasionally contribute to the project, mostly interacting with core developers, and rarely interacting with each other. To enhance readability of OpenStack graph, we

¹³<https://stat.ethz.ch/R-manual/R-devel/library/stats/html/kmeans.html>

¹⁴<http://www.cytoscape.org/>

604 removed the low-weight degrees (weight $\prec 5$) and isolated nodes.

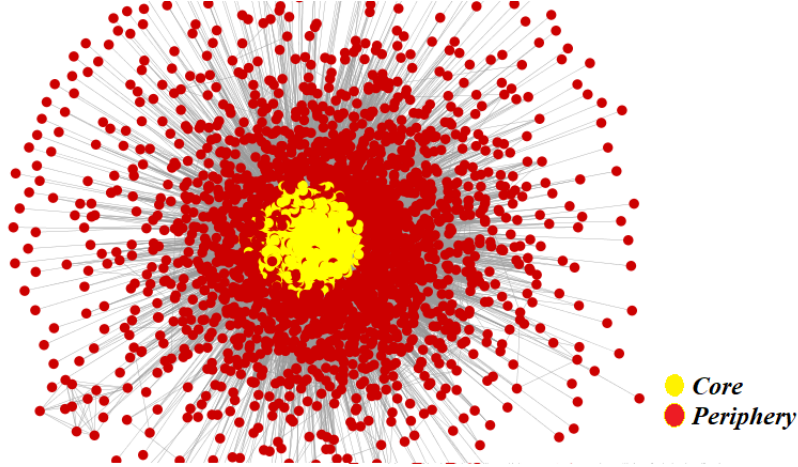


Figure 6: Code Review Social Network Diagram of Eclipse

605 After segregating core and peripheral contributors, a sentiment score av-
 606 erage for each contributor has been calculated based on the sentiment score of
 607 all the comments he made. Next, using Mann-Whitney U test [41], we com-
 608 pared the distributions of sentiment averages between groups of core and
 609 peripheral contributors. The test is applied following the commonly used
 610 confidence level of 95% (*i.e.*, $\alpha \prec 0.05$). Since we performed more than one
 611 comparison on the same dataset, to mitigate the risks of obtaining false po-
 612 sitive results, we use Bonferroni correction [41] to control the familywise error
 613 rate. Concretely, we calculated the adjusted p -value, which is multiplied by
 614 the number of comparisons. Whenever we obtained statistically significant
 615 differences between groups, we computed the Cliff's Delta effect size [41] to
 616 measure the magnitude of this difference.

617 • **Result.** Figure 7 shows the comparison of averages of sentiments between
 618 core and peripheral contributors. For the four studied projects, the distribu-
 619 tion of sentiment averages ranges between $[-1, 1]$. The Mann-Whitney test
 620 revealed a significant difference in the distribution of sentiment average of
 621 core and peripheral contributors. However, the effect size is small, except for
 622 the LibreOffice project where it is medium, as reported in Table 8.

623
 624 Surprisingly, we observed that the peripheral contributors in all four
 625 projects have clearly more outliers - *i.e.*, both positive and negative com-

Table 8: Mann-Whitney U test results

Project	U	p-value	Effect Size
Openstack	4,115,100	2.2e-16	small (0.26)
Eclipse	47,683	1.0e-06	small (0.26)
Android	536,920	2.2e-16	small (0.27)
LibreOffice	22,700	8.48e-12	medium (0.47)

pared to core ones whom sentiments remain concentrated around Neutral emotions (*i.e.*, value equal to zero). We hypothesize that the outliers segment are people participating by a single or a small amount of comments, which impacts the values of averages, whereas Core developers remain neutral while they comment on the source code revisions.

RQ2.3 Do Core and Peripheral Contributors Express Different Types of Sentiments According to their position in the review network?

Open source contributors do express different sentiments depending on the position within the peer review collaborative social network. Peripheral contributors in the four projects clearly have more outliers in expressing positive and negative sentiments, while Core developers remain neutral when commenting on source code revisions.

631

RQ3. How do the presence of sentiments in code reviews correlate with the outcome of the reviews?

Motivation. Code review is an essential practice to ensure the long-term quality of the code base. This modern practice could be influenced by expressed sentiment within contributors' comments. Intuitively, positive sentiments may improve the contributors mood, while negative ones may prove detrimental to their morale. Such change in morale can then impact both the time taken and the outcome of the review process. In particular, it is important to know how expressed sentiments can impact code review practices along the following two dimensions: (1) Code Review Time, and (2) Code Review Outcome. The duration of a source code review is an important factor for a software organization productivity [51]. We pose the following question:

- **RQ3.1 How do the sentiments expressed in the reviews correlate with the duration and the outcome of a review compared**

645
646

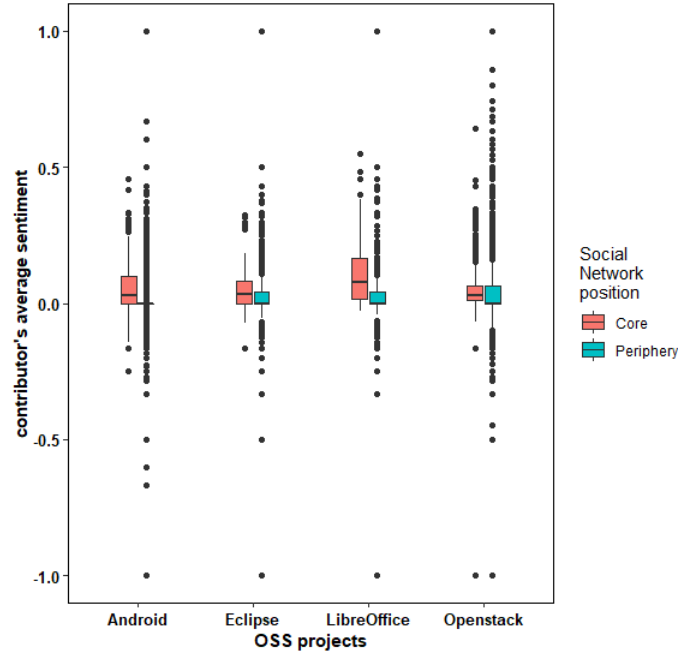


Figure 7: Comparing Sentiments Between Core and Peripheral Contributors.

to the reviews with no sentiments?

When a code review takes much longer than expected, the release of the software and the team productivity can suffer. A number of factors can contribute to such longer time, such as the absence or leave of the core developer responsible for the particular module related to the review or the change in priority. Another mitigating factor could be the presence of *controversy* in the review comments. Intuitively, a feature may be controversial if its code review attracts positive and negative comments almost equally. An empirical understanding of the extent to which such controversies can impact the code review outcome can offer a gain of awareness regarding positive/negative impact of code review practices. As a practical implication, we can motivate a new feature within Gerrit to proactively warn contributors involved in a review team about a risk of delaying the review due to controversies. We thus derive the following research question.

- **RQ3.2. Does the presence of controversies in the code reviews**

662 offers valuable insights into the outcome of those code reviews
663 compared to the reviews with non-controversial comments?

664 We are investigating outlier reviews (that took a very long time) in or-
665 der to determine whether the presence of controversial sentiments in their
666 comments could be the root cause of increased review time.

667 • **RQ3.3. Do sentiments expressed by core contributors impact**
668 **the review outcome differently than those expressed by pe-**
669 **ripheral contributors?**

670 In RQ2, we observed that core and peripheral contributors express dif-
671 ferent types of sentiments. Since the contribution of core and peripheral
672 contributors in reviews activities are likely different, *i.e.*, core contributors
673 are expected to be involved more closely than peripheral contributors in
674 code reviews. We are interested in examining whether sentiments expressed
675 by these two groups of contributors also affect the review process differently.
676 In the following we answer three sub questions.

677 *RQ3.1 How do the sentiments expressed in the reviews correlate with the*
678 *duration and the outcome of a review compared to the reviews with no senti-*
679 *ments?*

680 • **Approach.** An overview of review time distribution in studied projects,
681 pointed out that the slowest review took hundreds of days whereas the median
682 review was less than one day. To avoid bias due to skewed distributions, we
683 used Tukey’s outliers detection methods [52]. A review time is considered as
684 outlier if it is above an *Upper limit*. Tukey’s define this limit based on the
685 Lower and Upper quartiles [Q1, Q3] (*i.e.*, respectively the 25th and the 75th
686 percentiles of data distribution) such as:

$$Upper\ limit = Q3 + 1.5 * IQR \quad (1)$$

687 Where inter-quartile range (IQR) is the interval between Q1 and Q3.
688 Tukey’s method applied on review time distribution detected a distinct *Upper*
689 *limit* days for each project (13.86 for Eclipse, 10.02 for Android, 11.04 for
690 Openstack and 6.52 for LibreOffice). The new dataset contains a total of
691 114,546 reviews.

692 To assess the influence of positive or negative sentiments on the dura-
693 tion of a code review, we used the **Propensity Score Matching** (PSM)

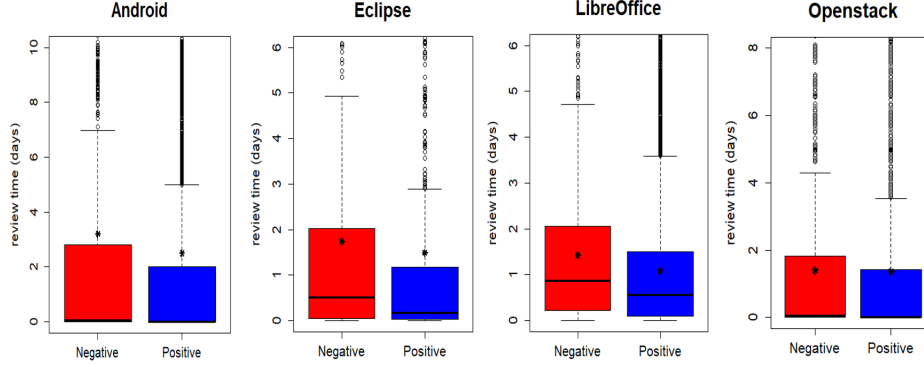


Figure 8: Box-plot of the Duration of Reviews (* : mean).

method [22], as described in Section 4.3. For practical applications, we compare only the reviews that are logically comparable in terms of technical characteristics: (1) *amount of comments for the review*; (2) *count of patch-Sets*, (3) *number of edited files*, (4) *Distinct involved Contributors*, (5) and *churn* (*i.e.*, sum of inserted and deleted lines of code to measure how large the change is). Also, to assess the influence of sentiments on the reviews' outcome, we mapped the sentiment summary of each review (Positive or Negative) with its final status (Merged or Abandoned).

• **Results.** Comparing a homogeneous group of reviews (obtained through PSM) reveals that positive reviews **took less time** to be closed than negative ones, as depicted in Figure 8. **Negative reviews required a supplementary time of 1.32 day on average to be closed than positive ones.** In other words, the average of durations for positive reviews is less than the average for negative reviews.

Also, as shown in Figure 8, positive reviews not only have the minimum median review time, but also, they have the lowest maximum number of days needed to be closed, compared to negative reviews. For instance, in the Eclipse project, positive reviews last a maximum of 2.89 days, while reviews containing negative comments took approximately 5 days of review. Also, the Mann-Whitney test revealed a significant difference in the distribution of reviews fixing times between positives and negatives reviews with a small

Table 9: The p-value and the effect size of review times in positive vs negative reviews

Project	p-value	Effect Size
Openstack	1.6e-4	small (0.02)
Eclipse	0.2e-4	small (0.09)
LibreOffice	2.8e-11	small (0.17)
Android	1.6e-4	small (0.08)

716 effect size¹⁵) for all studied projects (see Table 9).

717 Figure 9 shows mapping results of reviews types (Positive or Negative)
 718 with the final status of the review (Merged or Abandoned). For each project,
 719 the ratio values presents the distribution percentage of positive and negative
 720 reviews within merged review (first bar) and abandoned ones (second bar).
 721 Results show that, not only does the sentiment expressed by developers affect
 722 the duration of code review, but it also affects the outcome. For instance,
 723 in Eclipse project, over 93% of successfully merged reviews were tagged as
 724 positive, while 55% out of all abandoned reviews have negative sentiments
 725 into their comments.

RQ3.1 How do the sentiments expressed in the reviews correlate with the duration and the outcome of a review compared to the reviews with no sentiments?

The presence of positive sentiments in comments related to source code reviews seems to contribute to reducing the review time by an average of 0.4 day. It also seems to affect code reviews outcomes.

726

727 *RQ3.2. Does the presence of controversies in the code reviews offers valuable*
 728 *insights into the outcome of those code reviews compared to the reviews with*
 729 *non-controversial comments?*

730 • **Approach.** In the previews questions, we analyzed only reviews that took
 731 less than the Upper limit before being accepted or abandoned. In this re-
 732 search question, we examine reviews that took a very long time; *i.e.*, more
 733 than identified threshold. Our goal is to investigate whether the presence of

¹⁵p-value and effect size are measured using Mann-Whitney U test and the Cliff's Delta effect size as explained in RQ2.2

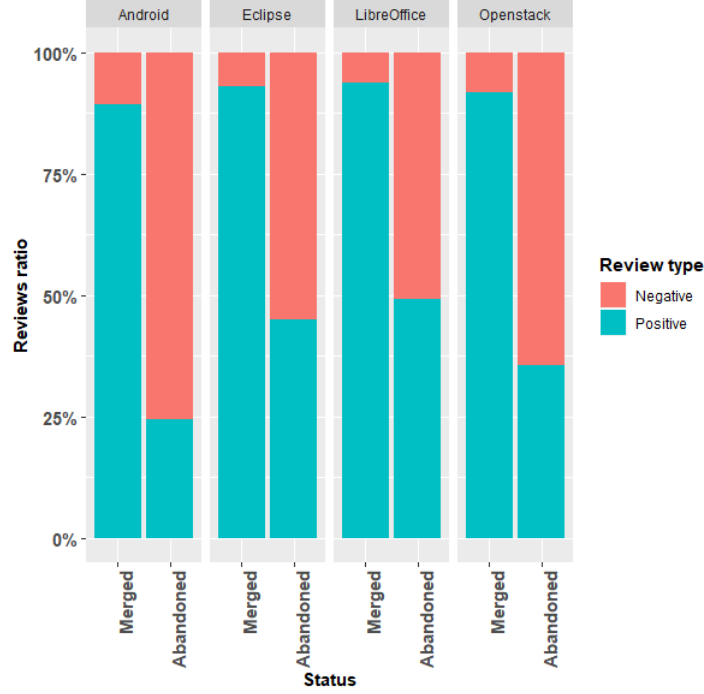


Figure 9: Ratio of Positive vs. Negative Reviews Regarding Reviews' Outcome.

734 controversy in reviews discussion is the root cause of the long delays. The
735 Merriam-Webster Dictionary defines controversy as a “strong disagreement
736 about something among a large group of people”. In our context, we classify
737 a review as *controversial* if the discussions about the submitted source code
738 contain controversial comments. We compute the degree of controversy us-
739 ing controversialMix [53], which is a score that estimates how many mixed
740 positive and negative comments are in a review discussion.

$$\text{ControversialMix} = \frac{(\text{Min}(|Pos|, |Neg|))}{(\text{Max}(|Pos|, |Neg|))} \frac{(|Pos| + |Neg|)}{(|Neu| + |Pos| + |Neg|)} \quad (2)$$

741 Where Pos, Neg and Neu are the sets of comments with positive, negative
742 and neutral polarity.

743 ControversialMix takes in consideration the amount of positive, negative
744 and neutral comments in order to capture the diversity on expressed senti-
745 ments within the same review. Before computing controversialMix we did

Table 10: Distribution of Controversial and Non Controversial comments within outliers reviews (Yes =controversial, No= Not controversial)

Project	Controversial	#Reviews	Avg_review_time(days)	P_value
Android	No	780	108.48	0.37
	Yes	31	120.38	
Eclipse	No	185	126.95	0.7
	Yes	4	132.46	
LibreOffice	No	442	44.18	0.01
	Yes	34	50.51	
Openstack	No	10,201	61.82	0.17
	Yes	95	67.30	

some data preprocessing by discarding : reviews with only one comment; reviews where all comments have the same tag (negative, positive, neutral) and reviews threads that have only positive or negative comments. Finally, a review is tagged as controversial if $\text{controversialMix} \geq 0.5$.

• **Results.** Table 10 shows the distribution of controversial reviews and the average review time needed to fix controversial and non controversial reviews.

Wilcoxon test applied on controversial and non-controversial reviews reveals that results about average review time (days) were significant for only one project: LibreOffice with a $p\text{-value} = 0.01$. For this particular project, one can see that controversial reviews required in average more days to be closed (+6.33 days). However, the limited amount of controversy identified for Eclipse, Android, Openstack respectively (4, 31, 95) compared to the amount of controversial reviews found in LibreOffice as shown in Table 10, could explain the non significant result obtained from the wilcoxon test on these projects. Consequently, we were not able to confirm this finding due to the lack of data.

RQ3.2. Does the presence of controversies in the code reviews offers valuable insights into the outcome of those code reviews compared to the reviews with non-controversial comments?

Controversy significantly increased the time taken to review code in the LibreOffice project (44.18 to 50.51). Unfortunately, we did not have enough controversial reviews in the other projects to confirm our finding.

762

763 *RQ3.3 How do the sentiments expressed by the core vs the peripheral contrib-*
764 *utors correlate with the outcomes of the code reviews?*

765 • **Approach.** We create two buckets for each project, one for each type
766 of contributors (*i.e.*, core and peripheral). For each class of contributors,
767 we create three polarity buckets, labeled as positive, negative, and neutral.
768 For instance, the positive bucket contains all the review times (in days) of
769 positive reviews. The negative bucket contains all the review times (in days)
770 of negative reviews. The neutral bucket contains all the review times (in days)
771 of neutral reviews. In each of these buckets we excluded: (1) reviews that
772 took less than one day, (2) reviews that took more than the thresholds for
773 each project that we determined using Tukey’s outliers detection algorithm
774 (see RQ3.1). Intuitively, from a productivity perspective, it is useless to
775 analyze the impact of sentiments for a review that took less than a day
776 (because it is already an impressive performance).

777 We then divide each bucket into two further buckets: (1) **Mixed.** We
778 put a review in this bucket if it has sentiments expressed by both core and
779 peripheral contributors. (2) **Exclusive.** We put a review in this bucket, if
780 it has sentiments expressed by either the core or the peripheral contributors,
781 but not by both in the same review. We compare the opinion impact of core
782 versus peripheral contributors for the reviews in the ‘Exclusive’ bucket for
783 each project.

784 • **Results.** In Table 11, we show the summary statistics of the review
785 time (in days) taken when the core or peripheral contributors offered posi-
786 tive or negative reviews. The ‘Neutral’ column under each contributor type
787 shows the time taken when the contributor offered neutral comments. For
788 all projects, the average review time increased when the peripheral contribu-
789 tors provided negative comments. For all project, the average review time is
790 larger when the contributors provided negative comments than when the core

Table 11: The impact of sentiments expressed by the core vs peripheral contributors on the code review elapsed time (in days)

Project	Reviewer Type	Review Time for Overall Sentiment Type			
		Time Metric	Positive	Negative	Neutral
Eclipse	Core	Average	4.8	5.1	4.8
		Std	3.3	3.2	3.3
		Median	4.0	4.9	3.9
	Peripheral	Average	4.6	4.8	4.8
		Std	3.2	3.4	3.3
		Median	3.7	3.9	3.7
Android	Core	Average	4.1	4.3	4.0
		Std	2.4	2.4	2.4
		Median	3.8	4.1	3.3
	Peripheral	Average	4.1	4.9	4.0
		Std	2.5	2.3	2.4
		Median	3.2	3.5	3.5
Libreoffice	Core	Average	3.1	3.3	3.1
		Std	1.6	1.7	1.6
		Median	2.9	3.0	2.8
	Peripheral	Average	3.1	4.1	3.2
		Std	1.6	1.6	1.6
		Median	2.9	3.2	2.8
Openstack	Core	Average	4.2	4.7	4.2
		Std	2.6	2.8	2.6
		Median	3.9	4.5	3.5
	Peripheral	Average	4.3	4.9	4.2
		Std	2.7	2.3	2.6
		Median	3.5	3.5	3.7

Table 12: The p-value and effect size of review times in the neutral vs non-neutral comments by the core and peripheral contributors

Project	Review Time	Core		Peripheral	
		p-value	δ	p-value	δ
Eclipse	Positive vs Neutral	0.40	0.009	0.17	N/A
	Negative vs Neutral	0.0002	0.27	0.48	0.151
Android	Positive vs Neutral	2.18E-07	0.18	1.61E-07	0.17
	Negative vs Neutral	2.28E-04	0.31	0.002	0.27
Libreoffice	Positive vs Neutral	9.67E-05	0.23	6.79E-05	0.24
	Negative vs Neutral	3.11E-01	0.38	8.52E-08	0.46
Openstack	Positive vs Neutral	7.82E-09	0.23	3.57E-03	0.27
	Negative vs Neutral	9.20E-04	0.42	6.39E-07	0.31

791 contributors provided neutral comments in the reviews. This trend is similar
792 between the core and peripheral contributors, *i.e.*, negative comments from
793 any contributors tend to increase the review time. Except for Eclipse, the
794 increase in average time taken due to the negative comments is *statistically*
795 *significant* (see Table 12¹⁶). However, we do not see such impact for posi-
796 tive comments. For peripheral contributors, the impact is more prominent.
797 For only one projects (Eclipse), the review time is less when the peripheral
798 contributors provided positive comments.

799 Both the core and peripheral contributors seem to equally impact the
800 review time when they provide positive comments in two projects (Android
801 and Libreoffice).

802 For one project (Eclipse), the positive comments from core contributors
803 seem to impact the review time more than the negative comments from pe-
804 ripheral contributors. For Openstack, the situation is reversed, *i.e.*, the neg-
805 ative comments from the peripheral contributors seem to impact the average
806 review time more than the core contributors.

807 On average, the review time is much less in the reviews where the pe-
808 ripheral contributors provided positive comments. This finding corroborates

¹⁶p-value and effect size are measured using Mann-Whitney U test and the Cliff's Delta effect size as explained in RQ3.1

809 our previous finding that peripheral contributors offer more sentiments in
810 the code reviews, because the core contributors tend to become more neutral
811 over time. Therefore, the happiness of the peripheral contributors seem to
812 be important to reduce the code review time.

RQ3.3 How do the sentiments expressed by the core vs the peripheral contributors correlate with the outcomes of the code reviews?

In all projects except Eclipse, the review times are impacted more by the negative comments from peripheral contributors than the negative comments from the core contributors. In all projects, the review times are longer when the peripheral contributors provided negative comments.

813

814 6. Future Possibilities

815 In all of our studied projects, the reviews with negative sentiments took
816 more time to complete. This observation leads to the question of how we can
817 leverage sentiment analysis to improve *productivity* in a code review process,
818 if the contributors participating in the code reviews can be both the provider
819 and receiver of such negative sentiments. One potential solution would be to
820 design automated sentiment-based monitors that can offer guidance to the
821 contributors. Although such solutions lack authoritativeness, they may nevertheless prove useful to guide the contributors through the different phases of a code review process by mitigating negativity in the review comments.
823 With a view to improve code review outcome and time based on sentiment analysis, we offer the following recommendations by taking cues from our
825 three research questions:
826

- 827 1. Sentiment analysis can be applied to find communities or sub-communities
828 within a project that may be affected by negative comments.
- 829 2. Harmful contributors, such as *bullies* can be detected to ensure that
830 they do not impact the review process negatively.
- 831 3. Controversial reviews can be identified to warn the project leaders
832 about potential controversial features or communities in a project.
- 833 4. Software Bots can be designed to warn the contributors participating
834 in a review when negativity in a review increases.

835 We now discuss the recommendations below.

836 6.1. Community-Based Analysis

837 In RQ3, we built social networks of contributors and observed two major
838 streams of contributors, *core* and *peripheral*. Compared to the peripheral
839 contributors, core members tend to remain with a project for longer time. A
840 deeper understanding of the interactions between the contributors based on
841 social network analysis can offer insights into whether *intrinsic* or *dynamic*
842 sub-communities do exist in modern Gerrit-based code review systems. The
843 identification of such communities can offer several benefits, such as promot-
844 ing a high-performance community to others, offering guidance to a commu-
845 nity that is exchanging negative sentiments but is not productive enough, for
846 examples.

847 6.2. Bullies Among Contributors

848 In all the four studied projects, the reviews with negative sentiments
849 took longer time to get accepted. One possible explanation is that a patch
850 with bugs is likely to be viewed negatively and thus will not be accepted or
851 will be iterated until fixed. However, it is not easily explainable why the
852 negative comments from peripheral contributors impacted the review time
853 more than the core contributors. One potential reason could be that the
854 peripheral contributors are mostly novices to the system compared to the
855 core contributors. Therefore, they would have expressed frustrations due to
856 their lack of understanding of the system.

857 Another possible reason is that there could be *bullies* among the contrib-
858 utors, who try to influence system design and code review outcome using
859 negative comments. Such negative comments can also impact the contrib-
860 utors. Indeed, Mäntylä *et al.* [54] observed that emotions expressed in Jira
861 issues can be correlated to the burnouts of the developers. Ortu *et al.* [36]
862 observed in Jira issues that despite being negative, *bullies* are not necessarily
863 more productive than other developers. An understanding of the role of po-
864 tential *bullies* in code reviews can offer benefits, such as their impact on the
865 code review outcome and productivity. Measures can be taken to detect the
866 *bullies* among the contributors and to remove them from the review process.

867 6.3. Impact of Controversies

868 As we observed in RQ4, regarding reviews taking a long time, the presence
869 of controversy can increase the review time even more. The analysis of con-
870 troversy has proved useful in social media, such as to detect fake news [55]. A

871 deeper analysis of the controversial code review comments can offer insights
872 into the specific reasons behind the comments. For example, it may happen
873 that the product feature (for which the patch is provided) may not be well
874 designed, such that the contributors debate during the review process. It
875 may also happen that the feature is not well-received, such that the contrib-
876 utors have different viewpoints on how to improve it. Therefore, measures
877 can be taken to mitigate the controversies and thus improve the code review
878 outcomes.

879 *6.4. Review Sentiment Bot*

880 Bots have been developed to assist in numerous software development ac-
881 tivities, such as automatically suggesting an answer from Stack Overflow
882 given a query [56], answering questions about an API from documenta-
883 tion [57], or warning developers in a GitHub project if they post negative
884 comments [58]. We can develop similar bots to automatically warn the con-
885 tributors in a code review system with the automated detection of negative
886 comments, their prevalence in the controversies and their proliferation by the
887 bullies. As a first step, we can start with the adaptation of Github sentiment
888 bot [58] for code reviews.

889 *6.5. Gender and cultural aspects*

890 We have investigated gender and cultural aspects bias concerns by defin-
891 ing the following null hypothesis:

892 H0: There is no significant difference in text-based sentiment between male
893 and female contributors.

894 H1: There is no significant difference in text-based sentiment between con-
895 tributors from different countries, which have different language and cultures.

896

897 **Difference between genders - Female and Male - may reveal**
898 **interesting facts under appropriate analysis.** Indeed, recent studies
899 discussed gender bias regarding productivity, in terms of commits, in OSS
900 projects [59, 60]. Moreover, Terrell et al. [61] reported that when new fe-
901 male contributors are identifiable, they have 12% lower chance of getting
902 their pull request accepted than other females whose gender was not iden-
903 tifiable from their profiles. Hence, we are interested in this work to know if
904 there is an association between developers' genders and their expressed senti-
905 ments. More specifically, we formulate the following research questions: **Are**
906 **females' contributors more likely to be positive/neutral/negative**

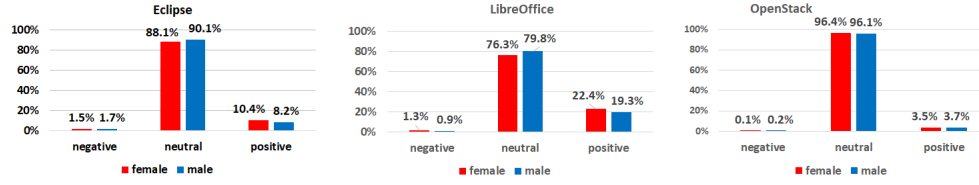


Figure 10: Distribution of Sentiments According to Gender.

907 **than males? Is the proportion of females that express negative**
 908 **sentiments the same as the proportion of males?** To answer these
 909 questions, we segregated contributors according to their gender. We used
 910 the NamSor¹⁷ API to classify contributors into binary gender given personal
 911 names, country of origin, and ethnicity. This API infers gender from the
 912 combination of first name, surname, and information of the country. We
 913 found that 6.8% of Eclipse contributors were females and 88.9% were male,
 914 and 4.4% unknown. LibreOffice respectively (9.4% ; 86.5% ; 4.1%) ; and
 915 OpenStack(10.9% ; 79.6% ; 9.5%). Unfortunately, we were not able to resolve
 916 genders for the Android project because of (encrypted name and email). Fig-
 917 ure 10 shows the distribution of sentiments across gender for three projects.

918 One first observation is that women and men seem to exhibit the same
 919 distribution of sentiments. We performed further statistical analysis to verify
 920 how genders differed in their expressed sentiments. Given that the variables
 921 do not exhibit a normal distribution, we performed a (non-parametric) Mann-
 922 Whitney-Wilcoxon test, with a confidence level of $\alpha = 0.05$. We found that,
 923 overall for the three projects, the tests are statistically significant ($p < 0.05$, Z
 924 statistic of -1.018) and thus we reject our null hypothesis H_0 . We claim that
 925 there is a significant difference in the distribution of text-based sentiment
 926 between male and female contributing to OSS projects. This result confirms
 927 previous findings by Paul et al. [34].

928 We also investigated the impact of the country origin for the top 5% core
 929 contributors within the three projects aiming to investigate the impact of
 930 the first language and cultural aspects of these contributors on code reviews.
 931 Figure 11 shows the geographical distribution of the top 5% contributors.
 932 We performed a Kruskal-Wallis statistical test to verify whether samples
 933 (i.e., different countries) have the same distribution of expressed sentiments.

¹⁷<https://www.namsor.com/>

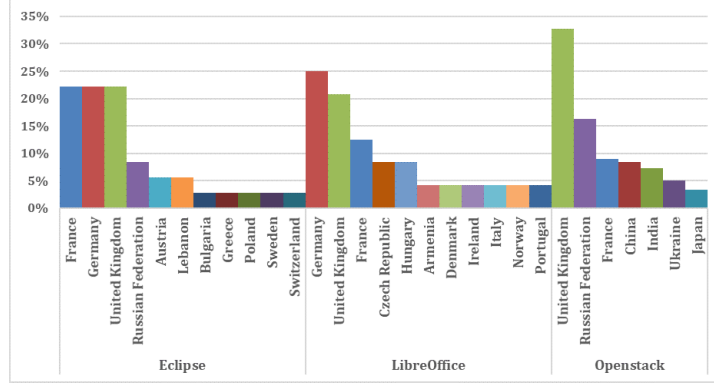


Figure 11: Countries distribution for top 5% contributors.

Kruskal-Wallis test results reveal that distribution differences are statistically significant ($p\text{-value} < 2.2e-16$), we reject our null hypothesis H1 and state that there is a statistically significant difference in expressed sentiments according to the country of origin.

However, studying the effect of gender cultural aspects on code reviews are beyond the scope of this study. We will address this concern in future work.

7. Threats to Validity

Threats to construct validity are mainly related to the accuracy of the tool used for sentiment analysis. We strengthen our sampling approach for manual annotation by using opportunistic sampling. The authors manually examined 2,220 comments. In general, we observed that the sentiments expressed in code reviews are easy to analyze due to the unambiguous nature of text-based sentiments expressed in the code reviews comments, which were understood by both coders with relative ease.

Threats to internal validity concern factors internal to our study that may affect our findings. The primary threat to internal validity in this study relates to project selection. One possible threat is that the retrieved dataset is too small and somehow is not representative enough. We were cautious to choose OSS projects with the following characteristics: (1) long-lived projects with dynamic communities around; (2) the community uses the review tools Gerrit to carry out code reviews activities. We also paid attention not to violate assumptions of the statistical tests, for example we

957 applied non-parametric tests that do not require making assumptions on the
958 normality of our data set.

959 In addition, we used propensity score matching [62] to eliminate the bias
960 that could be introduced by technical characteristics. We compared the
961 distribution of estimated propensity score between Positive and Negative
962 reviews in the matched sample and obtained an average of 96% of overlap,
963 which means that we are dealing with a homogeneous data set of reviews
964 based on the observed covariate values. This provides confidence that the
965 observed results are not due to structural differences in the patches (*i.e.*, we
966 are not comparing large patches with small patches, etc.). However other
967 technical characteristics can be considered such as the number of sentences
968 in the comments and code complexity.

969 *Threats to external validity* concern the generalization of our findings. In
970 the context of RQ1 we performed 2,220 manual classifications. We are aware
971 that the quality and size of the annotated set may impact the sentiment
972 classification accuracy. While the human raters are knowledgeable in mining
973 software repositories, sentiment analysis and empirical analysis, their judge-
974 ment may be impacted by the absence of related in-depth information of the
975 studied systems in the dataset, e.g., whether the reviewers in those studied
976 systems exhibited any latent communities as reported by Bird et al. [63].
977 Furthermore, our study involves only four projects. Thus, we should recog-
978 nize that our conclusions may not be generalizable to other systems. We are
979 also aware that the context of each project including the technical complexity
980 and organization are important factors that can limit generalization. How-
981 ever, these projects are among the most studied projects in the literature and
982 the system’s data are publicly available. We also have the opportunity to
983 perform a longitudinal study over more than five years, which mitigates the
984 risks related to cultural aspects. Yet, replication of our work on other open
985 and close source systems is desirable in order to generalize our conclusions.

986 *Threats to reliability validity* refers to the degree to which the same data
987 would lead to the same results when the study’s design is replicated. Our
988 research aims at investigating expressed sentiment by developers on reviews.
989 Our methodology for data analysis and results are well documented in this
990 paper. The tools are available [39] and our datasets are publicly available
991 online [43]. Also, all the participants of the manual tagging have a back-
992 ground in computer science; we are confident that reviews comments have
993 been interpreted according to the perspective of software engineers. We did
994 not involve raters with a different background, because they may overlook

995 or misinterpret the terms used by developers. However, RQ2.1 reveals that
 996 most comments in the dataset have neutral sentiments, while only less than
 997 3% of the comments are negative which may have an impact on our analysis
 998 and results. In RQ3.2, we assessed whether reviews with sentiments took a
 999 shorter/longer time than the reviews with neutral sentiments. We noticed
 1000 a low number of negative sentiments, similarly observed in previous studies
 1001 that used datasets from Stack Overflow (e.g., the Stack Overflow dataset by
 1002 Lin et al. [8] has around 75% neutral comments). Therefore, although the low
 1003 number of negative comments may introduce a threat to the generalizability
 1004 of our results across other systems, our analysis remains applicable to other
 1005 systems. In addition, we assessed the impact of sentiments on code review
 1006 outcomes in RQ3.3 by comparing the time taken for reviews with positive
 1007 comments vs the reviews with negative comments. We found on average
 1008 13.94% positive and 2.24% negative comments in the four studied systems.

1009 8. Conclusions

1010 We have analyzed developers' comments on reviews using historical data
 1011 from four open source projects. We aimed at investigating the influence of
 1012 text-based expressed sentiments on the code review duration and its outcome.
 1013 Using the best performing sentiment detection tool, we found that contrib-
 1014 utors do express sentiments when they are reviewing and commenting each
 1015 other's code. Also, we investigated the influence of expressed sentiments
 1016 within developers' comments on the time and outcome of the code review
 1017 process. We found that expressing positive sentiment when reviewing source
 1018 code have an influence on reviews duration time; in average it could save
 1019 1.32 days on the review completion time. Moreover, our findings indicate
 1020 that negative comments are likely to increase the proportion of unsuccessful
 1021 reviews.

1022 From a social network perspective, we used a K -means clustering ap-
 1023 proach based on SNA centrality measures to discern between core and pe-
 1024 ripheral contributors. We found that different contributors within the peer
 1025 review collaboration social network express different sentiments, with core
 1026 contributors expressing mostly neutral sentiments.

1027 Our work contributes theoretically and empirically to the body of OSS
 1028 research and has practical implications on sentiment awareness within OSS.
 1029 We hope that our work will inspire more studies on developing efficient tools
 1030 to help OSS contributors improve their productivity. As future work, we plan

to complement this quantitative study with a qualitative exploration aiming at gaining more understanding of the influence of expressed sentiments on code revision workflow. Also, we plan to investigate the effect of developers' expressed sentiment on contributor's engagement and/or turnover.

References

- [1] A. Bosu, J. C. Carver, Impact of peer code review on peer impression formation: A survey, in: Empirical Software Engineering and Measurement, 2013 ACM/IEEE International Symposium on, IEEE, pp. 133–142.
- [2] A. Murgia, P. Tourani, B. Adams, M. Ortu, Do developers feel emotions? an exploratory analysis of emotions in software artifacts, in: Proceedings of the 11th International Conference on Mining Software Repositories, MSR'14, pp. 262–271.
- [3] B. Pang, L. Lee, Opinion mining and sentiment analysis, Foundations and Trends in Information Retrieval 2 (2008) 1–135.
- [4] O. Kucuktunc, B. B. Cambazoglu, I. Weber, H. Ferhatosmanoglu, A large-scale sentiment analysis for yahoo! answers, in: Proceedings of the Fifth International Conference on Web Search and Data Mining, WSDM '12, pp. 633–642.
- [5] D. Garcia, M. S. Zanetti, F. Schweitzer, The role of emotions in contributors activity: A case study on the gentoo community, in: Cloud and green computing (CGC), 2013 third international conference on, IEEE, pp. 410–417.
- [6] O. Baysal, O. Kononenko, R. Holmes, M. Godfrey, The influence of non-technical factors on code review, in: proceedings of the 20th Working Conference on Reverse Engineering, pp. 122–131.
- [7] V. Efstathiou, D. Spinellis, Code review comments: Language matters, in: Proceedings of the 40th International Conference on Software Engineering (ICSE'18), p. 4.
- [8] B. Lin, F. Zampetti, G. Bavota, M. Di Penta, M. Lanza, R. Oliveto, Sentiment analysis for software engineering: How far can we go?, in:

1062 Proceedings of the 40th International Conference on Software Engineer-
1063 ing (ICSE'18), p. 11.

1064 [9] N. Imtiaz, J. Middleton, P. Girouard, E. Murphy-Hill, Sentiment and
1065 politeness analysis tools on developer discussions are unreliable, but
1066 so are people, in: Proceedings of the 3rd International Workshop on
1067 Emotion Awareness in Software Engineering SEmotion'18, pp. 55–61.

1068 [10] N. Novielli, D. Girardi, F. Lanubile, A benchmark study on sentiment
1069 analysis for software engineering research, in: Proceedings of the 15th
1070 International Conference on Mining Software Repositories, p. 12.

1071 [11] F. Calefato, F. Lanubile, F. Maiorano, N. Novielli, Sentiment polarity
1072 detection for software development, Empirical Software Engineering
1073 (2017) 31.

1074 [12] T. Ahmed, A. Bosu, A. Iqbal, S. Rahimi, Sentier: A customized senti-
1075 ment analysis tool for code review interactions, in: Proceedings of the
1076 32nd International Conference on Automated Software Engineering, pp.
1077 106–111.

1078 [13] M. R. Islam, M. F. Zibran, Leveraging automated sentiment analy-
1079 sis in software engineering, in: Proceedings of the 14th International
1080 Conference on Mining Software Repositories, MSR '17, pp. 203–214.

1081 [14] M. D. Munezero, C. S. Montero, E. Sutinen, J. Pajunen, Are they
1082 different? affect, feeling, emotion, sentiment, and opinion detection in
1083 text, IEEE transactions on affective computing 5 (2014) 101–111.

1084 [15] W. G. Parrott, Emotions in Social Psychology, Psychology Press, 2001.

1085 [16] B. Pang, L. Lee, S. Vaithyanathan, Thumbs up? sentiment classifica-
1086 tion using machine learning techniques, in: Conference on Empirical
1087 Methods in Natural Language Processing, pp. 79–86.

1088 [17] B. Liu, Sentiment Analysis and Opinion Mining, Morgan & Claypool
1089 Publishers, 2012.

1090 [18] P. Weissgerber, D. Neu, S. Diehl, Small patches get in!, in: Proceed-
1091 ings of the 2008 International Working Conference on Mining Software
1092 Repositories, MSR '08, pp. 67–76.

- 1093 [19] O. Kononenko, O. Baysal, L. Guerrouj, Y. Cao, M. Godfrey, Investi-
 1094 gating code review quality: Do people and participation matter?, in:
 1095 Proceedings on the International Conference on Software Maintenance
 1096 and Evolution (ICSME'15), pp. 111–120.
- 1097 [20] M. Beller, A. Bacchelli, A. Zaidman, E. Juergens, Modern code reviews
 1098 in open-source projects: Which problems do they fix?, in: Proceedings
 1099 of the 11th Working Conference on Mining Software Repositories, MSR
 1100 2014, pp. 202–211.
- 1101 [21] O. Kononenko, O. Baysal, M. W. Godfrey, Code review quality: How
 1102 developers see it, in: Proceedings of the 38th International Conference
 1103 on Software Engineering, ICSE '16, pp. 1028–1038.
- 1104 [22] A. Thavaneswaran, Propensity score matching in observational studies,
 1105 Manitoba Center for Health Policy. (2008).
- 1106 [23] M. E. Newman, The structure and function of complex networks, SIAM
 1107 review 45 (2003) 167–256.
- 1108 [24] L. C. Freeman, The development of social network analysis—with an em-
 1109 phasis on recent events, The SAGE handbook of social network analysis
 1110 21 (2011) 26–39.
- 1111 [25] X. Yang, Social network analysis in open source software peer review,
 1112 in: Proceedings of the 22Nd ACM SIGSOFT International Symposium
 1113 on Foundations of Software Engineering, FSE'14, pp. 820–822.
- 1114 [26] K. Crowston, K. Wei, Q. Li, J. Howison, Core and periphery in free/libre
 1115 and open source software team communications, in: Proceedings of the
 1116 39th International Conference on System Sciences, pp. 118.1–.
- 1117 [27] I. Robertson, C. Cooper, Well-being: Productivity and happiness at
 1118 work, Springer, 2011.
- 1119 [28] M. Thelwall, K. Buckley, G. Paltoglou, Sentiment strength detection for
 1120 the social web, Journal of the American Society for Information Science
 1121 and Technology 61 (2012) 2544–2558.
- 1122 [29] R. Jongeling, S. Datta, A. Serebrenik, Choosing your weapons: On
 1123 sentiment analysis tools for software engineering research, in: Software

- 1124 maintenance and evolution (ICSME), 2015 IEEE international confer-
1125 ence on, IEEE, pp. 531–535.
- 1126 [30] J. Guillory, J. Spiegel, M. Drislane, B. Weiss, W. Donner, J. Hancock,
1127 Upset now?: Emotion contagion in distributed groups, in: Proceedings
1128 of the SIGCHI Conference on Human Factors in Computing Systems,
1129 CHI '11, pp. 745–748.
- 1130 [31] E. Guzman, B. Bruegge, Towards emotional awareness in software de-
1131 velopment teams, in: Proceedings of the 2013 9th Joint Meeting on
1132 Foundations of Software Engineering, ESEC/FSE 2013, pp. 671–674.
- 1133 [32] V. Sinha, A. Lazar, B. Sharif, Analyzing developer sentiment in commit
1134 logs, in: Proceedings of the 13th International Conference on Mining
1135 Software Repositories, MSR '16, pp. 520–523.
- 1136 [33] E. Guzman, D. Azócar, Y. Li, Sentiment analysis of commit comments
1137 in github: An empirical study, in: Proceedings of the 11th Working
1138 Conference on Mining Software Repositories, MSR'14, pp. 352–355.
- 1139 [34] R. Paul, A. Bosu, K. Z. Sultana, Expressions of sentiments during code
1140 reviews: Male vs. female, in: Proceedings of the 16th International Con-
1141 ference on Software Analysis, Evolution and Reengineering SANER'19,
1142 pp. 15–26.
- 1143 [35] I. A. Khan, W.-P. Brinkman, R. M. Hierons, Do moods affect program-
1144 mers' debug performance?, *Cognition, Technology & Work* 13 (2011)
1145 245–258.
- 1146 [36] M. Ortu, B. Adams, G. Destefanis, P. Tourani, M. Marchesi, R. Tonelli,
1147 Are bullies more productive?: Empirical study of affectiveness vs. issue
1148 fixing time, in: Proceedings of the 12th Working Conference on Mining
1149 Software Repositories, MSR'15, pp. 303–313.
- 1150 [37] G. Destefanis, M. Ortu, S. Counsell, S. Swift, M. Marchesi, R. Tonelli,
1151 Software development: do good manners matter?, *PeerJ Computer Sci-*
1152 *ence* 2 (2016) e73.
- 1153 [38] M. M. Rahman, C. K. Roy, R. G. Kula, Predicting usefulness of code
1154 review comments using textual features and developer experience, in:

1155 Proceedings of the 14th International Conference on Mining Software
1156 Repositories, MSR '17, pp. 215–226.

1157 [39] X. Yang, R. G. Kula, N. Yoshida, H. Iida, Mining the modern code
1158 review repositories: A dataset of people, process and product, in: Pro-
1159 ceedings of the 13th International Conference on Mining Software Repos-
1160 itories, ACM, pp. 460–463.

1161 [40] L. Guo, P. Qu, R. Zhang, D. Zhao, H. Wang, R. Liu, B. Mi, H. Yan,
1162 S. Dang, Propensity score-matched analysis on the association between
1163 pregnancy infections and adverse birth outcomes in rural northwestern
1164 china, *Scientific reports* 8 (2018) 5154.

1165 [41] A. Dmitrienko, G. Molenberghs, C. Chuang-Stein, W. W. Offen, Anal-
1166 ysis of clinical trials using SAS: A practical guide, SAS Institute, 2005.

1167 [42] N. Novielli, F. Calefato, F. Lanubile, A gold standard for emotion an-
1168 notation in stack overflow, in: Proceedings of the 15th International
1169 Conference on Mining Software Repositories, MSR '18, pp. 14–17.

1170 [43] I. E. Asri, N. Kerzazi, G. Uddin, F. Khomh, An Empirical Study of Sen-
1171 timents in Code Reviews (online appendix), https://github.com/ikramElasri/SentiAnalysis_CodeReview, October 2018 (last ac-
1172 cessed).
1173

1174 [44] D. M. Christopher, R. Prabhakar, S. Hinrich, Introduction to infor-
1175 mation retrieval, *An Introduction To Information Retrieval* 151 (2008)
1176 5.

1177 [45] B. Bostanci, E. Bostanci, An evaluation of classification algorithms using
1178 mc nemar’s test, in: Proceedings of Seventh International Conference
1179 on Bio-Inspired Computing: Theories and Applications (BIC-TA 2012),
1180 Springer, pp. 15–26.

1181 [46] J. M. George, A. P. Brief, Feeling good-doing good: a conceptual anal-
1182 ysis of the mood at work-organizational spontaneity relationship., *Psy-
1183 chological bulletin* 112 (1992) 310.

1184 [47] G. Uddin, F. Khomh, Mining API Aspects in API Reviews, Tech-
1185 nical Report, Technical Report. 10 pages. <http://swat.polymtl.ca/data/opinionvalue>, 2017.
1186

- 1187 [48] G. Uddin, O. Baysal, L. Guerrouj, F. Khomh, Understanding how and
1188 why developers seek and analyze api-related opinions, *IEEE Transactions on Software Engineering* (2019).
1189
- 1190 [49] P. Tourani, Y. Jiang, B. Adams, Monitoring sentiment in open source
1191 mailing lists: Exploratory study on the apache ecosystem, in: *Proceedings of 24th Annual International Conference on Computer Science and Software Engineering, CASCON '14, IBM Corp., Riverton, NJ, USA, 2014*, pp. 34–44.
1192
1193
1194
- 1195 [50] A. Bosu, J. C. Carver, Impact of developer reputation on code review
1196 outcomes in oss projects: An empirical investigation, in: *Proceedings of the 8th International Symposium on Empirical Software Engineering and Measurement, ESEM '14*, pp. 33:1–33:10.
1197
1198
- 1199 [51] G. P. Sudhakar, A. Farooq, S. Patnaik, Measuring productivity of soft-
1200 ware development teams, *Journal of Management* 7 (2012) 65–75.
- 1201 [52] J. W. Tukey, *Exploratory data analysis*, volume 2, 1977.
- 1202 [53] A.-M. Popescu, M. Pennacchiotti, Detecting controversial events from
1203 twitter, in: *Proceedings of the 19th ACM International Conference on Information and Knowledge Management, CIKM '10*, pp. 1873–1876.
1204
- 1205 [54] M. Mäntylä, B. Adams, G. Destefanis, D. Graziotin, M. Ortu, Mining
1206 valence, arousal, and dominance – possibilities for detecting burnout
1207 and productivity?, in: *Proceedings of the 13th Working Conference on Mining Software Repositories*, pp. 247–258.
1208
- 1209 [55] K. Garimella, G. D. F. Morales, A. Gionis, M. Mathioudakis, Quanti-
1210 fying controversy on social media, *Transactions on Social Computing* 1
1211 (2018) Article no. 3.
- 1212 [56] S. Zamanirad, B. Benatallah, M. C. Barukh, F. Casati, Programming
1213 bots by synthesizing natural language expressions into api invocations,
1214 in: *Proceedings of the 32nd International Conference on Automated Software Engineering*, pp. 832–837.
1215
- 1216 [57] Y. Tian, F. Thung, A. Sharma, D. Lo, Apibot: question answering
1217 bot for api documentation, in: *Proceedings of the 32nd International Conference on Automated Software Engineering*, pp. 153–158.
1218

- 1219 [58] GitHub, Sentiment Bot, <https://github.com/apps/sentiment-bot>,
1220 18 May 2018 (last accessed).
- 1221 [59] B. Vasilescu, D. Posnett, B. Ray, M. G. van den Brand, A. Serebrenik,
1222 P. Devanbu, V. Filkov, Gender and tenure diversity in github teams,
1223 in: Proceedings of the 33rd Annual ACM Conference on Human Factors
1224 in Computing Systems, CHI '15, ACM, New York, NY, USA, 2015, pp.
1225 3789–3798.
- 1226 [60] C. Mendez, H. S. Padala, Z. Steine-Hanson, C. Hilderbrand, A. Horvath,
1227 C. Hill, L. Simpson, N. Patil, A. Sarma, M. Burnett, Open source
1228 barriers to entry, revisited: A sociotechnical perspective, in: Proceedings
1229 of the 40th International Conference on Software Engineering, ICSE '18,
1230 ACM, New York, NY, USA, 2018, pp. 1004–1015.
- 1231 [61] J. Terrell, A. Kofink, J. Middleton, C. Raine, E. Murphy-Hill,
1232 C. Parnin, J. Stallings, Gender differences and bias in open source:
1233 pull request acceptance of women versus men, *PeerJ Computer Science*
1234 3 (2017) e111.
- 1235 [62] Z. Luo, J. C. Gardiner, C. J. Bradley, Applying propensity score meth-
1236 ods in medical research: pitfalls and prospects, *Medical Care Research*
1237 and Review 67 (2010) 528–554.
- 1238 [63] C. Bird, D. Pattison, R. D’Souza, V. Filkov, P. Devanbu, Latent social
1239 structure in open source projects, in: Proceedings of the 26th ACM
1240 SIGSOFT International Symposium on Foundations of software engi-
1241 neering, ACM, pp. 24–35.