

A peer-reviewed version of this preprint was published in PeerJ on 13 May 2019.

[View the peer-reviewed version](https://peerj.com/articles/cs-193) (peerj.com/articles/cs-193), which is the preferred citable publication unless you specifically need to cite this preprint.

di Biase M, Bruntink M, van Deursen A, Bacchelli A. 2019. The effects of change decomposition on code review—a controlled experiment. PeerJ Computer Science 5:e193 <https://doi.org/10.7717/peerj-cs.193>

The effects of change decomposition on code review - a controlled experiment

Marco di Biase ^{Corresp., 1, 2}, Magiel Bruntink ², Arie van Deursen ¹, Alberto Bacchelli ³

¹ Delft University of Technology, Delft, The Netherlands

² Software Improvement Group, Amsterdam, The Netherlands

³ University of Zurich, Zurich, Switzerland

Corresponding Author: Marco di Biase

Email address: m.dibiase@tudelft.nl

Background. Code review is a cognitively demanding and time-consuming process. Previous qualitative studies hinted at how decomposing change sets into multiple yet internally coherent ones would improve the reviewing process. So far, no quantitative analysis of this hypothesis has been provided.

Aims. (1) Quantitatively measure the effects of change decomposition on the outcome of code review (in terms of number of found defects, wrongly reported issues, suggested improvements, time, and understanding); (2) Qualitatively analyze how subjects approach the review and navigate the code, building knowledge and addressing existing issues, in large vs. decomposed changes.

Method. Controlled experiment using the pull-based development model involving 28 software developers among professionals and graduate students.

Results. Change decomposition leads to fewer wrongly reported issues, influences how subjects approach and conduct the review activity (by increasing context-seeking), yet impacts neither understanding the change rationale nor the number of found defects.

Conclusions. Change decomposition reduces the noise for subsequent data analyses but also significantly supports the tasks of the developers in charge of reviewing the changes. As such, commits belonging to different concepts should be separated, adopting this as a best practice in software engineering.

The Effects of Change Decomposition on Code Review - A Controlled Experiment

Marco di Biase^{1,2}, Magiel Bruntink², Arie van Deursen¹, and Alberto Bacchelli³

¹Delft University of Technology - Delft, The Netherlands

²Software Improvement Group - Amsterdam, The Netherlands

³University of Zurich - Zurich, Switzerland

Corresponding author:

Marco di Biase¹

Email address: m.dibiase@tudelft.nl

ABSTRACT

Background. Code review is a cognitively demanding and time-consuming process. Previous qualitative studies hinted at how decomposing change sets into multiple yet internally coherent ones would improve the reviewing process. So far, no quantitative analysis of this hypothesis has been provided.

Aims. (1) Quantitatively measure the effects of change decomposition on the outcome of code review (in terms of number of found defects, wrongly reported issues, suggested improvements, time, and understanding); (2) Qualitatively analyze how subjects approach the review and navigate the code, building knowledge and addressing existing issues, in large vs. decomposed changes.

Method. Controlled experiment using the pull-based development model involving 28 software developers among professionals and graduate students.

Results. Change decomposition leads to fewer wrongly reported issues, influences how subjects approach and conduct the review activity (by increasing context-seeking), yet impacts neither understanding the change rationale nor the number of found defects.

Conclusions. Change decomposition reduces the noise for subsequent data analyses but also significantly supports the tasks of the developers in charge of reviewing the changes. As such, commits belonging to different concepts should be separated, adopting this as a best practice in software engineering.

1 INTRODUCTION

Code review is the activity performed by software teams to check code quality, with the purpose of identifying issues and shortcomings (Bacchelli and Bird, 2013). Nowadays, reviews are mostly performed in an iterative, informal, change- and tool-based fashion, also known as Modern Code Review (MCR) (Cohen, 2010). Both open-source and industry software teams employ MCR to check code changes before being integrated in their codebases (Rigby and Bird, 2013). Past research has provided evidence that MCR is associated with improved key software quality aspects, such as maintainability (Morales et al., 2015) and security (di Biase et al., 2016), as well as with less defects (McIntosh et al., 2016).

Reviewing a source code change is a cognitively demanding process. Researchers provided evidence that understanding the code change under review is among the most challenging tasks for reviewers (Bacchelli and Bird, 2013). In this light, past studies have argued that code changes that—at the same time—address multiple, possibly unrelated concerns (also known as *noisy* (Murphy-Hill et al., 2012) or *tangled changes* (Herzig and Zeller, 2013)) can hinder the review process (Herzig and Zeller, 2013; Kirinuki et al., 2014), by increasing the cognitive load for reviewers. Indeed, it is reasonable to think that grasping the rationale behind a change that spans multiple concepts in a system requires more effort than the same patch committed separately. Moreover, the noise could put a reviewer on a wrong track, thus leading to missing defects (*false negatives*) or to raising unfounded issues in sound code (*false positives* in this paper).

Qualitative studies reported that professional developers perceive tangled code changes as problematic

47 and asked for tools to automatically decompose them (Tao et al., 2012; Barnett et al., 2015). Accordingly,
48 change untangling mechanisms have been proposed (Tao and Kim, 2015; Dias et al., 2015; Barnett et al.,
49 2015).

50 Although such tools are expectedly useful, the effects of change decomposition on review is an open
51 research problem. Tao and Kim presented the earliest and most relevant results in this area (Tao and Kim,
52 2015), showing that change decomposition allows practitioners to achieve their tasks better in a similar
53 amount of time.

54 In this paper, we continue on this research line and focus on evaluating the effects of change decom-
55 position on code review. We aim at answering questions, such as: Is change decomposition beneficial for
56 understanding the rationale of the change? Does it have an impact on the number/types of issues raised?
57 Are there differences in time to review? Are there variations with respect to defect lifetime?

58 To this end, we designed a controlled experiment focusing on pull requests, a widespread approach to
59 submit and review changes (Gousios et al., 2015). Our work investigates whether the results from Tao and
60 Kim (Tao and Kim, 2015) can be replicated, and extend the knowledge on the topic. With a Java system
61 as a subject, we asked 28 software developers among professionals and graduate students to review a
62 refactoring and a new feature (according to professional developers (Tao et al., 2012), these are the most
63 difficult to review when tangled). We measure how the partitioning vs. non-partitioning of the changes
64 impacts defects found, false positive issues, suggested improvements, time to review, and understanding
65 the change rationale. We also perform qualitative observations on how subjects conduct the review and
66 address defects or raise false positives, in the two scenarios.

67 This paper makes the following contributions:

- 68 • the design of an asynchronous controlled experiment to assess the benefits of change decomposition
69 in code review using pull requests, available for replication (di Biase et al., 2018);
- 70 • empirical evidence that change decomposition in the pull-based review environment leads to fewer
71 false positives.

72 The paper proceeds as follows: Section 2 illustrates the related work; Section 3 describes our research
73 objectives; the design of our experiment is described in Section 4; threats to validity are discussed in
74 Section 5; results are presented in Section 6; Section 7 reports the discussion based on the results; finally,
75 Section 8 summarizes our study.

76 2 RELATED WORK

77 Several studies explored tangled changes and concern separation in code reviews. Tao and Kim in-
78 vestigated the role of understanding code changes during the software development process, exploring
79 practitioners' needs (Tao et al., 2012). Their study outlined that grasping the rationale when dealing with
80 the process of code review is indispensable. Moreover, to understand a composite change, it is useful
81 to break it into smaller ones each concerning a single issue. Rigby et al. empirically studied the peer
82 review process for six large, mature OSS projects, showing that small change size is essential to the more
83 fine-grained style of peer review (Rigby et al., 2014). Kirinuki et al. provided evidence about problems
84 with the presence of multiple concepts in a single code change (Kirinuki et al., 2014). They showed that
85 these are unsuitable for merging code from different branches, and that tangled changes are different to
86 review because practitioners have to seek the changes for the specified task in the commit.

87 Regarding empirical controlled experiments on the topic of code reviews, the most relevant work is
88 by Uwano et al. (2006). They used an eye-tracker to characterize the performance of subjects reviewing
89 source code. Their experimentation environment enabled them to identify a pattern called *scan*, consisting
90 of the reviewer reading the entire code before investigating the details of each line. In addition, their
91 qualitative analysis found that participants who did not spend enough time during the *scan* took more
92 time to find defects. Uwano's experiment was replicated by Sharif et al. (2012). Their results indicated
93 that the longer participants spent in the *scan*, the quicker they were able to find the defect. Conversely,
94 review performance decreases when participants did not spend sufficient time on the scan, because they
95 find irrelevant lines.

96 Even if MCR is now a mainstream process, adopted in both open source and industrial projects, we
97 found only two studies on change partitioning and its benefits for code review. The work by Barnett et al.
98 (2015) analyzed the usefulness of an automatic technique for decomposing changesets. They found a

99 positive association between change decomposition and the level of understanding of the changesets.
100 According to their results, this would help time to review as the different contexts are separated. Tao and
101 Kim (2015) proposed a heuristic-based approach to decompose changeset with multiple concepts. They
102 conducted a user study with students investigating whether their untangling approach affected the time
103 and the correctness in performing review-related tasks. Results were promising: Participants completed
104 the tasks better with untangled changes in a similar amount of time. In spite of the innovative techniques
105 they proposed to untangle code changes and in these promising results, the evaluation of effects of change
106 decomposition was preliminary.

107 In contrast, our research focuses on setting up and running an experiment to empirically assess the
108 benefits of change decomposition for the process of code review, rather than evaluating the performances
109 of an approach.

110 3 MOTIVATION AND RESEARCH OBJECTIVES

111 3.1 Experiment definition and context

112 Our analysis of the literature showed that there is only preliminary empirical evidence on how code
113 review decomposition affects its outcomes, its change understanding, time to completion, effectiveness
114 (i.e., number of defects found), false positives (issues mistakenly identified as defect by the reviewer),
115 and suggested improvements. This motivates us in setting up a controlled experiment, exploiting the
116 popular pull-based development model, to assess the conjecture that a proper separation of concerns in
117 code review is beneficial to the efficiency and effectiveness of the review.

118 Pull requests feature asynchronous, tool-based activities in the bigger scope of pull-based software
119 development (Gousios et al., 2014). The pull-based software process features a distributed environment
120 where changes to a system are proposed through patch submissions that are pulled and merged locally,
121 rather than being directly pushed to a central repository.

122 Pull requests are the way contributors submit changes for review in GitHub. Change acceptance has
123 to be granted by other team members called integrators (Gousios et al., 2015). They have the crucial role
124 of managing and integrating contributions and are responsible for inspecting the changes for functional
125 and non-functional requirements. 80% of integrators use pull requests as the means to review changes
126 proposed to a system (Gousios et al., 2015).

127 In the context of distributed software development and change integration, GitHub is one of the most
128 popular code hosting sites with support for pull-based development. GitHub pull requests contain a branch
129 from which changes are compared by an automatic discovery of commits to be merged. Changes are then
130 reviewed online. If further changes are requested, the pull request can be updated with new commits to
131 address the comments. The inspection can be repeated and, when the patch set fits the requirements, the
132 pull request can be merged to the master branch.

133 3.2 Research questions

134 The motivation behind modern code review is to find defects and improve code quality (Bacchelli and
135 Bird, 2013). We are interested in checking if reviewers are able to address *defects* (referred in this paper
136 as *effectiveness*). Furthermore, we focus on comments pointing out *false positives* (wrongly reported
137 defects), and *suggested improvements* (non-critical non-functional issues such as suggested refactorings).
138 Suggested improvements highlight reviewer participation (McIntosh et al., 2014) and these comments are
139 generally considered very useful (Bosu et al., 2015). Our first research question is:

140 **RQ1.** Do tangled pull requests influence *effectiveness* (i.e., number of defects found), *false positives*,
141 and *suggested improvements* of reviewers, when compared to untangled pull requests?

Based on the first research question, we formulate the following null-hypotheses for (statistical) testing:

Tangled pull requests do not reduce:

H_{0e} the effectiveness of the reviewers during peer-review

H_{0f} the false positives detected by the reviewers during peer-review

H_{0c} the suggested improvements written by the reviewers during peer-review

142

143 Given the structure and the settings of our experimentation, we can also measure the time spent on
144 review activity and defect lifetime. Thus, our next research question is:

145 **RQ2.** Do tangled pull requests influence the time necessary for a review and defect lifetime, when
146 compared to untangled pull requests?

For the second research question, we formulate the following null-hypotheses:

Tangled pull requests do not reduce:
H_{0t1} time to review
H_{0t2} defect lifetime

147 Further details on how we measure time and define defect lifetime are described in Section 4.7.

148 In our study, we aim to measure whether change decomposition has an effect on understanding the
149 rationale of the change under review. Understanding the rationale is the most important information need
150 when analyzing a change, according to professional software developers (Tao et al., 2012). As such, the
151 question we set to answer is:
152

153 **RQ3.** Do tangled pull requests influence the reviewers' understanding of the change rationale, when
154 compared to untangled ones?

For our third research question, we test the following null-hypotheses:

Tangled pull requests do not reduce:
H_{0u} change-understanding of reviewers during peer-review
when compared to untangled pull requests

155 Finally, we qualitatively investigate how participants individually perform the review to understand
156 how they address defects or potentially raise false positives. Our last research question is then:
157

158 **RQ4.** What are the differences in patterns and features used between reviews of tangled and
159 untangled pull requests?

160 4 EXPERIMENTAL DESIGN AND METHOD

161 In this section, we detail how we designed the experiment and the research method that we followed.

162 4.1 Object system chosen for the experiment

163 The system that was used for reviews in the experiment is JPacman, an open-source Java system available
164 on GitHub¹ that emulates a popular arcade game used at Delft University of Technology to teach software
165 testing.

166 The system has about 3,000 lines of code and was selected because a more complex and larger project
167 would require participants to grasp the rationale of a more elaborate system. In addition, the training phase
168 required for the experiment would imply hours of effort, increasing the consequent fatigue that participants
169 might experience. In the end, the experiment targets assessing differences in review partitioning and is
170 tailored for a process rather than a product.

¹<https://github.com/SERG-Delft/jpacman-framework>

171 4.2 Recruiting of the subject participants

172 The study was conducted with 28 participants recruited by means of convenience sampling (Wohlin et al.,
 173 2012) among experienced and professional software developers, PhD, and MSc students.² They were
 174 drawn from a population sample that volunteered to participate. The voluntary nature of participation
 175 implies the consent to use data gathered in the context of this study. Software developers belong to three
 176 software companies, PhD students belong to three universities, and MSc students to different faculties
 177 despite being from the Delft University of Technology. We involved as many different roles as possible
 178 to have a larger sample for our study and increase its external validity. Using a questionnaire, we asked
 179 development experience, language-specific skills, and review experience as number of reviews per week.
 180 We also included a question that asked if a participant knew the source code of the game. Table 1 reports
 181 the results of the questionnaire, which are used to characterize our population and to identify key attributes
 182 of each subject participant.

TABLE 1. DESCRIPTIVE DATA OF THE SUBJECT PARTICIPANTS

Group	# of subjects		Role	FTE Experience		Reviews per week			
	total	with system knowledge		μ	σ	per role		per group	
						μ	σ	μ	σ
Control (tangled changes)	6	2 (33%)	SW Developer	4.3	4.8	4.8	3.3		
	3	1 (33%)	PhD Student	5.0	2.9	3.0	2.9	3.6	3.6
	5	3 (60%)	MSc Student	2.2	0.7	2.6	3.8		
Treatment (untangled changes)	6	2 (33%)	SW Developer	4.8	2.9	3.3	3.4		
	3	1 (33%)	PhD Student	6.0	6.6	2.0	0.8	4.0	6.4
	5	3 (60%)	MSc Student	2.2	1.1	6.0	9.0		

183 4.3 Monitoring versus realism

184 In line with the nature of pull-based software development and its peer review with pull requests, we
 185 designed the experimentation phase to be executed asynchronously. This implies that participants could
 186 run the experiment when and where they felt most comfortable, with no explicit constraints for place,
 187 time or equipment.

188 With this choice, we purposefully gave up some degree of control to increase realism. Having a more
 189 strictly controlled experimental environment would not replicate the usual way of running such tasks (that
 190 is, asynchronous and informal). Besides, an experiment run synchronously in a laboratory would still
 191 raise some control challenges: It might be distracting for some participants, or even induce some *follow*
 192 *the crowd* behavior, thus leading to people rushing to finish their tasks.

193 To regain some degree of control, participants ran all the tasks in a provided virtual machine available
 194 in our replication package (di Biase et al., 2018). Moreover, we recorded the screencast of the experiment,
 195 therefore not leaving space to misaligned results and mitigating issues of incorrect interpretation. Subjects
 196 were provided with instructions on how to use the virtual machine, but no time window was set.

197 4.4 Independent variable, group assignment, and duration

198 The independent variable of our study is change decomposition in pull requests. We split our subjects
 199 between a *control* group and a *treatment* group: The control group received one pull request containing
 200 a single commit with all the changes tied together; the treatment group received two pull requests with
 201 changes separated according to a logical decomposition.

202 Participants were randomly assigned to either the control group or the treatment using strata based on
 203 experience as developers and previous knowledge. Previous research has shown that these factors have
 204 an impact on review outcome (Rigby et al., 2012; Bosu et al., 2015): Developers who previously made
 205 changes to files to be reviewed had a higher proportion of useful comments.

206 All subjects were asked to run the experiment in a single session so that external distracting factors
 207 could be eliminated as much as possible. If a participant needed a pause, the pause is considered and
 208 excluded from the final result as we measure and monitor for periods of inactivity. We seek to reduce the
 209 impact of fatigue by limiting the expected time required for the experiment to an average of 60 minutes;

²Delft University of Technology Human Research Committee approved our study with IRB approval #578. University of Zurich authorized the research with IRB approval #2018-024.

210 this value is closer to the minimum rather than the median for similar experiments (Ko et al., 2015). As
211 stated before, though, we did not suggest or force any strict limit on the duration of the experiment to the
212 ends of replicating the code review informal scenario. No learning effect is present as every participant
213 runned the experiment only once.

214 **4.5 Pilot experiments**

215 We ran two pilot experiments to assess the settings. The first subject (a developer with 5 FTE³ years of
216 experience) took too long to complete the training and showed some issues with the virtual machine. Con-
217 sequently, we restructured the training phase addressing the potential environment issues in the material
218 provided to participants. The second subject (a MSc student with little experience) successfully completed
219 the experiment in 50 minutes with no issues. Both pilot experiments were executed asynchronously in the
220 same way as the actual experiment.

221 **4.6 Tasks of the experiment**

222 The participants were asked to conduct the following four tasks. Further details are available in the online
223 appendix (di Biase et al., 2018).

224 **1 - Preparing the environment.** Participants were given precise and detailed instructions on how to
225 set-up the environment for running the experiment. These entailed installing the virtual machine, setting
226 up the recording of the screen during the experiment, and troubleshooting common problems, such as
227 network or screen resolution issues.

228 **2 - Training the participants.** Before starting with the review phase, we first ensured that the participants
229 were sufficiently familiar with the system. It is likely that the participants had never seen the codebase
230 before: this situation would limit the realism of the subsequent review task.

231 To train our participants we asked subjects to implement three different features in the system:

- 232 1. Change the way the player moves on the board game, using different keys,
- 233 2. check if the game board has null squares (a board is made of multiple squares) and perform this
234 check when the board is created, and
- 235 3. implement a new enemy in the game, with similar artificial intelligence to another enemy but
236 different parameters.

237 This learning by doing approach is expected to have higher effectiveness than providing training material
238 to participants (Slavin, 1987). By definition, this approach is a method of instruction where the focus is
239 on the role of feedback in learning. The desired features required change across the system's codebase.
240 The third feature to be implemented targeted the classes and components of the game that would be object
241 of the review tasks. The choice of using this feature as the last one is to progressively increment the level
242 of difficulty.

243 No time window was given to participants, aiming for a more realistic scenario. As explicitly
244 mentioned in the provided instructions, participants were allowed to use any source for retrieving
245 information about something they did not know. This was permitted as the study does not want to assess
246 skills in implementing some functionality in a programming language. The only limitation is that the
247 participants must use the tools within the virtual machine.

248 The virtual machine provided the participants with the Eclipse Java IDE. The setup already had the
249 project imported in Eclipse's workspace. We used a plugin in Eclipse, WatchDog (Beller et al., 2015), to
250 monitor development activity. With this plugin, we measured how much time participants spent reading,
251 typing, or using the IDE. The purpose was to quantify the time to understand code among participants
252 and whether this relates to a different outcome in the following phases. Results for this phase are shown
253 in Figure 1, which contains boxplots depicting the data. It shows that there is no significant difference
254 between the two groups.

³A full-time employee (FTE) works the equivalent of 40 hours a week. We consider 1 FTE-year when a person has worked the equivalent of 40 hours a week for one year. For example, an individual working two years as a developer for 20 hours a week would have 1 FTE-year experience.

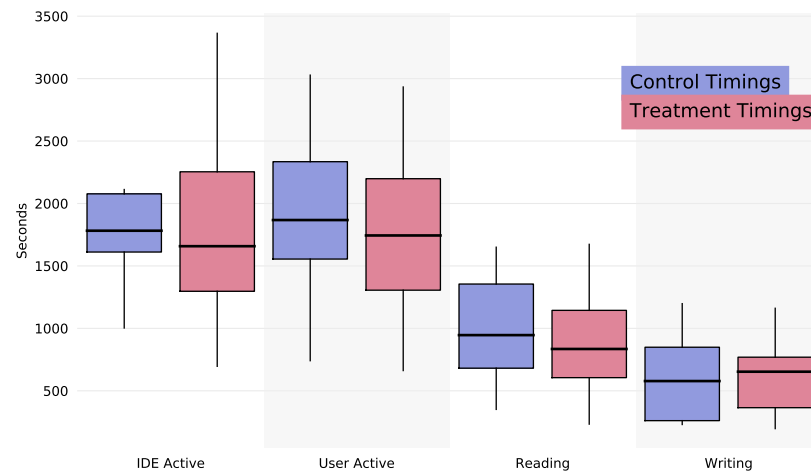


FIGURE 1. BOXPLOTS FOR TRAINING PHASE MEASUREMENTS. THE RESULTS HIGHLIGHT NO DIFFERENCES BETWEEN THE TWO GROUPS.

255 **3 - Perform code review on proposed change(s).** Participants were asked to review two changes
 256 made to the system:

- 257 1. the implementation of the artificial intelligence for one of the enemies
- 258 2. the refactoring of a method in all enemy classes (moving its logic to the parent class).

259 These changes can be inspected in the online appendix (di Biase et al., 2018) and have been chosen to
 260 meet the same criteria used by Herzig et al. (2016) when choosing tangled changes. Changes proposed
 261 can be classified as *refactoring* and *enhancement*. Previous literature gave insight as to how these two
 262 kinds of changes, when tangled together, are the hardest to review (Tao et al., 2012). Although recent
 263 research proposed a theory for the optimal ordering of code changes in a review (Baum et al., 2017),
 264 we used the default ordering and presentation provided by GitHub, because it is the de-facto standard.
 265 Changesets were included in pull requests on private GitHub repositories so that participants performed
 266 the tasks in a real-world review environment. Pull requests had identical descriptions for both the control
 267 and the treatment, with no additional information except their descriptive title. While research showed
 268 that a short description may lead to poor review participation (Thongtanunam et al., 2017), this does not
 269 apply to our experiment as there is no interaction among subjects.

270 Subjects were instructed to understand the change and check its functional correctness. We asked
 271 the participants to comment on the pull request(s) if they found any problem in the code, such as any
 272 functional error related to correctness and issues with code quality. The changes proposed had three
 273 different functional issues that were intentionally injected into the source code. Participants could see
 274 the source code of the whole project in case they needed more context, but only through GitHub's
 275 browser-based UI.

276 The size of the changeset was around 100 lines of code and it involved seven files. Gousios et al.
 277 showed that the number of total lines changed by pull requests is on average less than 500, with a median
 278 of 20 (Gousios et al., 2014). Thus, the number of lines of the changeset used in this study is between the
 279 median and the average.

280 **4 - Post-experiment questionnaire.** In the last phase participants were asked to answer the questions
 281 shown in Table 4. Questions Q1 to Q4 were about change-understanding, while Q5 to Q12 involved
 282 subjects' opinions about changeset comprehension and its correctness, rationale, understanding, etc. Q5 to
 283 Q12 were a summary of interesting aspects that developers need to grasp in a code change, as mentioned
 284 in the study of Tao et al. (2012). The answers must be provided in a Likert scale (Oppenheim, 2000)
 285 ranging from 'Strongly disagree' (1) to 'Strongly agree' (5).

286 4.7 Outcome measurements

287 **Effectiveness, false positives, suggested improvements.** Subjects were asked to comment a pull
 288 request in the pull request discussion or in-line comment in a commit belonging to that pull request.

289 The number of comments addressing functional issues was counted as the effectiveness. At the same
290 time, we also measured false positives (i.e., comments in pull request that do not address a real issue in
291 the code) and suggested improvements (i.e., remarks on other non-critical non-functional issues). We
292 distinguished suggested improvements and false positives from the comments that correctly addressed
293 an issue because the three functional defects were intentionally put in the source code. Comments that
294 did not directly and correctly tackle one of these three issues were classified either as false positives or
295 suggested improvements. They were classified by the first author by looking at the description provided
296 by the subject. A correctly identified issue needs to highlight the problem, and optionally provide a short
297 description.

298 **Time.** Having the screencast of the whole experiment, as well as using tools that give time measures,
299 we gathered the following measurements:

300 • Time for Task 2, in particular:

- 301 – total time Eclipse is [opened/active]
- 302 – total time the user is [active/reading/typing];

303 as collected by WatchDog (Section 4.6).

- 304 • Total net time for Task 3, defined as from when the subject opens a pull request until when (s)he
305 completes the review, purged of eventual breaks.
- 306 • Defect lifetime, defined as the period during which a defect continues to exist. It is measured
307 from the moment the subject opens a pull request to when (s)he writes a comment that correctly
308 identifies the issue. For the case of multiple comments on the same pull request, this is the time
309 between finishing with one defect and addressing the next. A similar measure was previously used
310 by Prechelt and Tichy (1998).

311 All the above measures are collected in seconds elapsed.

312 **Change-understanding.** In this experiment, change understanding was measured by means of a ques-
313 tionnaire submitted to participants post review activity, as mentioned in Task 4 in Section 4.6. Questions
314 are shown in Table 4 from Q1 to Q4. Its aim is to evaluate differences in change-understanding. A similar
315 technique was used by Binkley et al. (2013).

316 **Final Survey.** Lastly, participants were asked to give their opinion on statements targeting the perception
317 of correctness, understanding, rationale, logical partitioning of the changeset, difficulty in navigating the
318 changeset in the pull request, comprehensibility, and the structure of the changes. This phase, as well as
319 the previous one, was included in Task 4, corresponding to questions Q5 to Q12 (Table 4). Results were
320 given on a Likert scale from “Strongly disagree” (1) to “Strongly agree” (5) (Oppenheim, 2000), reported
321 as mean, median and standard deviation over the two groups, and tested for statistical significance with
322 the Mann-Whitney U-test.

323 4.8 Research method for RQ4

324 For our last research question, we aimed to build some initial hypothesis to explain the results from the
325 previous research questions. We sought what actions and patterns led a reviewer in finding an issue or
326 raising false positive, as well as other comments. This method was applied only to the review phase,
327 without analyzing actions and patterns concerning the training phase. The method to map actions to
328 concepts started by annotating the screencasts retrieved after the conclusion of the experimental phase.
329 Subjects performed a series of actions that defined and characterized both the outcome and the execution
330 of the review. The first author inserted notes regarding actions performed by participants to build a
331 knowledge base of steps (i.e., participant opens `fileName`, participant uses GitHub search box with the
332 `keyword`, etc.).

333 Using the methodology for qualitative content analysis delineated by Schreier (2013), we firstly
334 defined the coding frame. Our goal was to characterize the review activity based on patterns and behaviors.
335 As previous studies already tackled this problem and came up with reliable categories, we used the
336 investigations by Tao et al. (2012) and Sillito et al. (2006) as the base for our frame. We used the concepts
337 from Tao et al. (2012) regarding *Information needs for reasoning and assessing the change* and *Exploring*
338 *the context and impact of the change*, as well as the *Initial focus points* and *Building on initial focus points*
339 steps from Sillito et al. (2006).

340 To code the transcriptions, we used the deductive category application, resembling the data-driven
341 content analysis technique by Mayring (2000). We read the material transcribed, checking whether a
342 concept covers that action transcribed (e.g, participant opens file `fileName` so that (s)he is looking for
343 context). We grouped actions covered by the same concept (e.g, a participant opens three files, but always
344 for context purpose) and continued until we built a pattern that led to a specific outcome (i.e., addressing
345 a defect or a false positive). We split the patterns according to their concept ordering such that those that
346 led to more defects found or false positive issues were visible.

347 5 THREATS TO VALIDITY AND LIMITATIONS

348 **Internal validity** The sample size comprised in our experiment poses an inherent threat to the internal
349 validity of our experiment. Furthermore, the design and asynchronous execution of the experimental
350 phase creates uncertainty regarding possible external interactions. We could not control random changes
351 in the experimental setting, and this translates to possible disturbances coming from the surrounding
352 environment, that could cause skewed results.

353 Moreover, our experiment settings could not control if participants interacted among them, despite
354 participants did not have any information about each other.

355 Regarding the statistical regression (Wohlin et al., 2012), tests used in our study were not performed
356 with the Bonferroni correction, following the advice by Perneger: “*Adjustments are, at best, unnecessary*
357 *and, at worst, deleterious to sound statistical inference*” (Perneger, 1998).

358 **Construct validity** Relatively to the restricted generalizability across constructs (Wohlin et al., 2012),
359 in our experiment we uniquely aim to measure the values presented in Section 4.7. The treatment might
360 influence direct values we measure, but it could potentially cause negative effects on concepts that our
361 study does not capture. Additionally, we acknowledge threats regarding the time measures taken by
362 the first author regarding RQ2. Clearly, manual measures are suboptimal, that were adopted to avoid
363 participants having to perform such measures themselves.

364 When running an experiment, participants might try to guess what is the purpose of the experimentation
365 phase. Therefore, we could not control their behavior based on the guesses that either positively or
366 negatively affected the outcome.

367 Finally, we acknowledge threats to construct validity when designing the questionnaires used for RQ3,
368 despite designed using standard ways and scales (Oppenheim, 2000).

369 **External validity** Threats to external validity for this experiment concern the selection of participants to
370 the experimentation phase. Volunteers selected with convenience sampling could have an impact on the
371 generalizability of results, which we tried to mitigate by sampling multiple roles for the task. If the group
372 is very heterogeneous, there is a risk that the variation due to individual differences is larger than due to
373 the treatment (Cook and Campbell, 1979).

374 Furthermore, we acknowledge and discuss the possible threat regarding the system selection for the
375 experimental phase. Naturally, the system used is not fully representative of a real-world scenario. Our
376 choice, however, as explained in Section 4.1, aims to reduce the training phase effort required from
377 participants and to encourage the completion of the experiment.

378 Finally, our experiment was designed considering only a single programming language, using the
379 pull-based methodology to review and accept the changes proposed using GitHub as platform. Therefore,
380 threats for our experiment are related to mono-operation and mono-method bias (Wohlin et al., 2012).

381 6 RESULTS

382 RQ1. Effectiveness, false positives, and suggestions

383 For our first research question, descriptive statistics about results are shown in Table 2. It contains data
384 about effectiveness of participants (i.e., correct number of issues addressed), false positives, and number
385 of suggested improvements. Given the sample size, we applied a non-parametric test and performed a
386 Mann-Whitney U-test to test for differences between the control and the treatment group. This test, unlike
387 a t-test, does not require the assumption of a normal distribution of the samples. Results of the statistical
388 test are intended to be significant for a confidence level of 95%.

389 Results indicate a significant difference between the control and the treatment group regarding the
 390 number of false positives, with a *p-value* of 0.03. On the contrary, there is no statistically significant
 391 difference regarding the number of defects found (effectiveness) and number of suggested improvements.

TABLE 2. RQ1 - NUMBER OF DEFECTS FOUND (EFFECTIVENESS), FALSE POSITIVES AND SUGGESTED IMPROVEMENTS – STATISTICALLY SIGNIFICANT P-VALUES IN BOLDFACE.

	Group	# of subjects	Total	Median	Mean	σ	Confidence Interval 95%	p-value
Effectiveness (defects found)	Control	14	20	1.0	1.42	0.72	(0, 2.85)	0.6
	Treatment	14	17	1.0	1.21	0.77	(-0.30, 2.72)	
False Positives	Control	14	6	0	0.42	0.5	(-0.54, 1.40)	0.03
	Treatment	14	1	0	0.07	0.25	(-0.43, 0.57)	
Suggested Improvements	Control	14	7	0	0.5	0.62	(-1.22, 1.22)	0.4
	Treatment	14	19	1.0	1.36	1.84	(-2.17, 5.03)	

392 The example of a false positive is when one of the subjects of the control group writes: “*This doesn’t*
 393 *sound correct to me. Might want to fix the for, as the variable varName is never used*”. This is not a
 394 defect, as `varName` is used to check how many times the `for`-statement has to be executed, despite not
 395 being used inside the statement. This is also written in a code comment. Another false positive example is
 396 provided from a participant in the control group who, reading the *refactoring* proposed by the changeset
 397 under review, writes: “*The method methodName is used only in Class ClassName, so fix this*”. This
 398 is not a defect as the same `methodName` is used by the other classes in the hierarchy. As such, we can
 399 reject only the null hypothesis H_{0f} regarding the false positives, while we cannot provide statistically
 400 significant evidence about the other two variables tested in H_{0e} and H_{0c} .

401 The statistical significance alone for the false positives does not provide a measure to the actual
 402 impact of the treatment. To measure the effect size of the factor over the dependent variable we chose
 403 the Cliff’s Delta (Cliff, 1993), a non-parametric measure for effect size. The calculation is given by
 404 comparing each of the scores in one group to each of the scores in the other, with the following formula:
 405 $\delta = \frac{\#(x_1 > x_2) - \#(x_1 < x_2)}{n_1 n_2}$ where x_1, x_2 are values for the two groups and n_1, n_2 are their sample size. For data
 406 with skewed marginal distribution it is a more robust measure if compared to Cohen standardized effect
 407 size (Cohen, 1992). The computed value shows a positive (i.e., tangled pull requests lead to more false
 408 positives) effect size ($\delta = 0.36$), revealing a medium effect. The effect size is considered negligible for
 409 $|\delta| < 0.147$, small for $|\delta| < 0.33$, medium for $|\delta| < 0.474$, large otherwise (Romano et al., 2006).

Result 1: *Untangled pull requests (treatment) lead to fewer false positives with a statistically significant, medium size effect.*

410 Given the presence of *suggested improvements* in our results, we found that the control group writes
 411 in total seven, while the participants in the treatment write nineteen. This difference is interesting, calling
 412 for further classification of the suggestions. For the control group, participants wrote respectively three
 413 improvements regarding *code readability*, two concerning *functional checks*, one regarding *understanding*
 414 of source code and one regarding other code issues. For the treatment group, we classified five suggestions
 415 for *code readability*, eight for *functional checks* and seven for *maintainability*. Although subjects have been
 416 explicitly given the goal to find and comment exclusively functional issues (Section 4.6), they wrote these
 417 suggestions spontaneously. The suggested improvements are included in the online appendix (di Biase
 418 et al., 2018) along with their classification.

419 RQ2. Review time and defect lifetime

420 To answer RQ2, we measured and analyzed the time subjects took to review the pull requests, as well
 421 as the amount of time they used to fix each of the issues present. Descriptive statistics about results for
 422 our second research question are shown in Table 3. It contains data about the time participants used to
 423 review the patch, completed by the measurements of how much they took to fix respectively two of the
 424 three issues present in the changeset. All measures are in seconds. We exclude data relatively to the third
 425 defect as only one participant detected it. To perform the data analysis, we used the same statistical means
 426 described for the previous research question. When computing the review net time used by the subjects,
 427 results show an insignificant difference, thus we are not able to reject null-hypothesis H_{0t1} . This indicates

428 that the average case of the treatment group takes the same time to deliver the review, despite having two
 429 pull requests to deal with instead of one. However, analyzing results regarding the *defect lifetime* we also
 430 see no significant difference and cannot reject H_{02} . Data show that the mean time to address the first
 431 issue is about 14% faster in the treatment group if compared with the control. This is because subjects
 432 have to deal with less code that concerns a single concept, rather than having to extrapolate context
 433 information from a tangled change. At the same time the treatment group is taking longer (median) to
 434 address the second defect. We believe that this is due to the presence of two pull requests, and as such, the
 435 context switch has an overhead effect on that. From the screencast recordings we found no reviewer using
 436 multi-screen setup, therefore subjects had to close a pull-request and then review the next, where they
 437 need to gain knowledge on different code changes.

Result 2: *Our experiment was not able to provide evidence for a difference in net review time between untangled pull requests (treatment) and the tangled one (control); this despite the additional overhead of dealing with two separate pull requests in the treatment group.*

TABLE 3. RQ2 - REVIEW TIME, FIRST AND SECOND DEFECT LIFETIME - MEASUREMENTS IN SECONDS ELAPSED

	Group	# of subjects	Median	Mean	σ	Conf. Interval 95%	p-value
Review net time	Control (Tangled changes)	14	831	853	385	(99, 1607)	0.66
	Treatment (Untangled changes)	14	759	802	337	(140, 1463)	
1st defect lifetime	Control	11	304	349	174	(8, 691)	0.79
	Treatment	11	297	301	109	(86, 516)	
2nd defect lifetime	Control	6	222	263	149	(-28, 555)	0.17
	Treatment	6	375	388	122	(148, 657)	

438 RQ3. Understanding The Change's Rationale

439 For our third research question, we seek to measure whether subjects are affected by the dependent
 440 variable in their understanding of the rationale of the change. Rationale understanding questions are Q1 to
 441 Q4 (Table 4) and Figure 2 reports the results. Higher scores for Q1, Q2, and Q4 mean better understanding,
 442 whereas for Q3 a lower score signifies a correct understanding. As for the previous research questions,
 443 we test our hypothesis with a non-parametrical statistical test. Given the result we cannot reject the null
 444 hypothesis H_{0u} of tangled pull requests reducing change understanding. Participants are in fact able to
 445 answer the questions correctly, independent of their experimental group.

446 After the review, our experimentation also provided a final survey (Q5 to Q12 in Table 4) that
 447 participants filled in at the end. Results shown in Figure 2 indicate that subjects judge equally the
 448 changeset (Q5), found no difficulty in understanding the changeset (Q6), agree on having understood the
 449 rationale behind the changeset (Q7). This results shows that our experiment cannot provide evidence of
 450 differences in change understanding between the two groups.

451 Participants did not find the changeset hard to navigate (Q9), and believe that the changeset was
 452 comprehensible (Q11). Answers to questions Q9 and Q11 are surprising to us, as we would expect
 453 dissimilar results for code navigation and comprehension. In fact, change decomposition should allow
 454 subjects to navigate code easier, as well as improve source comprehension.

455 On the other hand, subjects from the control and treatment group judge differently when asked if the
 456 changeset was partitioned according to a logical separation of concerns (Q8), if the relationships among
 457 the changes were well structured (Q10) and if the changes were spanning too many features (Q12). These
 458 answers are in line with what we would expect, given the different structure of the code to be reviewed.
 459 The answers are different with a statistical significance for Q8, Q10 and Q12.

Result 3: *Our experiment was not able to provide evidence of a difference in understanding the rationale of the changeset between the experimental groups. Subjects reviewing the untangled pull requests (treatment) recognize the benefits of untangled pull requests, as they evaluate the changeset as being (1) better divided according to a logical separation of concerns, (2) better structured, and (3) not spanning too many features.*

TABLE 4. RQ3 - POST-EXPERIMENT QUESTIONNAIRE.
QUESTIONS WITH * HAVE $p < 0.05$

Questions on understanding the rationale of the changeset	
The purpose of this changeset entails ...	
Q1	... changing a method for the enemy AI
Q2	... the refactoring of some methods
Q3	... changing the game UI panel
Q4	... changing some method signature
Questions on participant's perception on the changeset	
Q5	The changeset was functionally correct
Q6	I found no difficulty in understanding the changeset
Q7	The rationale of this changeset was perfectly clear
Q8 *	The changeset [showed] a logical separation of concerns
Q9	Navigating the changeset was hard
Q10 *	The relations among the changes were well structured
Q11	The changeset was comprehensible
Q12 *	Code changes were spanning too many features

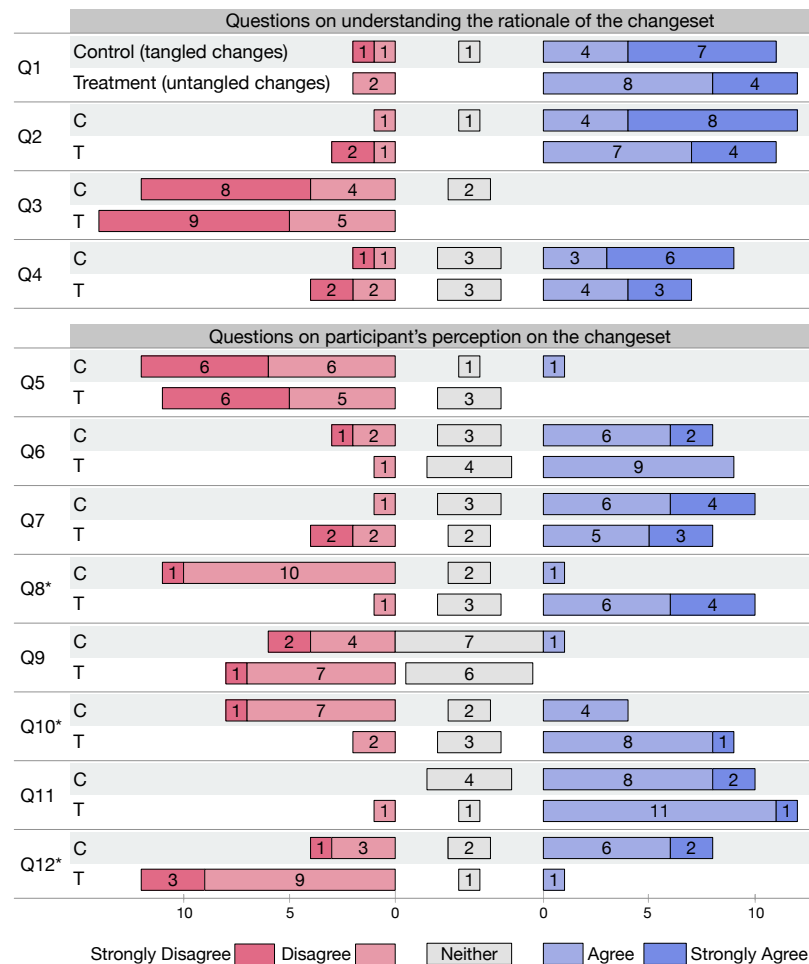


FIGURE 2. RQ3 - ANSWERS TO QUESTIONS IN TABLE 4.
QUESTIONS WITH * HAVE $p < 0.05$

TABLE 5. RQ4 - CONCEPTS FROM LITERATURE AND THEIR MAPPED KEYWORD

Concept	Mapped keyword
What is the rationale behind this code change? (Tao et al., 2012)	Rationale
Is this change correct? Does it work as expected? (Tao et al., 2012)	Correctness
Who references the changed classes/methods/fields? (Tao et al., 2012)	Context
How does the caller method adapt to the change of its callees? (Tao et al., 2012)	Caller/Callee
Is there a precedent or exemplar for this? (Sillito et al., 2006)	Similar/Precedent

460 **RQ4. Tangled vs. Untangled review patterns**

461 For our last research question, we seek to identify differences in patterns and features during review,
 462 and their association to quantitative results. We derived such patterns from Tao et al. (2012) and Sillito
 463 et al. (2006). These two studies are relevant as they investigated the role of understanding code during
 464 the software development process. Tao et al. (2012) laid out a series of information needs derived from
 465 state-of-the-art research in software engineering, while Sillito et al. (2006) focused on questions asked by
 466 professional experienced developers while working on implementing a change. The mapping found in the
 467 screencasts is shown in Table 5.

468 Table 6 contains the qualitative characterization, ordered by the sum of defects found. Values in each
 469 row correspond to how many times a participant in either group used that pattern to address a defect or
 470 point to a false positive.

TABLE 6. RQ4 - PATTERNS IN REVIEW TO ADDRESS A DEFECT OR LEADING TO A FALSE POSITIVE

ID	Pattern			Control		Treatment	
	1 st concept	2 nd concept	3 rd concept	Defect	FP	Defect	FP
P1	Rationale	Correctness		8	3	4	0
P2	Rationale	Context	Correctness	4	0	5	0
P3	Context	Rationale	Correctness	3	2	3	0
P4	Context	Correctness	Caller/Callee	1	0	2	0
P5	Context	Correctness		2	1	0	0
P6	Correctness	Context		0	0	2	0
P7	Rationale	Correctness	Context	0	0	1	0
P8	Correctness	Context	Caller/Callee	1	0	0	0
P9	Correctness	Context	Similar/Precedent	1	0	0	1

471 Results indicate that pattern *P1* is the one that led to most issues being addressed in the control group
 472 (eight), but at the same time is the most imprecise one (three false positives). We conjecture that this is
 473 related to the lack of context-seeking concept. Patterns *P1* and *P3* have most false positives addressed in
 474 the control group. In the treatment group, pattern *P2* led to more issues being addressed (five), followed
 475 by the previously mentioned *P1* (four).

476 Analyzing the transcribed screencasts, we note an overall trend of reviewing code changes in the
 477 control group, exploring the changeset using less context exploration than in the treatment. Among the
 478 participants belonging to the treatment, we witnessed a much more structured way of conducting the
 479 review. The overall behavior is that of getting the context of the single change, looking for the files
 480 involved, called, or referenced by the changeset, in order to grasp the rationale. All of the subjects except
 481 three repeated this step multiple times to explore a chain of method calls, or to seek for more context
 482 in that same file opening it in GitHub. We consider this the main reason to explain that untangled pull
 483 requests lead to more precise (fewer false positives) results.

Result 4: *Our experiment revealed that review patterns for untangled pull requests (treatment) show more context-seeking steps, in which the participants open more referenced/related classes to review the changeset.*

484 7 DISCUSSION

485 7.1 Implications for Researchers

486 In past studies, researchers found that developers call for tool and research support for decomposing a
487 composite change (Tao et al., 2012). For this reason, we were surprised that our experiment was not able
488 to highlight differences in terms of reviewers' effectiveness (number of defects found) and reviewers'
489 understanding of the change rationale, when the subjects were presented with smaller, self-contained
490 changes. Further research with additional participants is needed to corroborate our findings.

491 If we exclude latent problems with the experiment design that we did not account for, this result may
492 indicate that reviewers are still able to conduct their work properly, even when presented with tangled
493 changes. However, the results may change in different contexts. For example, the cognitive load for
494 reviewers may be higher with tangled changes, thus the negative effects in terms of effectiveness could be
495 visible when a reviewer has to assess a large number of changes every day, as it happens with integrators
496 of popular projects in GitHub (Gousios et al., 2015). Moreover, the changes we considered are of average
497 size and difficulty, yet results may be impacted by larger changes and/or more complex tasks. Finally,
498 participants were not core developers of the considered software system; it is possible that core developers
499 would be more surprised by tangled changes, find them more convoluted or less "natural," thus rejecting
500 them (Hellendoorn et al., 2015). We did not investigate these scenarios further, but studies can be designed
501 and carried out to determine whether and how these aspects influence the results of the code review effort.

502 Given the remarks and comments of professional developers on tangled changes (Tao et al., 2012),
503 we were also surprised that the experiment did not highlight any differences in the net review time
504 between the treatment groups. Barring experimental design issues, this result can be explained by
505 the additional context switch, which does not happen in the tangled pull request (control) because the
506 changes are done in the same files. An alternative explanation could be that the reviewers with the
507 untangled pull requests (treatment) spent more time "wondering around" and pinpointing small issues
508 because they found the important defects quicker; this would be in line with the cognitive bias known as
509 Parkinson's Law (Parkinson and Osborn, 1957) (all the available time is consumed). However, time to
510 find the first and second defects (3) is the same for both experimental groups thus voiding this hypothesis.
511 Moreover, similarly to us, Tao and Kim also did not find difference with respect to time to completion in
512 their preliminary user study (Tao and Kim, 2015). Further studies should be designed to replicate our
513 experiment and, if results are confirmed, to derive a theory on why there is no reduction in review time.

514 Our initial hypothesis on why time does not decrease with untangled code changes is that reviewers
515 of untangled changes (control) may be more willing to build a more appropriate context for the change.
516 This behavior seems to be backed up by our qualitative analysis (Section 6), through the context-seeking
517 actions that we witnessed for the treatment group. If our hypothesis is not refuted by further research,
518 this could indicate that untangled changes may lead to a more thorough low-level understanding of the
519 codebase. Despite we did not measure this in the current study, it may explain the lower number of
520 false positives with untangled changes. Finally, untangled changes may lead to better transfer of code
521 knowledge, one of the positive effects of code review (Bacchelli and Bird, 2013).

522 7.2 Recommendation for Practitioners

523 Our experiment is not able to show no negative effects when changes are presented as separate, untangled
524 changesets, despite the fact that reviewers have to deal with two pull requests instead of one, with
525 the subsequent added overhead and a more prominent context switch. With untangled changesets, our
526 experiment highlighted an increased number of suggested improvements, more context-seeking actions
527 (which, it is reasonable to assume, increase the knowledge transfer created by the review), and a lower
528 number of wrongly reported issues.

529 For the aforementioned reasons, we support the recommendation that change authors prepare self-
530 contained, untangled changeset when they need a review. In fact, untangled changesets are not detrimental
531 to code review (despite the overhead of having more pull-requests to review), but we found evidence of
532 positive effects. We expect the untangling of code changes to be minimal in terms of cognitive effort and
533 time for the author. This practice, in fact, is supported by answers Q8, Q10, Q12 to the questionnaire and
534 by comments written by reviewers in the control group (i.e., "Please make different commit for these two
535 features", "I would prefer having two pull requests instead of one if you are fixing two issues").

536 8 CONCLUSION

537 The goal of the study presented in this paper is to investigate the effects of change decomposition on mod-
538 ern code review (Cohen, 2010), particularly in the context of the pull-based development model (Gousios
539 et al., 2014).

540 We involved 28 subjects, who performed a review of pull request(s) pertaining to (1) a refactoring
541 and (2) the addition of a new feature in a Java system. The control group received a single pull request
542 with both changes tangled together, while the treatment group received two pull requests (one per type of
543 change). We compared control and treatment groups in terms of effectiveness (number of defects found),
544 number of false positives (wrongly reported issues), number of suggested improvements, time to complete
545 the review(s), and level of understanding the rationale of the change. Our investigation involved also a
546 qualitative analysis of the review performed by subjects involved in our study.

547 Our results suggests that untangled changes (treatment group) lead to:

- 548 1. fewer reported false positives defects,
- 549 2. more suggested improvements for the changeset,
- 550 3. same time to review (despite the overhead of two different pull requests),
- 551 4. same level of understanding the rationale behind the change,
- 552 5. and more context-seeking patterns during review.

553 Our results support the case that committing changes belonging to different concepts separately should
554 be an adopted best practice in contemporary software engineering. In fact, untangled changes not only
555 reduce the noise for subsequent data analyses (Herzig et al., 2016), but also support the tasks of the
556 developers in charge of reviewing the changes by increasing context-seeking patterns.

557 ACKNOWLEDGMENTS

558 The authors would like to thank all participants of the experiment and the pilot. We furthermore thank the
559 fellow researchers who gave critical suggestion to help strengthening the methodology of our study.

560 REFERENCES

- 561 Bacchelli, A. and Bird, C. (2013). Expectations, outcomes, and challenges of modern code review. In
562 *Proceedings of the 2013 International Conference on Software Engineering, ICSE '13*, pages 712–721,
563 Piscataway, NJ, USA. IEEE Press.
- 564 Barnett, M., Bird, C., Brunet, J., and Lahiri, S. (2015). Helping developers help themselves: Automatic
565 decomposition of code review changesets. In *Proceedings of the 37th International Conference on*
566 *Software Engineering - Volume 1, ICSE '15*, pages 134–144, Piscataway, NJ, USA. IEEE Press.
- 567 Baum, T., Schneider, K., and Bacchelli, A. (2017). On the optimal order of reading source code changes
568 for review. In *2017 IEEE International Conference on Software Maintenance and Evolution (ICSME)*,
569 pages 329–340.
- 570 Beller, M., Gousios, G., Panichella, A., and Zaidman, A. (2015). When, how, and why developers (do
571 not) test in their ides. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software*
572 *Engineering, ESEC/FSE 2015*, pages 179–190, New York, NY, USA. ACM.
- 573 Binkley, D., Davis, M., Lawrie, D., Maletic, J., Morrell, C., and Sharif, B. (2013). The impact of identifier
574 style on effort and comprehension. *Empirical Software Engineering*, 18(2):219–276.
- 575 Bosu, A., Greiler, M., and Bird, C. (2015). Characteristics of useful code reviews: An empirical study
576 at microsoft. In *2015 IEEE/ACM 12th Working Conference on Mining Software Repositories*, pages
577 146–156.
- 578 Cliff, N. (1993). Dominance statistics: Ordinal analyses to answer ordinal questions. *Psychological*
579 *Bulletin*, 114(3):494.
- 580 Cohen, J. (1992). Statistical power analysis. *Current directions in psychological science*, 1(3):98–101.
- 581 Cohen, J. (2010). Modern code review. In Oram, A. and Wilson, G., editors, *Making Software*, chapter 18,
582 pages 329–338. O'Reilly.
- 583 Cook, T. D. and Campbell, D. T. (1979). *Quasi-experimentation: Design and analysis for field settings*,
584 volume 3. Rand McNally Chicago.

- 585 di Biase, M., Bruntink, M., and Bacchelli, A. (2016). A Security Perspective on Code Review: The
586 Case of Chromium. In *16th IEEE International Working Conference on Source Code Analysis and*
587 *Manipulation, SCAM 2016, Raleigh, NC, USA, October 2-3, 2016*, pages 21–30. IEEE Press.
- 588 di Biase, M., Bruntink, M., van Deursen, A., and Bacchelli, A. (2018). The Effects
589 of Change Decomposition on Code Review - A Controlled Experiment - Online appendix.
590 <https://data.4tu.nl/repository/uuid:826f7051-35f6-4696-b648-8e56d3ea5931>.
- 591 Dias, M., Bacchelli, A., Gousios, G., Cassou, D., and Ducasse, S. (2015). Untangling fine-grained code
592 changes. In *Proceedings of the 22nd International Conference on Software Analysis, Evolution, and*
593 *Reengineering, SANER 2015*, pages 341–350. IEEE Computer Society.
- 594 Gousios, G., Pinzger, M., and van Deursen, A. (2014). An exploratory study of the pull-based software
595 development model. *Proceedings of the 36th International Conference on Software Engineering - ICSE*
596 *2014*, (May 2014):345–355.
- 597 Gousios, G., Zaidman, A., Storey, M., and van Deursen, A. (2015). Work practices and challenges
598 in pull-based development: The integrator’s perspective. In *Proceedings of the 37th International*
599 *Conference on Software Engineering - Volume 1, ICSE ’15*, pages 358–368, Piscataway, NJ, USA.
600 IEEE Press.
- 601 Hellendoorn, V. J., Devanbu, P. T., and Bacchelli, A. (2015). Will they like this? Evaluating code
602 contributions with language models. In *Proceedings of the 12th Working Conference on Mining*
603 *Software Repositories*, pages 157–167. IEEE Press.
- 604 Herzig, K., Just, S., and Zeller, A. (2016). The impact of tangled code changes on defect prediction
605 models. *Empirical Software Engineering*, 21(2):303–336.
- 606 Herzig, K. and Zeller, A. (2013). The impact of tangled code changes. In *Mining Software Repositories*
607 *(MSR), 2013 10th IEEE Working Conference on*, pages 121–130. IEEE.
- 608 Kirinuki, H., Higo, Y., Hotta, K., and Kusumoto, S. (2014). Hey! are you committing tangled changes?
609 In *Proceedings of the 22nd International Conference on Program Comprehension, ICPC 2014*, pages
610 262–265, New York, NY, USA. ACM.
- 611 Ko, A., LaToza, T., and Burnett, M. (2015). A practical guide to controlled experiments of software
612 engineering tools with human participants. *Empirical Software Engineering*, 20(1):110–141.
- 613 Mayring, P. (2000). Qualitative content analysis. *Forum: Qualitative Social Research*.
- 614 McIntosh, S., Kamei, Y., Adams, B., and Hassan, A. (2014). The impact of code review coverage and code
615 review participation on software quality: A case study of the qt, vtk, and itk projects. In *Proceedings of*
616 *the 11th Working Conference on Mining Software Repositories, MSR 2014*, pages 192–201, New York,
617 NY, USA. ACM.
- 618 McIntosh, S., Kamei, Y., Adams, B., and Hassan, A. E. (2016). An empirical study of the impact of
619 modern code review practices on software quality. *Empirical Software Engineering*, 21(5):2146–2189.
- 620 Morales, R., McIntosh, S., and Khomh, F. (2015). Do code review practices impact design quality? A
621 case study of the Qt, Vtk, and Itk projects. In *Proceedings of the 22nd International Conference on*
622 *Software Analysis, Evolution and Reengineering, SANER 2015*, pages 171–180. IEEE.
- 623 Murphy-Hill, E., Parnin, C., and Black, A. (2012). How we refactor, and how we know it. *IEEE*
624 *Transactions on Software Engineering*, 38(1):5–18.
- 625 Oppenheim, A. (2000). *Questionnaire design, interviewing and attitude measurement*. Bloomsbury
626 Publishing.
- 627 Parkinson, C. N. and Osborn, R. C. (1957). *Parkinson’s law, and other studies in administration*,
628 volume 24. Houghton Mifflin Boston.
- 629 Perneger, T. V. (1998). What’s wrong with bonferroni adjustments. *British Medical Journal*,
630 316(7139):1236.
- 631 Prechelt, L. and Tichy, W. (1998). A controlled experiment to assess the benefits of procedure argument
632 type checking. *IEEE Transactions on Software Engineering*, 24(4):302–312.
- 633 Rigby, P., Cleary, B., Painchaud, F., Storey, M., and German, D. (2012). Contemporary peer review in
634 action: Lessons from open source development. *IEEE software*, 29(6):56–61.
- 635 Rigby, P., German, D., Cowen, L., and Storey, M. (2014). Peer Review on Open-Source Software Projects.
636 *ACM Transactions on Software Engineering and Methodology*, 23(4):1–33.
- 637 Rigby, P. C. and Bird, C. (2013). Convergent contemporary software peer review practices. In *Proceedings*
638 *of the 2013 9th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2013*, pages 202–212.
639 ACM.

- 640 Romano, J., Kromrey, J., Coraggio, J., and Skowronek, J. (2006). Appropriate statistics for ordinal level
641 data: Should we really be using t-test and cohen'sd for evaluating group differences on the nsse and
642 other surveys. In *Annual Meeting of the Florida Association of Institutional Research*, pages 1–33.
- 643 Schreier, M. (2013). Qualitative content analysis. In *The SAGE handbook of qualitative data analysis*,
644 pages 170–183. SAGE.
- 645 Sharif, B., Falcone, M., and Maletic, J. (2012). An eye-tracking study on the role of scan time in finding
646 source code defects. In *Proceedings of the Symposium on Eye Tracking Research and Applications*,
647 pages 381–384. ACM.
- 648 Sillito, J., Murphy, G., and De Volder, K. (2006). Questions programmers ask during software evolution
649 tasks. In *Proceedings of the 14th ACM SIGSOFT international symposium on Foundations of software
650 engineering*, pages 23–34. ACM.
- 651 Slavin, R. (1987). Mastery learning reconsidered. *Review of educational research*, 57(2):175–213.
- 652 Tao, Y., Dang, Y., Xie, T., Zhang, D., and Kim, S. (2012). How do software engineers understand code
653 changes?: An exploratory study in industry. In *Proceedings of the ACM SIGSOFT 20th International
654 Symposium on the Foundations of Software Engineering, FSE '12*, pages 1–11, New York, NY, USA.
655 ACM.
- 656 Tao, Y. and Kim, S. (2015). Partitioning composite code changes to facilitate code review. In *Proceedings
657 of the 12th Working Conference on Mining Software Repositories*, pages 180–190. IEEE.
- 658 Thongtanunam, P., McIntosh, S., Hassan, A. E., and Iida, H. (2017). Review participation in modern code
659 review. *Empirical Software Engineering*, 22(2):768–817.
- 660 Uwano, H., Nakamura, M., Monden, A., and Matsumoto, K. (2006). Analyzing individual performance
661 of source code review using reviewers' eye movement. In *Proceedings of the 2006 symposium on Eye
662 tracking research & applications*, pages 133–140. ACM.
- 663 Wohlin, C., Runeson, P., Höst, M., Ohlsson, M., Regnell, B., and Wesslén, A. (2012). *Experimentation in
664 software engineering*. Springer Science & Business Media.