# Investigating the Impact of Continuous Integration Practices on the Productivity and Quality of Open-Source Projects

Jadson Santos
Federal University of Rio Grande do Norte
Natal, Brazil
jadson.santos@ufrn.br

Daniel Alencar da Costa
University of Otago
Dunedin, New Zealand
danielcalencar@otago.ac.nz

Uirá Kulesza
Federal University of Rio Grande do Norte
Natal, Brazil
uira@dimap.ufrn.br

## ABSTRACT

**Background:** Much research has been conducted to investigate the impact of Continuous Integration (CI) on the productivity and quality of open-source projects. Most of studies have analyzed the impact of adopting a CI server service (e.g, Travis-CI) but did not analyze CI sub-practices. **Aims:** We aim to evaluate the impact of five CI sub-practices with respect to the productivity and quality of GitHub open-source projects. **Method:** We collect CI sub-practices of 90 relevant open-source projects for a period of 2 years. We use regression models to analyze whether projects upholding the CI sub-practices are more productive and/or generate fewer bugs. We also perform a qualitative document analysis to understand whether CI best practices are related to a higher quality of projects. **Results:** Our findings reveal a correlation between the Build Activity and Commit Activity sub-practices and the number of merged pull requests. We also observe a correlation between the Build Activity, Build Health and Time to Fix Broken Builds sub-practices and number of bug-related issues. The qualitative analysis reveals that projects with the best values for CI sub-practices face fewer CI-related problems compared to projects that exhibit the worst values for CI sub-practices. **Conclusions:** We recommend that projects should strive to uphold the several CI sub-practices as they can impact in the productivity and quality of projects.

## CCS CONCEPTS

• **Software and its engineering** → **Agile software development**; Software evolution.

## KEYWORDS

Continuous Integration, CI Sub-Practices, CI Maturity, Software Productivity, Software Quality

## 1 INTRODUCTION

Continuous Integration (CI) is a software development practice that aims to reduce the complexity of integrating code produced by different developers. It can also decrease the risks [9] of software projects by promoting the automated build and testing of their code. The automated process of CI can bring benefits such as the improvement in the delivery time of a new software release as well as decreasing the number of errors potentially generated by manual tasks. The adoption of CI requires the adoption of a set of sub-practices [9, 11, 27], such as to perform code commits frequently, build the software frequently and develop and perform automated tests.

Over the last years, researchers have empirically studied CI from different perspectives [25] [20] [32] [18] [24]. Recent research works have investigated the impact of CI on productivity and quality of software projects [30] [2]. Vasilescu et al. [30] analyzed a dataset of 246 GitHub open-source projects that adopted Travis-CI at some point in their lifetime. They found that after adopting Travis-CI, core developers in those projects have significantly increased the number of merged pull requests (PRs) as well as discovered more bugs-related issue reports (IRs). Helis et al. [2] conducted an empirical study that analyzed 162,653 PRs of 87 GitHub projects to evaluate the impact of adopting CI on the delivery time of merged PRs. They found a large increase in the number of submitted, merged and delivered PRs per release after CI was adopted. These empirical studies derive their results based on projects that have adopted a CI service, such as Travis-CI.[1] Nevertheless, Felidré et al. [10] observed that many projects that adopt a CI service, do not necessarily follow fundamental CI practices. They refer to this problem as *Continuous Integration Theater* [10]. CI Theater occurs when a project has adopted a CI service to automate their build and testing process but, on the other hand, do not dedicate much attention to other existing CI sub-practices, such as frequent code commits, frequent builds and automated tests.

In this context, this paper investigates the impact of different CI sub-practices on the productivity and quality outcomes of open-source projects. Similar to previous works [30] [2], we define productivity in terms of merged PRs and quality in terms of closed bug-related IRs. Beyond selecting projects solely on the basis that they have used a CI service, we extract 5 CI sub-practices (build duration, build activity, build health, time to fix a broken build,

---

[1]https://travis-ci.org/

and commit activity) and investigate the correlation of those sub-practices with the productivity (merged pull requests) and quality (closed bug-related issue reports) of a software project. Furthermore, we perform a qualitative analysis to investigate how much these CI sub-practices can reflect the maturity of the projects concerning the adoption of CI. We use a document analysis approach to analyze and compare the comments of PRs from two groups of projects - the *Best CI Sub-Practices Group* and the *Worst CI Sub-Practices Group*. Thereby, we asked the following research questions:

- **RQ1: Which CI sub-practices contribute to higher productivity outcomes?** We apply linear regression models to analyze the influence of CI sub-practices in the productivity of projects. Can a specific CI sub-practice increase the productivity of a project?
- **RQ2: Which CI sub-practice contribute to higher quality outcomes?** We apply linear regression models to analyze the influence of CI sub-practices in the quality of projects. Can a specific CI sub-practice increase the quality of a project?
- **RQ3: Can consistently applied CI sub-practices be an indicator of process quality?** The goal of this research question is to verify whether the projects that best follow the 5 CI sub-practices are more mature regarding CI. Do developers face fewer problems in the development process than the developers of projects that have the worst values for CI sub-practices?

Based on our results, we find a positive correlation between 2 CI sub-practices and a increase in the number of merged PRs. We find also a positive correlation between frequency of builds and the increase of bugs-related IRs and our model shows that projects with poor build health tend to generate more bug-related IRs. We find also evidence that projects that uphold the CI sub-practices can indicate a better quality in the development process, while they face fewer problems regarding CI, compared to projects that do not uphold the CI sub-practices.

**Paper organization**. The rest of this paper is organized as follows: In Section 2, we explain the design of our empirical study, describing the motivation and methodology used in each RQ. In Section 3 we present the results of our empirical study, while we discuss the impact of these results in Section 4. In Section 5 we discuss the threats to the validity. In Section 6, we discuss the related work. Finally, we concluded the paper in Section 7.

## 2 STUDY SETTINGS & RESEARCH QUESTIONS

In this section, we first describe how we selected our projects and then we explain the motivation and methodology of each research question (RQ).

### 2.1 Projects Selection

To select the projects of our study, we searched for popular GitHub projects that use Travis-CI as their CI service. The projects must also have had a substantial number of builds, pull requests (PRs), issue reports (IRs) and commits. Lastly, our studied projects would need to be active (in terms of builds, PRs, IRs and commits) to perform our investigations. Similar to other studies [2, 30, 31],

we selected projects using Travis-CI as Travis-CI is one of the most popular CI services on GitHub [30]. Differently from other CI servers (e.g., Jenkins), Travis-CI provides the entire build history of a project. Lastly, although Github-Actions has become popular, its addition to GitHub is still fairly recent, [2] which could hinder our goal of finding a 2-year history of builds in our projects.

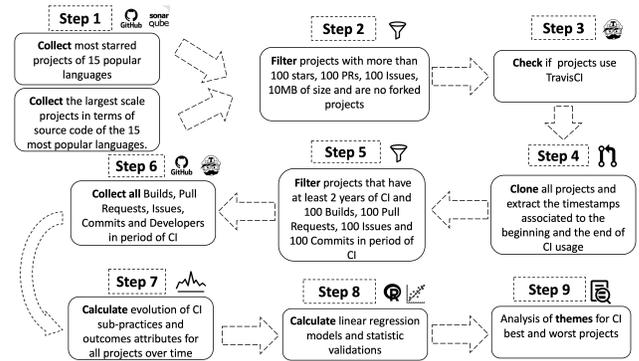Figure 1 provides an overview of all steps involved in our project selection approach.



**Figure 1: An overview of our methodology**

Our selection strategy was inspired by Zhao et al.'s work [34]. In **Step 1**, we collected the most starred repositories for a set of 15 popular programming languages on GitHub: Java, JavaScript, C#, Python, PHP, TypeScript, C, Go, C++, Kotlin, Ruby, Rust, Swift, Scala and Objective-C using the GitHub Search API,[3] yielding 11,671 projects. Next, to increase the number of projects in our study, we searched for additional projects on SonarQube using its web API.[4] We only kept SonarQube projects that were also available on GitHub, obtaining an additional set of 16,212 projects. In total, we collected 27,883 projects from the two searches.

In **Step 2**, our goal was to eliminate non-relevant projects (e.g., toy projects). Considering the initial 27,883 projects, we first verified whether they survived the following thresholds: more than 100 stars, 100 PRs, 100 Issues, 10MB of size and were not forked projects. Figure 2 depicts the selection process for steps 2, 3, 4 and 5.

The number "100" in our thresholds was inspired by previous studies [2, 30]. Additionally, similarly to prior research works [1, 7, 8], we used *10MB of size* as a threshold to eliminate starred projects that do not represent meaningful projects. For example, the hello-world[5] GitHub project has 1.7K stars, more than 700 Issues and 300 PRs, consisting only of a single README.md file for a demo.

After applying these first filters, 3,143 projects remained. From this set of projects, in **Step 3**, we checked whether they have used Travis-CI. We did so by using a URL pattern match [31] extracted from the project's name on GitHub. 2,029 projects remained.

In **Step 4**, similarly to Zhao et al.'s work [34], we cloned the GitHub repositories of all 2,029 projects. We identified the main Travis-CI branch of the repositories by identifying the earliest

---

[2]https://github.blog/2019-08-08-github-actions-now-supports-ci-cd/

[3]https://docs.github.com/en/rest/reference/search#search-repositories

[4]https://community.sonarsource.com/t/list-of-all-public-projects-on-sonarcloud-using-api/33551

[5]https://github.com/octocat/hello-world

Figure 2: Overview of filtering non-relevant projects



Figure 3: Periods of Analysis

commit to the `.travis.yml` configuration file. We used the timestamp of the first commit to the `.travis.yml` file as the starting date of the CI usage within the projects. Afterwards, we retrieved the timestamp of the last build of the projects from their TRAVIS-CI service. We considered the timestamp of the last build as the ending date of CI usage in the projects.

We filter for projects that have at least **2 years** of activity on TRAVIS-CI in **Step 5**. We obtain a total of 1,751 projects. Next, we performed a new filtering step, selecting projects with at least 100 Builds, 100 Pull Requests, 100 Issues and 100 Commits, but now considering the period of 2 years of TRAVIS-CI activity. We obtained a total of 776 relevant projects.

In **Step 6**, we collected information from PRs, IRs, commits, and builds from both GITHUB and TRAVIS-CI within the period of 2 years of CI usage. We collected only closed PRs and IRs. Regarding to closed PRs, we only analyzed merged PRs, i.e., PRs that have a non-null value within the *merged_at* field as retrieved by the GITHUB API. In total, considering the 776 relevant projects, we collected 403,403 PRs, 983,460 IRs, 394,063 builds and 1,486,429 commits.

In the **Step 7**, we analyze the evolution of CI sub-practices (explained in Section 2.2) within all the 776 relevant projects over time. To do so, we divided the project history into time intervals. Similar to Zhao et al.'s work [34] and Vasilescu et al.'s work [30], we use an interval of **1 month** to segment our projects' history into periods of analysis. Table 1 shows the computed attributes with their respective definitions. Instead of using the period of 24 months centered around the adoption of TRAVIS-CI (as in prior studies [34]), to pursue our goals, we analyze the latest 24 months of CI usage in our projects. Our goal is to evaluate CI projects that have better recent values for the CI sub-practices. The latest 24 months is the period where projects should have had the most stable values for the CI sub-practices. Figure 3 illustrates the information collected for each period of analysis.

After segmenting our projects' history in periods, we noted that several projects had several periods without activity, i.e., 0 builds and 0 commits in a specific month. To only consider projects that remained active during the period of analysis, we filtered for projects that had at least 20 periods (80% of the total number of periods) with at least 1 build and 1 commit. Thus, we have a final set of **90 relevant and active projects** used in our study.
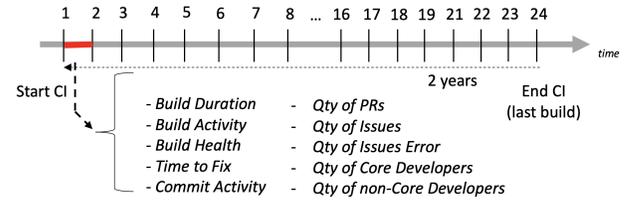
In the last two steps of our method, we built our linear regression models **(Step 8)** to verify the potential influence of each CI sub-practice in the *productivity* and *quality* [30] outcomes of our projects. In **Step 9**, we perform a qualitative analysis to verify whether projects that uphold the CI sub-practices also show signs of quality in their development process. All collected information as well as the database and scripts used in this study are available in our online appendix.[6]

## 2.2 CI sub-practices

Existing research still lacks a well-defined set of criteria to analyze the potential impact of CI on software development [25]. For instance, several studies consider the use of a CI service (e.g., TRAVIS-CI) as the sole criterion to indicate whether projects use CI or not [25]. However, Duvall et al.[9] and Fowler [11] have defined a set of CI sub-practices that, combined together, can be used to determine whether a project uses CI or not. Additionally, Felidré et al. [10] revealed that many CI projects do not necessarily follow CI sub-practices consistently (e.g., infrequent commits are frequent). Therefore, in our work, instead of investigating whether the usage of TRAVIS-CI can be associated with better quality and productivity outcomes [34], we study **5 CI sub-practices** in a more separate fashion to understand their potential contribution to quality and productivity outcomes.

- **Build Duration** [10, 13] measures the duration of the build *(build finished at timestamp - build started at timestamp)*. The build duration was retrieved directly from TRAVIS-CI. To fit our models, we used the median build duration per period of analysis (1 month).
- **Build Activity** [9] is a unit interval (i.e., a closed interval [0,1]) representing the rate of builds across days, i.e, if builds were made every day in the period of analysis, the value would be 1. If builds were made in half of the days, the value would be 0.5. If there were no builds, the value would be 0.
- **Build Health** [25] is a unit interval representing the rate of build failures across days. If there were build failures every day, the value would be 0. if there were no build failures, the value would be 1.
- **Time to Fix a Broken Build** [10] consists of the median time in a period (1 month) that builds remained broken. When a build breaks, we compute the time in seconds until the build returns to the *"passed"* status. If the analysis period ends and the build did not return to the *"passed"* status, we consider the time since it was broken until the end of the

**Table 1: Projects outcomes attributes collected for each of 24 months of study**

| Atributes | Definition |
| --- | --- |
| qty_pull_request | number of PRs in the period |
| qty_pull_request_core | Number of PRs in the period associated with core developers |
| qty_pull_request_non_core | Number of PRs in the period associated with non-core developers |
| qty_developers_prs | Number of distinct PR authors in period |
| qty_core_developers_prs | Number of distinct PR core developers authors in period |
| qty_non_core_developers_prs | Number of distinct PR non-core developers authors in period |
| qty_issue | number of Issues in the period |
| qty_issue_core | Number of Issues in the period associated with core developers |
| qty_issue_non_core | Number of Issues in the period associated with non-core developers |
| qty_issue_error | Number of Issues in the period associated with error labels |
| qty_issue_error_core | Number of Issues in the period associated with error labels and core developers |
| qty_issue_error_non_core | Number of Issues in the period associated with error labels and non-core developers |
| qty_developers_issues | number of distinct Issues authors in period |
| qty_core_developers_issues | Number of distinct Issues core developers authors in period |
| qty_non_core_developers_issues | Number of distinct Issues non-core developers authors in period |

period. When a period has no broken builds, the value would be 0.

- **Commit Activity** [10] is a unit interval representing the rate of commits across days. If commits were made every day in a period (1 month), the value would be 1. If commits were made in half of days the value would be 0.5. If there were no commits in a period the value would be 0.

In this paper, we focus on investigating the build and commit related sub-practices. These sub-practices were chosen because they cover most of the practices defined by Duvall et al [9] and Fowler [11] as well as they cover most of the sub-practices analyzed by Felidré et al [10]. The "Code Coverage" sub-practice analyzed by Felidré et al. [10] was initially considered in this study. However, we found few projects that collect and store coverage data for a long period of time on a centralized platform like SonarQube.[7]

## 2.3 Research Questions

*RQ1: Which CI sub-practices contribute to higher productivity outcomes?*

**Motivation:** Vasilescu et al. [30] observed that CI is associated with higher productivity outcomes (i.e., number of merged PRs). However, we still do not know whether the increase in productivity is associated to only certain sub-practices of CI (as opposed to all of them). Better understanding which CI sub-practices have stronger associations with productivity outcomes (i.e., number of merged PRs [30]) may help us to optimize the usage of CI.

**Approach:** We fit linear regression models (Ordinary Least Squares, OLS [17]), to find associations between the number of merged PRs (Table 1) and the 5 CI sub-practices. We would expect that the number of merged PRs should be higher in projects where CI sub-practices are employed consistently.

Because the number of developers in a project can impact the number of merged PRs, we control the influence of the number developers in our models. We create another model where we divide the number of merged PRs by the number of active developers in each period of analysis. Inspired by the findings of Vasilescu et al.'s work [30], we decided to divide our analysis in two groups of developers: (i) core developers and (ii) non-core developers. Our goal was to minimize the effects of the growth of internal and external contributors on our models.

To identify core developers, we combined strategies used by Vasilescu et al. [30] and Poncin et al. [23]. A core developer is a developer who: (i) *"either had write access to a project's code repository or had closed issues and pull requests submitted by others"* [30] or (ii) *"has been involved in the project for a relatively long period of time and has more revisions in the version control system than average"* [23].

As it is common for open-source developers to use different aliases when performing contributions, we used a strategy to identify core developers from author information. Differently from Zhao et al. [34] who stated: *"we used heuristics that match first and last names, email prefixes, and email domains"*, after collecting such information from GitHub, we noticed that emails were not informed by many developers. Thus, we only used the first and last name of developers. To identify whether two developers have similar names, we applied a combination of 2 algorithms: (i) Jaro Winkler Similarity [12] and (ii) Levenshtein Distance [6]. We empirically found that, with Jaro Index $\geq$ **0.85** and Levenshtein Distance < **5**, we could accurately compare developer names. To validate our method, we implemented tests that randomly match the names of the developers of our projects. We then manually verified whether the matches were made correctly before performing our analyses. Regarding the strategy used by stated by Poncin et al. [23], we considered at least 12 months of commits in a repository as the *"relatively long period of time"* to identify core developers.

---

[7]https://www.sonarqube.org/

Before fitting the OLS model to our data, we verified the correlation among the independent variables of the model (i.e., the 5 CI sub-practices). First, we applied Shapiro-Wilk's tests [29] to verify the normality of our data. Having observed that none of the independent variables follow a normal distribution, we chose the Kendall non-parametric test [28] to measure the correlation between independent variables.

### RQ2: Which CI sub-practice contribute to higher quality outcomes?

**_Motivation:_** In RQ2, we aim to better understand which CI sub-practices have a stronger impact on the quality outcomes of a project. Although Vasilescu et al. [30] observed that the usage of CI likely holds a relationship with quality outcomes (i.e., number of bug-related IRs), we still do not know whether only certain CI practices have a significant association with quality outcomes. Similar to RQ1, better understanding which CI sub-practices have stronger associations with quality outcomes (i.e., number bug-related IRs [30]) may help ups to optimize the usage of CI.

**_Approach:_** Similar to Vasilescu et al. [30], we manually checked our projects to understand how they used tags to indicate the presence of bugs. Similarly to Vasilescu et al.'s work [30], we selected projects that have at least 75% of their IRs tagged with labels in our dataset. For the sake of understanding how labels were used, we apply this selection criterion in the broader group of 776 relevant projects. We found a total of 121 projects that were manually checked. We found 2 main tags that were used to indicate the presence of bugs (which were not present in the original list [30]): (i) "crash" and (ii) "regression". For example, the _mui/material-ui_ project[8] describes the _"regression"_ tag as _"A bug, but worse"_. Thus, we added these 2 new tags to the original list of tags defined by Vasilescu et al. [30]. The resulting list of labels indicating bug-related IRs was: _"defect", "error", "bug", "issue", "mistake", "incorrect", "fault" , "flaw", "crash" and "regression"_. We then searched in the group of 90 relevant and active projects (i.e., the final group of projects where we applied our analysis) for these words in IRs' labels. To do so, we performed lowercasing and applied the Porter Stemming Algorithm.[9]

Our data for quality outcomes, i.e., number of bug-related IR, is similar in many aspects to the data used in Vasilescu et al.'s work [30]: (i) most of our variables are counts (e.g., qty of bugs reported, qty of issues, qty of developers, etc.); (ii) some variables are over-dispersed, i.e., the variance is much larger than the mean; and (iii) some response variables present an excess number of zeros. Because of these characteristics, in RQ2, we applied zero-inflated negative binomial regression models (ZINB) [15].

To ensure that using ZINB models was an appropriate choice, we first compared poisson regression models with negative binomial regression models. This comparison shows that the residuals are more spread out for the poisson, furthermore, we applied the likelihood ratio test,[10] which shows a statistically significant difference in favor to negative binomial model. Thus a negative binomial offers

a significantly better fit to our data compared to a poisson regression model. Next, we fit a zero-inflated negative binomial model with the ordinary negative binomial model using Vuong's test[11] of a non-nested model. Our test provides evidence of the superiority of a zero-inflated model over an ordinary model (AIC-corrected -3.2992634, p-value < 0.05 ).

### RQ3: Can consistently applied CI sub-practices be an indicator of process quality?

**_Motivation:_** Felidré et al. [10] observed that not all projects follow all CI sub-practices consistently. For instance, while some projects may have acceptable build duration, they may not perform frequent commits. It is important to know whether projects that perform most of CI sub-practices in a consistent manner are also associated with a higher quality in their development process. Given that it would be challenging to quantitatively measure "quality of development process", we perform a qualitative analysis in RQ3.

**_Approach:_** To answer this research question, we separate the projects into 2 groups: (i) projects that follow all the best CI sub-practices analyzed in this study; and (ii) projects that do not effectively follow all the best CI sub-practices. Figure 4 shows the steps to generate these two groups of projects.
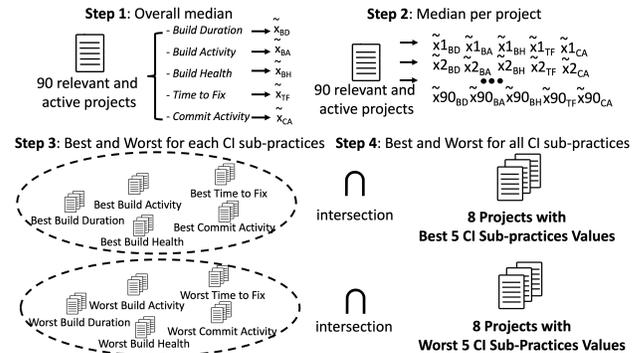


**Figure 4: Best and Worst CI sub-practices groups selection**

First, we computed the overall median of each of the 5 CI sub-practices for the 90 relevant and active projects in our dataset. Next, we computed the median of CI sub-practices per project. We then selected projects that had the best and worst values compared to the overall median for each of the 5 CI sub-practices. Lastly, we selected projects in which their medians were all (i) best or (ii) worst the overall median for all 5 CI sub-practices. This segmentation resulted in a group of 8 projects whose medians were above the overall median for all CI sub-practices (the _Best CI Practices_ group) and another group of equally 8 projects whose medians were below the overall median for all CI sub-practices (the _Worst CI Practices_ group). We emphasize that the values of the sub-practices have different meanings, depending on the way in which they were calculated. For example, in the case of the Build Duration sub-practice, the best values are lower values. In the case of the Commit Activity sub-practice, the best values are the higher values. These

---

[8]https://github.com/mui/material-ui/labels?page=5&sort=name-asc
[9]https://tartarus.org/martin/PorterStemmer/
[10]https://www.statisticshowto.com/likelihood-ratio-tests/

[11]https://www.rdocumentation.org/packages/pscl/versions/1.5.5/topics/vuong

differences were considered when splitting our projects into two groups.

Afterwards, we applied a qualitative approach known as Document Analysis [4]. This analysis consists of coding documents into themes. We used the comments of the PRs of the projects as inputs to this analysis. We obtained representative samples of PRs to perform our analyses[19]. The *Best CI Practices* group has 1,560 PRs, and the *Worst CI Practices* group have 989 PRs. Considering a confidence level of 95% and a confidence interval of 10%, we created two representative samples containing 91 PRs and 88 PRs, respectively.

We randomly selected 91 and 88 PRs under the condition that each project would contribute with the same number of PRs inside representative sample. We rounded up the number of PRs and randomly selected 12 PRs for each of the 8 projects within each group. For each of the selected PRs, all comments were collected and save in 2 different files with randomly generated names, which were then used in the thematic analysis. This was important to not identify to which group of projects the PRs belonged. As this is a subjective process, the thematic analysis was performed by two coders–a main coder and a reviewer–to minimize the bias. After the thematic analysis was done, we revealed the data to identify group to which the comments belonged.

## 3 RESULTS

In this section, we present the results of each research question.

## RQ1: Which CI sub-practices contribute to higher productivity outcomes?

Before fitting our regression models, we applied the Kendall test to verify the correlation among the independent variables. The Null Hypothesis is that variables are uncorrelated. If p-value < 0.05, we reject the Null Hypothesis. The $\tau$ value shows how strong the correlation is. *Builds Activity* and *Commit Activity* practices were correlated with $\tau = 0.5109093$ (p-value < 2.2e-16), which is to be expected, since in a CI environment, more commits generate more builds. *Builds Activity* and *Time to Fix a Broken Build* also demonstrated a correlation with $\tau 0.2374206$ (p-value < 2.2e-16). With more builds, projects spend more time to fix them. *Build Health* and *Time to Fix a Broken Build* revealed a negative correlation with $tau - 0.4995392 (p - value < 2.2e - 16)$, the longer it takes for projects to fix their builds, the less the *Build Health*. Given that our $\tau$ coefficients were not close to 1 (the perfect correlation), we decided to keep all independent variables in our model.

Our regression model obtained a R-squared of **0.353** (R-squared adj. **0.347**). It is challenging to determine what is a "good" R-squared value as it will depend on the goal of the research. For example, if the main goal is prediction, R-squared values should be high (e.g., 0.7 to 0.9) [5]. However, lower R-squared values (e.g., around 0.20) may also provide important insights in psychology or social sciences [3]. Given that our goal is to better understand associations between CI sub-practices and quality/productivity outcomes, we believe our R-squares value are acceptable. Figure 5 shows the relationships between the 5 CI sub-practices and the number of merged PRs (as fit by our OLS model).
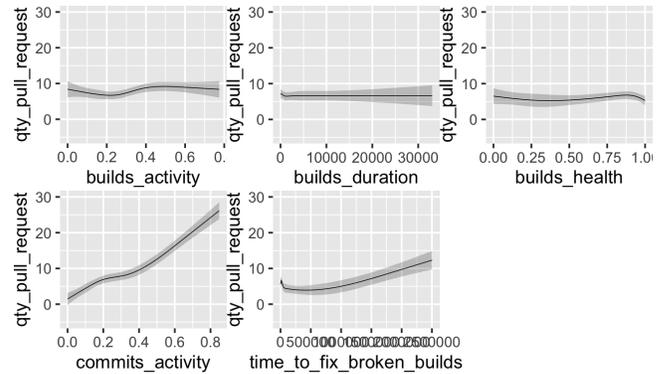


**Figure 5: Regression model between the 5 CI sub-practices and the number of merged PRs**

To verify the strength of the associations between the CI sub-practices and the response variable, we computed Spearman's tests [26]. This test measures the strength and direction of associations between two ranked variables by the correlation coefficient *rho*. Our tests reveal that *Commit Activity*, with $\rho = 0.388$, and *Build Activity*, with $\rho = 0.274$, are the two most influential sub-practices when it comes to explain the number of merged PRs. Projects with high *Commit Activity* and *Build Activity* values, tend to be more productive, i.e., more PRs are merged. Although these results are expected, it is interesting to note that sub-practices such as *Time to Fix a Broken Build* do not have much influence on productivity. For example, this may be due to the fact that merges may still occur when the build is broken [14]. Table 2 shows the values of our Spearman's tests.

**Table 2: Spearman tests: CI sub-practices vs. merged PRs**

| CI sub-practice | rho2 |
|---|---|
| Builds Duration | 0.007 |
| Builds Activity | **0.274** |
| Builds Health | 0.000 |
| Time to Fix Broken Builds | 0.014 |
| Commits Activity | **0.388** |

We fitted another model for the same 5 CI sub-practices in which we divided the number of merged PRs by the number of core and non-core developers in each period of the analysis. We recalculate the regression model to normalize the number of merged PRs by the number of developers and to minimize the effects of the number of developers on the productivity of the projects. This new models obtained lower R-square values (**0.296** for core developers and **0.210** for non-core developers) but still revealed significant correlations with CI sub-practices. This result suggests that the number of core and non-core developers may also influence the number of merged PRs but they do not eliminate the potential effect of CI sub-practices on the productivity of outcomes. According to our model, the number of non-core developers has more influence on productivity outcomes than core developers.

> ***Commit Activity*** and ***Build Activity*** are strongly associated with the number of merged Pull Requests. When considering the number of core and non-core developers, the correlations becomes weaker but they are still significant.

## RQ2: Which CI sub-practice contribute to higher quality outcomes?

Table 3 shows our zero-inflated negative binomial models. We highlight the CI sub-practices that obtain statistical significance.

**Table 3: ZINB model for bug-related IRs**

| Count model coefficients (negbin with log link) | | | | |
|---|---|---|---|---|
| | Est. Std. | Error | z value | Pr(>\|z\|) |
| (Intercept) | -0.1316 | 0.2530 | -0.520 | 0.6027 |
| **builds activity** | **1.9846** | **0.4163** | **4.767** | **1.87e-06 *** |
| commits activity | 0.7552 | 0.3941 | 1.916 | 0.0553 . |
| **builds health** | **-0.7466** | **0.2710** | **-2.754** | **0.0058 ** |
| **time to fix** | **-2.4175** | **0.5730** | **-4.219** | **2.46e-05 *** |
| builds duration | 3.4774 | 1.8004 | 1.931 | 0.0534 . |
| Log(theta) | -1.5574 | 0.0603 | -25.819 | < 2e-16 *** |
| Zero-inflation model coefficients (binomial with logit link) | | | | |
| | Est. Std. | Error | z value | Pr(>\|z\|) |
| (Intercept) | 0.1061 | 0.8241 | 0.129 | 0.8975 |
| builds activity | -10.4734 | 10.0846 | -1.039 | 0.2990 |
| **commits activity** | **-11.1607** | **5.5541** | **-2.009** | **0.0445 *** |
| builds health | -11.4912 | 7.0698 | -1.625 | 0.1041 |
| **time to fix** | **8.0927** | **3.9048** | **2.073** | **0.0382 *** |
| builds duration | 75.1436 | 43.4784 | 1.728 | 0.0839 . |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

The *Build Activity* sub-practice has a significant correlation with the number of bug-related IRs (Table 3). This result may be because more builds may be a consequence of the necessity of fixing more bug-related IRs. Another explanation is that more builds may indicate projects that interact more with their end-users who, in turn, report more bug-related IRs. On the other hand, *Build Health* has an inverse correlation with bug-related IRs. The better the *Build Health* the lower the number of bug-related IRs. This result suggests that projects with less broken builds are less likely to deliver bugs to the end users. Surprisingly, the *Time to Fix a Broken Build* sub-practice has an inverse correlation with bug-related IRs. Projects taking a longer time to fix broken builds may have fewer bug-related IRs.

In our dataset, taking a longer time to fix broken builds is correlated with a reduced *Build Activity* (see Section RQ1). We conjecture that fewer builds may be associated with a low activity within the projects (e.g., receiving fewer updates or contributions from external users). Consequently, less bug-related IRs could have been reported in projects with a longer *Time to Fix Broken Builds*. In fact, the *Time to Fix Broken Builds* is inversely correlated with *Build Health* (R-squared of 0.2324), meaning that projects that take a longer time to fix a broken build also tend to have more build breakages in general, which may make the project less attractive to users.

For example, project stability (e.g., less breakages) may be related to project attractiveness [33].

Our model also found (in the coefficient parts of the Zero-inflation model) that an increase in Commit Activity is associated with a higher probability of generating bug-related IRs. The Time to Fix Broken Builds increases the odds of being in the group without bug-related IRs. This result reinforces the Count model coefficient part of the ZINB.

> ***Build Activity*** has a significant correlation with the the number of bug-related issue reports. Higher values of ***Commit Activity*** also increase the probability of generating bug-related IRs. Maintaining a good build health seems to decrease the number of bug-related issue reports. However, surprisingly, large values of ***Time to Fix a Broken Build*** are associated with less bug-related IRs.

## RQ3: Can consistently applied CI sub-practices be an indicator of process quality?

Figure 6 illustrates the themes that emerged from our qualitative analysis for the *Worst CI projects*, whereas Figure 7 shows the themes that emerged from our qualitative analysis for the *Best CI projects*. We grouped the themes into 5 CI-related high-level themes. The high-level themes and their definitions are shown in Table 4.
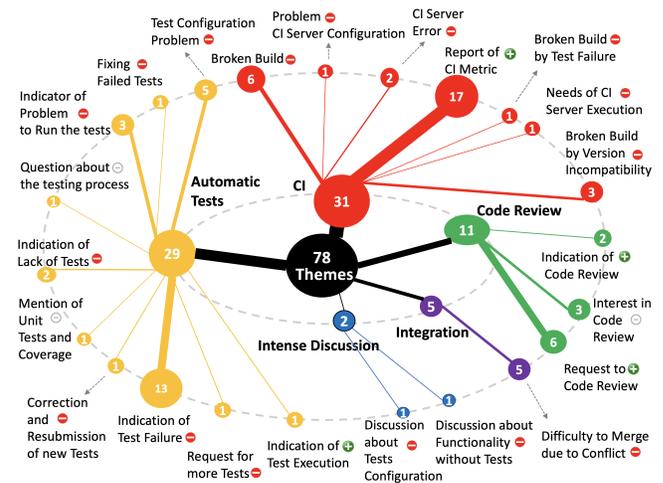


**Figure 6: Themes that emerged from the Worst CI projects**

The number in the center of each node indicates the number of occurrences of each theme. The thickness of the edges and nodes is based on the number of times the theme emerged during the thematic analysis. For example, the *"Discussion about tests configuration"* theme has only one occurrence in the Worst CI projects. The themes were also classified as having a positive (➕), neutral (➖) or negative (⛔) meaning in the project under analysis.

As an example of a PR comment that we consider as having a negative meaning, we quote: *"Can you please resolve merging conflicts with changelog and then I will proceed with merging"*. This comment was associated with the **"Difficulty to Merge due to**
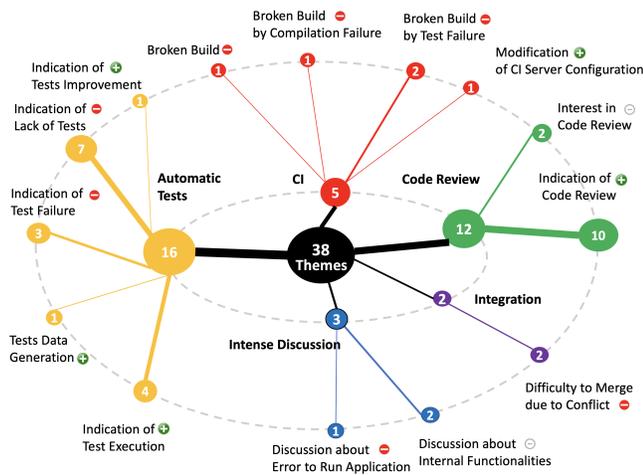
**Figure 7: Themes that emerged from the Best CI projects**

**Table 4: CI high-level Themes emerged from the Document Analysis**

| CI high-level themes | Description |
| --- | --- |
| Automatic Tests | Denotes evidence of automated tests |
| Integration | Denotes integration issues |
| Intense Discussion | More than 10 comments in a PR |
| Code Review | Indicates that the source code has been reviewed |
| CI | Issues, practices or tools directly related to CI |

**Conflict"** theme under the **"Integration"** high-level theme. Building the software constantly and making frequent commits should reduce conflicts and, consequently, integration effort.

As an example of a PR comment that we consider as having a positive meaning, we quote: *"fast yet thorough review! I've implemented all of your suggestions"*. This comment was associated with theme **Indication of Code Review** of the **Code Review** high-level theme. According to Steve McConnell [21]: *"Code review is one of the most effective ways to find bugs and improve code quality"*, which is a practice recommended in CI. If we find a comment indicating that a code review was accomplished on the project at some point during the implementation of a PR, then we consider that the revised code tends to have fewer bugs [21]. Thus, **Indication of Code Review** is classified as having a positive meaning. In the opposite case, we quote: *"Thanks for looking into the test failure!"* This comment was associated with theme **Indication of Test Failure** of the **Automatic Tests** high-level theme. According to the grey literature: [22]: *"Developers must write automated builds with tests that pass 100% of the time, and not get or commit broken code from/to the version control repository"*. If a PR comment indicates that "tests failed" it has a negative aspect for the project. Thus, we classified it as having a negative meaning. We explain our themes more thoroughly below.

**Automated Tests.** We obtained 29 citations related to *Automated Tests* in the Worst CI projects, of which 13 citations were related to the **Indication of Test Failure (negative)** theme. As examples of comments indicating test failures in the Worst CI projects, we quote the following comments: *"This works great in my tests with Postman against my rails app but the specs I added fail in a mysterious way."* In another PR, a developers comments: *"On the tests: The first two fail because 'PORO::Employee#positions' will return an empty array by default"*. We also observed occurrences of specific tests failures in comments such as: *"tests that are failing are being run with Python 2.7"*. Regarding the Best CI projects, we only found 3 citations related to the **Indication of Test Failure** theme. As an example, we quote: *"Init test for hash type 6211. Could not load test module: m06211.pm"* and *"the test fail is unrelated to your changes."*.

**Code Review.** We found only two citations of the **Indication of Code Review (positive)** theme in the Worst CI projects, whereas we found 10 citations of the same theme in the Best CI projects. As examples, we quote the following excerpts for the Worst CI projects: *"I couldn't figure it out. Otherwise it is ready for review"*. For the Best CI projects we quote: *"I started reviewing some of the code "*, *"I see more changes that are needed but it looks already very good. Good job. "*, *"another question: did you @xxxxxx try with much longer salt size"*, and *"Processed your review comments..."*.

**Integration.** We found five citations related to the **Difficulty to Merge due to Conflict (negative)** theme in the Worst CI projects, whereas we found two citations of the same theme in the Best CI projects. As examples, we quote the following excerpts for the Worst CI projects: *"I rebased this branch on master resolved conflicts and merged then force pushed"*, *"Didn't need to bump the examples version but it doesn't really matter. Another conflict popped up however."*. For the Best CI projects we quote: *"Looks good now but I can't merge because of docs/credits.txt conflict."*.

**Intense Discussion.** We found two main themes in the Worst CI projects: **Discussion about tests configuration (negative)** and **Discussion on functionalities without tests (negative)**. On the other hand, we found one negative theme in the Best CI projects: **Discussion on errors when running the application**. As examples, we quote the following excerpts for the Worst CI projects: *"This behavior has no tests and is not documented."* and *"Let's say hypothetically I moved the tests out of the package. On a scale from -1 to +1 how annoyed would you be to re-sync this again?."*. For the Best CI projects we quote: *"I believe the line 'texexpand: include titlepag.tex failed.' could be the cause for the index.html being malformed."*

**Continuous Integration (CI).** The Worst CI projects obtained 7 different themes directly related to the **CI** high-level theme, whereas the Best CI projects obtained only four **CI** themes. We thought this result was counter-intuitive. However, from the 7 different themes in the Worst CI projects, 6 of them were negative. Some examples of negative comments were: *"The build is failing – would you take a look to see if you can get it working"*, *"Travis CI is being derpy again so whenever that decides to work."*, *"Removing the archive makes our builds reproducible without setting that environment variable."*, *"Travis CI failed only ruby 2.2.10. I will rebuild Travis CI."*, *"I see that this is failing to pass the tests in Travis - I'll have a go at fixing that up now"*, *"Closed and reopened to wake up Travis CI."*. Regarding the Best CI Projects, from the four themes that we found, three themes were negative and one theme was positive.

The positive theme was the **Modification of CI Server Configuration** theme: *"I modified the Travis configuration anyways to make builds 30 seconds faster."*

The theme **Reporting of a CI Metric** has a greater occurrence in the Worst CI projects and no occurrence at all in the Best CI projects. Indeed, this happened because the *"graphite-project/graphite-web"* project (classified under the Worst CI projects) used the Code-Cov tool[12], which automatically reports coverage values in the comments of each PR. Despite being a positive practice, unfortunately, the usage of such a tool does match the behavior of the project, i.e., the project obtained lower values within all CI sub-practices.

In general, the Worst CI projects have more citations of CI-related themes (78). However, most of these comments (60%) were classified as having a negative meaning. If we do not consider the **Reporting of a CI Metric** theme, the percentage would rise to (77%). The Best CI projects obtained 38 CI-related themes with 17 (44%) of the themes classified as having a negative meaning. For all CI high-level themes, the results were favorable for Best CI projects. Therefore, our data and results suggest that consistently employing CI sub-practices may indeed be associated with an improved development process.

> **Our results reveal that projects with best CI sub-practices values have fewer citations related to the CI high-level themes and these citations have a more positive meaning. It suggests that higher values in the 5 CI sub-practices indeed can indicate a better quality in the development process.**

## 4  DISCUSSION

The main goal of our work was to highlight the importance of studying CI by means of its sub-practices. We were inspired by Felidré et al.'s work on the *CI Theater* phenomenon [10], which revealed that many projects do not follow CI practices. Therefore, instead of considering the usage of a CI server as the sole criterion to assume that projects use CI (which is a problem also highlighted by Soares et al.'s work [25]), we investigated the potential impact of CI on quality and productivity outcomes by means of 5 CI sub-practices.

All analyzed projects in our study adopt a CI service to automate the build and test execution process. However, only the most active projects — the projects making frequent commits and builds — are more productive. Furthermore, the increase in commit and build frequency does not come at the cost of quality, meaning that frequent commits and builds are more effective to generate more delivery value into the project. Additionally, CI can make developers more confident as they witness automated builds and tests providing continuous feedback. Thus pull requests are likely only created when they are in a near-approval state.

On the other hand, more active projects (in terms of builds and commits) tend to have more bug-related issue reports open. This can be an indication that such projects have a higher demand and receive more requests from external users or internally from the developers. Our results suggests that maintaining a good build health has a correlation with code quality. A better build health

is associated with more stable versions of the project, generating fewer errors for users and consequently fewer bug-related issue reports.

Our work also found that projects with the best values for CI sub-practices generated fewer CI-related themes when we analyzed their development processes. In fact, the CI-related themes found in projects adhering to the sub-practices usually have a more positive connotation within the project. It reveals that projects that follow CI best sub-practices can be a reflection of a shared understanding between developers on how CI must be used, allowing developers to focus more on the business rules instead of discussing how CI should be used. Perhaps, if a project can document its "Continuous Integration Principles" it would make it easier for new contributors to understand the philosophy of the project when it comes to CI usage.

Next we discuss the implications of our work for researchers and practitioners.

**Implications for researchers:**  Our study shows that we need to consider, measure and evaluate the different CI sub-practices to better determine that CI is being used in a project. Especially in RQ3, when we segment our projects into different groups (i.e., Best and Worst CI practices), we have evidence that the developers' experience with CI is quite different. Future research should consider evaluating several CI sub-practices before evaluating the potential influence of CI adoption on software development aspects.

**Implications for practitioners:**  The results of RQ1 and RQ2 indicate that certain benefits of CI (e.g., more merged PRs) are not equally shared by all of the CI sub-practices. For instance, our models in RQ1 reveal that practices such as *Build Health* and *Build Duration* do not have much influence on the number of merged PRs. Our results suggest that, instead of fully adopting CI from the beginning, a staged adoption (e.g., only certain practices) may be more wise given the goals of the project. For example, starting with frequent commits and builds Moreover, given that *culture* is an important factor for the successful adoption of new methodologies [16], the staged CI adoption can be a CI-enabler. With respect to quality, our recommendation is to strive to maintain a sound build health, i.e., to generate fewer broken builds. Build health obtained a significant and inverse correlation with the number of bug-related issue reports.

## 5  THREATS TO VALIDITY

In this section, we discuss the threats to the validity of our study.

**Construct Validity Threats:**  Although Travis-CI is widely used in the context of GitHub open-source projects, existing projects may use other CI services or servers, which might have led us to discard relevant projects.

With respect to the analysis regarding core developers, there is no unique identifier to compare GitHub developers (as mentioned by Zhao et al. [34]). Due to this issue, we applied a heuristic based on the comparison of developer names to collect the data associated with a specific core or non-core developer. This heuristic may not be the most accurate, generating false negatives or positives in our data. We used sub-samples of the investigated projects to assess the quality of the extracted information related to developers.

---

[12]https://about.codecov.io/product/feature/pull-request-comments/

**Internal Validity Threats:** We collected data regarding productivity and quality outcomes and CI sub-practices using implementation strategies reported by previous work. The productivity and quality outcomes were quantified in our study by means of merged PRs and bug-related closed issue reports, respectively, inspired by Vasilescu et al.'s work [30]. However, these measures only represent a specific perspective of productivity and quality. Other existing perspectives (e.g., effort) could be explored in future work. Our study also explored a total of 5 sub-practices related to CI. Although they are expressive, this set could be extended with the analysis of other existing sub-practices such as code coverage, for example. In fact, we did not include the code coverage analysis because we only found a few GitHub projects maintaining a history of this metric in existing public repositories. Felidré et al.'s work [10] found the same obstacle when analyzing existing CI practices on GitHub projects.

In the qualitative analysis of RQ3, we acknowledge the potential for bias from the authors' subjective interpretations of PR comments related to CI sub-practices. We mitigate this threat through the peer review of the codes and themes extracted using document analysis. As mentioned in Section 2, we also omitted the groups of projects of the PR comments that were analyzed to avoid bias related to some specific group (i.e., *The Best and Worst CI Practices groups*.)

**External Validity Threats:** All analyzed open-source projects are hosted on the GitHub platform and use Travis-CI as their CI service. We acknowledge that our results are restricted to the context of the analyzed projects. Additional replication studies are necessary in future work to generalize the results for projects of different nature. We made our dataset and results available to allow future replications of our study.[13]

## 6 RELATED WORK

In this section, we situate our study with respect to previous works that investigate the relationship between CI sub-practices and the productivity and quality of open-source projects.

Vasilescu et al. [30] developed a study to discern the effects of CI adoption in quality and productivity outcomes. They collected a dataset of 246 GitHub projects which at some point in their history added the Travis CI to the development process. They found that after adopting the Travis CI, teams are significantly more effective at merging pull requests submitted by core members. They also report that core developers in teams using CI are able to discover significantly more bugs than in teams not using CI. Similar to Vasilescu's et al. work, our study investigated the CI impact on the productivity and quality of popular open-source projects by using the merged PR and closed bug-related issue outcomes, respectively. On the other hand, we analyzed the influence of different CI sub-practices on these outcomes. Our study found similar results for: (i) productivity - in which the *Commit Activity and Build Activity* CI sub-practices have impacted positively in the number of merged PRs; and (ii) quality - where the *Build Activity* CI sub-practice is correlated with an increase in the number of closed bug-related issues.

Bernardo et al. [2] analyzed 162,653 pull requests of 87 GitHub projects that are implemented in 5 different programming languages

to evaluate the impact of adopting CI on the time to deliver merged PRs. Their work also considers the beginning of the usage of Travis online service as the starting point for the CI adoption. The study reports a large increase in the number of submitted, merged and delivered PRs per release after the CI adoption. Similar to Bernardo et al. work [2], our work also found an increase in the number of merged PRs but correlated with the *Commit Activity and Build Activity* CI sub-practices. In addition to that, we investigated the improvement of the quality of the project in terms of closed bug-related issues.

Felidré et al. [10] analyzed 1,270 open-source projects that use Travis CI and quantitatively studied the behavior of the adoption of existing CI sub-practices. They found that: (i) 60% of the projects have infrequent commits; (ii) 85% of the projects have at least one broken build that take a long time to be fixed; and (ii) most of projects have a build that executes with more than 10 minutes. Similar to Felidré et al. [10], we found indication that many Travis CI projects do not pay attention to the behavior of CI sub-practices during their development and evolution. However, the main goal of our work was assessing the benefits of these practices for the productivity and quality of projects. As our results showed, tracking these sub-practices is important to merge more PRs, generate fewer errors and face fewer problems during the development process.

## 7 CONCLUSION

We conducted an exploratory quantitative and qualitative study investigating how Continuous Integration practices may influence the productivity and quality outcomes of open-source projects. We analyzed a set of 90 relevant and active open-source projects for a period of 2 years.

We applied linear regression models to study potential associations between five CI sub-practices and the number of merged PRs (productivity). In addition, we studied the potential associations between the five CI sub-practices and the number of bug-related issue reports (quality).

Our findings revealed a positive correlation between the *Commit Activity* and *Build Activity* CI sub-practices and the increase in the number of merged pull requests. We also observed that *Build Activity* has a significant correlation with the number of bug-related issue reports. Moreover, a sound *Build Health* is associated with a decrease in the number of bug-related issue reports. Lastly, larger values of *Time to Fix a Broken Build* are associated with less bug-related issues.

To complement our quantitative analysis, we performed a qualitative Document Analysis to identify CI-related themes in the comments of pull requests. Our analysis suggest that higher values in the five analyzed CI sub-practices can indicate a better quality in the development process.

---

[13]https://zenodo.org/record/6513155

# REFERENCES

[1] Miltiadis Allamanis, Hao Peng, and Charles Sutton. 2016. A Convolutional Attention Network for Extreme Summarization of Source Code. *CoRR* abs/1602.03001 (2016). arXiv:1602.03001 http://arxiv.org/abs/1602.03001

[2] João Helis Bernardo, Daniel Alencar da Costa, and Uirá Kulesza. 2018. Studying the Impact of Adopting Continuous Integration on the Delivery Time of Pull Requests. In *Proceedings of the 15th International Conference on Mining Software Repositories*. Association for Computing Machinery, New York, NY, USA, 131–141. https://doi.org/10.1145/3196398.3196421

[3] Francesco S Bersani, Daniel Lindqvist, Synthia H Mellon, Elissa S Epel, Rachel Yehuda, Janine Flory, Clare Henn-Hasse, Linda M Bierer, Iouri Makotkine, Duna Abu-Amara, et al. 2016. Association of dimensional psychological health measures with telomere length in male war veterans. *Journal of affective disorders* 190 (2016), 537–542.

[4] Glenn Bowen. 2009. Document Analysis as a Qualitative Research Method. *Qualitative Research Journal* 9 (08 2009), 27–40. https://doi.org/10.3316/QRJ0902027

[5] Hyunyoung Choi and Hal Varian. 2012. Predicting the present with Google Trends. *Economic record* 88 (2012), 2–9.

[6] Cuelogic. 2017. *The Levenshtein Algorithm*. Retrieved december 02, 2021 from https://www.cuelogic.com/blog/the-levenshtein-algorithm

[7] Marcos César de Oliveira. 2017. DRACO: Discovering Refactorings That Improve Architecture Using Fine-Grained Co-Change Dependencies. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering* (Paderborn, Germany) *(ESEC/FSE 2017)*. Association for Computing Machinery, New York, NY, USA, 1018–1021. https://doi.org/10.1145/3106237.3119872

[8] Marcos César de Oliveira, Davi Freitas, Rodrigo Bonifácio, Gustavo Pinto, and David Lo. 2019. Finding needles in a haystack: Leveraging co-change dependencies to recommend refactorings. *Journal of Systems and Software* 158 (2019), 110420. https://doi.org/10.1016/j.jss.2019.110420

[9] Paul Duvall, Stephen M. Matyas, and Andrew Glover. 2007. *Continuous Integration: Improving Software Quality and Reducing Risk (The Addison-Wesley Signature Series)*. Addison-Wesley Professional.

[10] Wagner Felidré, Leonardo Furtado, Daniel A. da Costa, Bruno Cartaxo, and Gustavo Pinto. 2019. Continuous Integration Theater. In *2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. 1–10. https://doi.org/10.1109/ESEM.2019.8870152

[11] Martin Fowler. 2006. *Continuous Integration*. Retrieved november 24, 2021 from https://martinfowler.com/articles/continuousIntegration.html

[12] GeeksforGeeks. 2021. *Jaro and Jaro-Winkler similarity*. Retrieved december 02, 2021 from https://www.geeksforgeeks.org/jaro-and-jaro-winkler-similarity/

[13] Taher Ghaleb, Daniel Costa, and Ying Zou. 2019. An Empirical Study of the Long Duration of Continuous Integration Builds. *Empirical Software Engineering* 24 (08 2019). https://doi.org/10.1007/s10664-019-09695-9

[14] Taher Ahmed Ghaleb, Daniel Alencar da Costa, Ying Zou, and Ahmed E Hassan. 2019. Studying the impact of noises in build breakage data. *IEEE Transactions on Software Engineering* (2019).

[15] UCLA: Statistical Consulting Group. 2022. *ZERO-INFLATED NEGATIVE BINOMIAL REGRESSION | R DATA ANALYSIS EXAMPLES*. Retrieved may 02, 2022 from https://stats.oarc.ucla.edu/r/dae/zinb/

[16] Manjul Gupta, Joey F George, and Weidong Xia. 2019. Relationships between IT department culture and agile software development practices: An empirical investigation. *International Journal of Information Management* 44 (2019), 13–24.

[17] Frank Harrell. 2022. *ols: Linear Model Estimation Using Ordinary Least Squares*. Retrieved may 02, 2022 from https://rdrr.io/cran/rms/man/ols.html

[18] Michael Hilton, Nicholas Nelson, Timothy Tunnell, Darko Marinov, and Danny Dig. 2017. Trade-Offs in Continuous Integration: Assurance, Security, and Flexibility. Association for Computing Machinery, New York, NY, USA. https://doi.org/10.1145/3106237.3106270

[19] FRANKLINE KIBUACHA. 2021. How to Determine Sample Size for a Research Study. (2021). https://www.geopoll.com/blog/sample-size-research/

[20] Eero Laukkanen, Maria Paasivaara, and Teemu Arvonen. 2015. Stakeholder Perceptions of the Adoption of Continuous Integration – A Case Study. In *2015 Agile Conference*. 11–20. https://doi.org/10.1109/Agile.2015.15

[21] Steve McConnell. 2004. *Code Complete: A Practical Handbook of Software Construction* (2 ed.). Microsoft Press, Redmond, WA. https://www.safaribooksonline.com/library/view/code-complete-second/0735619670/

[22] Andrew Glover Paul Duvall, Steve Matyas. 2007. *Introducing continuous integration*. Retrieved june 21, 2022 from https://www.infoworld.com/article/2077731/introducing-continuous-integration.html?page=3

[23] Wouter Poncin, Alexander Serebrenik, and Mark van den Brand. 2011. Process Mining Software Repositories. In *2011 15th European Conference on Software Maintenance and Reengineering*. 5–14. https://doi.org/10.1109/CSMR.2011.5

[24] Gustavo Sizilio Nery, Daniel Alencar da Costa, and Uirá Kulesza. 2019. An Empirical Study of the Relationship between Continuous Integration and Test Code Evolution. In *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. 426–436. https://doi.org/10.1109/ICSME.2019.00075

[25] Eliezio Soares, Gustavo Sizílio, Jadson Santos, Daniel Alencar da Costa, and Uirá Kulesza. 2021. The Effects of Continuous Integration on Software Development: a Systematic Literature Review. *CoRR* abs/2103.05451 (2021). arXiv:2103.05451 https://arxiv.org/abs/2103.05451

[26] Laerd Statistics. 2021. *Spearman's Rank-Order Correlation*. Retrieved december 02, 2021 from https://statistics.laerd.com/statistical-guides/spearmans-rank-order-correlation-statistical-guide.php

[27] Daniel Ståhl and Jan Bosch. 2014. Modeling continuous integration practice differences in industry software development. *Journal of Systems and Software* 87 (2014). https://doi.org/10.1016/j.jss.2013.08.032

[28] STHDA Statistical tools for high-throughput data analysis. 2021. *Correlation Test Between Two Variables in R*. Retrieved december 02, 2021 from http://www.sthda.com/english/wiki/wiki.php?id_contents=7307

[29] STHDA Statistical tools for high-throughput data analysis. 2021. *Normality Test in R*. Retrieved december 02, 2021 from http://www.sthda.com/english/wiki/normality-test-in-r

[30] Bogdan Vasilescu, Yue Yu, Huaimin Wang, Premkumar Devanbu, and Vladimir Filkov. 2015. Quality and Productivity Outcomes Relating to Continuous Integration in GitHub. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. Association for Computing Machinery, New York, NY, USA, 805–816. https://doi.org/10.1145/2786805.2786850

[31] Carmine Vassallo, Fabio Palomba, Alberto Bacchelli, and Harald C. Gall. 2018. Continuous Code Quality: Are We (Really) Doing That?. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering* (Montpellier, France) *(ASE 2018)*. Association for Computing Machinery, New York, NY, USA, 790–795. https://doi.org/10.1145/3238147.3240729

[32] David Gray Widder, Michael Hilton, Christian Kästner, and Bogdan Vasilescu. 2019. A Conceptual Replication of Continuous Integration Pain Points in the Context of Travis CI. Association for Computing Machinery, New York, NY, USA. https://doi.org/10.1145/3338906.3338922

[33] Kazuhiro Yamashita, Shane McIntosh, Yasutaka Kamei, and Naoyasu Ubayashi. 2014. Magnet or sticky? an oss project-by-project typology. In *Proceedings of the 11th working conference on mining software repositories*. 344–347.

[34] Yangyang Zhao, Alexander Serebrenik, Yuming Zhou, Vladimir Filkov, and Bogdan Vasilescu. 2017. The impact of continuous integration on other software development practices: A large-scale empirical study. In *2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*. 60–71. https://doi.org/10.1109/ASE.2017.8115619