

1 Usage and structure of continuous integra-
2 tion as configuration?

3 Joseph Ling
j1653@kent.ac.uk



School of Computing
University of Kent
United Kingdom

Word Count: around 4,400

4 March 10, 2020

6 Continuous integration (CI) is becoming more popular as software develop-
7 ment moves to an Agile fast paced development cycle. Most CI is run based
8 of configuration therefore how much is it acutally being used? As well as
9 how are these files being structured? We got XXXX open source projects
10 from Github to answer these questions. In doing so compared our results
11 against Michael Hilton, Marinov and Dig (2016) work to see if their has been
12 a increase in usage. We found a shift in CI services being used and were
13 able to get similar results to their study. In terms of structure we found that
14 configuration files are written with no comments normally. We suggest at the
15 end further research is needed to get a better understanding of this growing
16 field.

1 Introduction

Continuous integration (CI) is becoming more popular over the last few years. This can be seen by how major version control hosting services Github, Bitbucket and Gitlab have all started to or have been improving their CI product. In terms of research, configuration as code Rahman, Mahdavi-Hezaveh and Williams (2019) and continuous integration Copeland (2010) with Shahin, Ali Babar and Zhu (2017) demonstrating breadth of the research.

Continuous integration is a process of automatically running compiling, running tests and checking that the product works. This can be combined with Continuous Delivery where the product is deployed or released after it has gone through CI.

This can get complicated quickly therefore configuration as code (or infrastructure as code) is used to configure it. The main kind of configuration format used for this is yaml (reference to what it is??) followed by xml and java based scripting formats.

In terms of looking at usage we are going to do a similar look at the data as did Michael Hilton, Marinov and Dig (2016). The important aspect will be looking at how usage has changed over the last 5 years along with looking more closely at which repositories are more likely to use CI/CD. For this we are going to focus on the following research questions:

- What percentage of open-source projects use CI?
- multiple CI used
- what is the breakdown of usage of different services?
- Do certain types of projects use CI more than others?

This should give us a better understanding of the sample of repositories from Github. From there we look at the structure of the configuration files to understand how certain aspects of it are used.

- 45 • configuratizon errors when loading the config (just yaml parsing errors
46 atm)
- 47 • how are comments used in the configuration?
- 48 • Are external scripts used within the configuration?

49 A key aspect is that these questions do not look too deeply into the
50 individual implementation of each CI system. This is because there are
51 already some good papers looking Gallaba and McIntosh (2018) at this but
52 in order to be able to compare the different configuration types it is important
53 to compare similar attributes (there is also a time factor in here as well).

54 2 Previous Works

55 2.1 Continous integration

56 Continous integration is the frequent submission of work normally tied into
57 a feedback loop. For example using version control daily committing changes.
58 That then a server builds and tests the changes informing you of status of
59 those cahnges. The generally agree upon detailed definition is Fowler (2010).

60 2.2 Usage of continous integration

61 The actual usage of continous integration as configuration was looked at
62 by Michael Hilton, Marinov and Dig (2016). In this they use three source
63 of information github repositories, travis builds and a survery. In order to
64 be do a more systematic study of CI usage than Vasilescu et al. (2015). In
65 analysing that data they found that "The trends that we discovered point
66 to an expected growth of CI. In the future, CI will have an even greater
67 influence than it has today.". As we are looking at the same question we will
68 use four of the same research questions out of the fourteen. In order to see
69 what difference four years has made to the growth of usage of CI.

70 2.3 Config as code

71 Configuration as code or infrastructure as code has been an increasing area
72 of research over the last few years. There seems to be slightly more research
73 in infrastructure as code Rahman, Mahdavi-Hezaveh and Williams (2019).
74 There has been a focus on Puppet and Chef, for example in Sharma, Frangkoulis
75 and Spinellis (2016) looks at code quality by the measure of "code smell" of
76 Puppet code. This tackles the problem by defining by best practices and
77 analyzing the code against that. In the case of Cito et al. (2017) it uses
78 the docker linter in order to be able to analyse the files. For the continuous
79 integration systems we pick we will look into the tooling around that to aid
80 the analysis.

81 3 Methodology

82 In order to get repositories with CI/CD configuration from Github we have a
83 number of approaches. The first is to use the search for particular files but
84 this is limited to only 1000 results. The alternative is to search for repositories
85 and we bypass the 1000 result limit to an extent by getting results for every
86 'star' count (stars are used to like or upvote a repository). Although this will
87 be giving us a lot of results it will still only be a sample of the population but
88 will give us a wider range of results. As there is rate limiting multiple github
89 api keys can be used to speed up the scraping of data (gitter could also
90 be used to speed up the process I think).

91 After we have got a repository we need to get the CI/CD files from it.
92 This is fairly easy as the CI/CD systems normally require a strict naming
93 convention and location within the repository. However as most of them are
94 yaml based you can have ".yaml" and ".yml" and users can use all sorts of
95 mixtures of upper and lower case. We try to account for this but won't get
96 every scenario. This combined with the fact that we are only looking for
97 top configuration files based on github (2017) along with github actions and
98 azure pipelines. Is why we also check repositories for their README.md file

99 to check if it has a build tag.

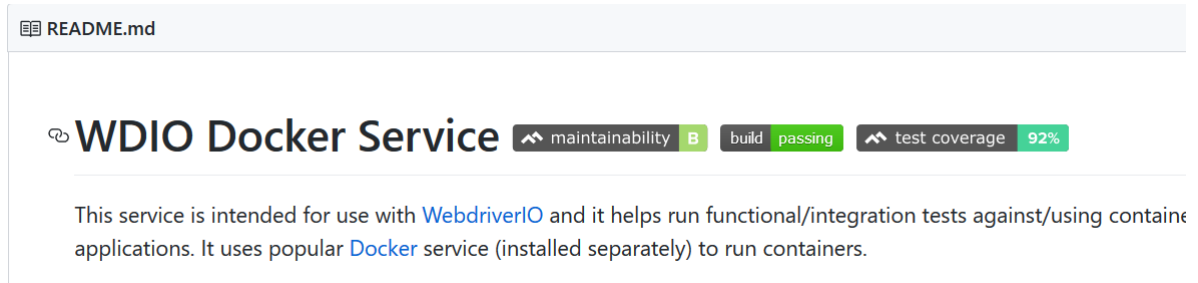


Figure 1: Example of CI tag for Github ReadMe

where did this image come from??
reference it man

100
101 In doing so it should give a wider net when sampling and help to under-
102 stand when a CI system is either not using configuration as code or using a
103 different CI system.

104 There are dangers in scraping data off github in terms of assumptions to
105 do with the population as found in Kalliamvakou et al. (2014). Our dataset
106 does not contain any forked repositories. But due to time constraints number
107 of commits and frequency of recent commits has not been looked at. This
108 would be an interesting area of further research in order to improve the
109 quality of the sample but also to look at how that affects the frequency of
110 CI usage.

111 Additionally the assumption that all repositories are of programming
112 projects with code in them is wrong. A number of repositories can be used
113 for storage, experimental, academic and other things. However they to all
114 some extent can use CI/CD for their work as a number of books were found
115 when looking through the dataset could use CI/CD.

116 Tooling for the configuration files, I looked into Travis, Github Actions
117 and Jenkins to work out whether or not it could aid in the research or not.
118 As a key part of understanding the first relies on knowing whether or not it is
119 valid. In terms for travis there is currently two parsers to validate the config-
120 uration. One which is deprecated since 2017 travis (2017) the other which is
121 currently in development travis (2020). Both didn't provided the necessary
122 results with the most recent one not being able to handle default fields. For

123 Github Actions as it's still a new tooling for it hasn't been developed out-
124 side of the Github editor web page ([https://github.community/t5/GitHub-](https://github.community/t5/GitHub-Actions/YAML-validator-for-Github-Actions-possible-expansion-of/td-p/29557)
125 [Actions/YAML-validator-for-Github-Actions-possible-expansion-of/td-p/29557](https://github.community/t5/GitHub-Actions/YAML-validator-for-Github-Actions-possible-expansion-of/td-p/29557)).
126 For Jenkins which is older solution allows validation through http/ssh request
127 to the Jenkins server (Gitlab follows this style as well) Jenkins (2020) Gitlab
128 (2020). This could work well although would require setting up a server for
129 each configuration type and might not validate if variables from the config
130 aren't defined on the server. As well as it would be best to be able to validate
131 them all or none of them in terms of being able to compare results easily.

132 4 Usage of CI

133 4.1 What percentage of open-source projects use CI?

134 Based a search for configuration as configuration files for the following CI
135 systems: Travis, Gitlab, Azure, App Veyor, Drone, Jenkins, Github, Circleci,
136 Semaphore, Teamcity and buildkite. Wrecker got bought by Oracle and from
137 doing a search on Github for what I think based on the docs (docs: Wrecker
138 and Oracle (2018) and search: GitHub (2020)) for their config file naming
139 convention. I was only able to find 20 results so did not include in the scraping
140 script to speed up the process of searching for the other configuration file
141 formats.

CI/CD	count	repos with config	no. multiple	multiple percent
config file(s)	12128	38.51%	1675	13.81%
found in ReadMe	873	2.77%		
none found	18493	58.72%		

Table 1: Percentage of CI used for projects

142 Our sample of repositories is 31,494 in comparison to Michael Hilton,
143 Marinov and Dig (2016) which had a sample of 34,544. The percentage of
144 CI projects they had was 40.27%. As if you combined the "config file(s)"

145 and "found in ReadMe". However in order to work out if a project might be
146 using CI but the config file wasn't picked a search string is used. Therefore
147 it is not as accurate as finding a config file as their could be false postives.

148 However that doesn't give us too much insight into the dataset. Here is a
149 graph showing the subscribers plotted against the number of stars. The key
150 here to understand is not potentially any correlation but to see the spread
151 of data that the table is showing.

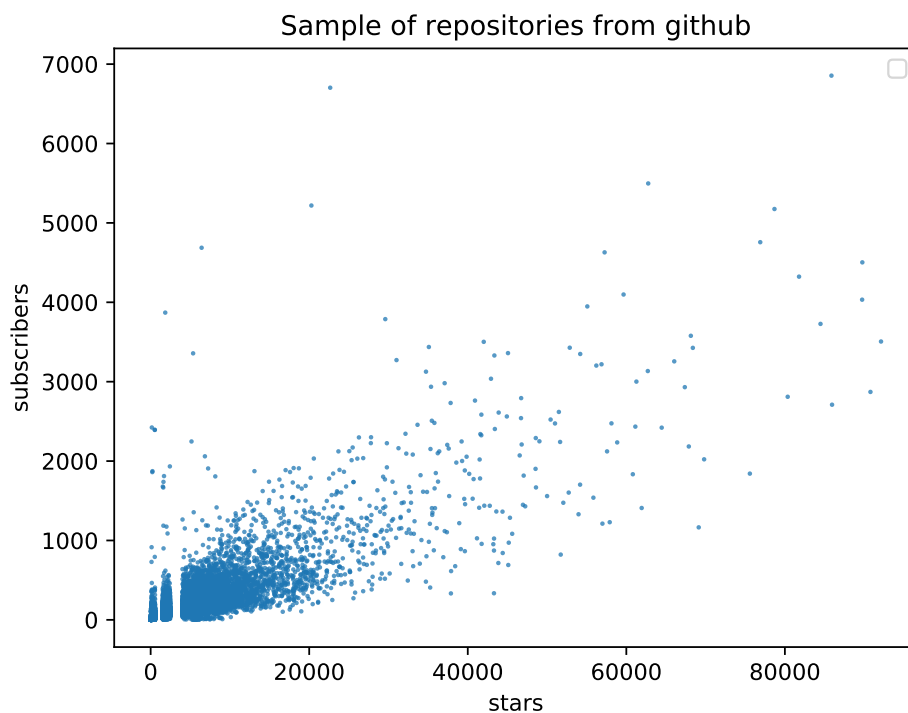


Figure 2: Scatter graph of Github stars against subscribers

152 Figure 2 helps give a understanding to the give a depth of the data for
153 where the graph is just blue. This is because on Github you get more repos-
154 itories with smaller star counts than large ones.

155 Figure 3 provides insight into the density of the data for between 0 to
156 25000.

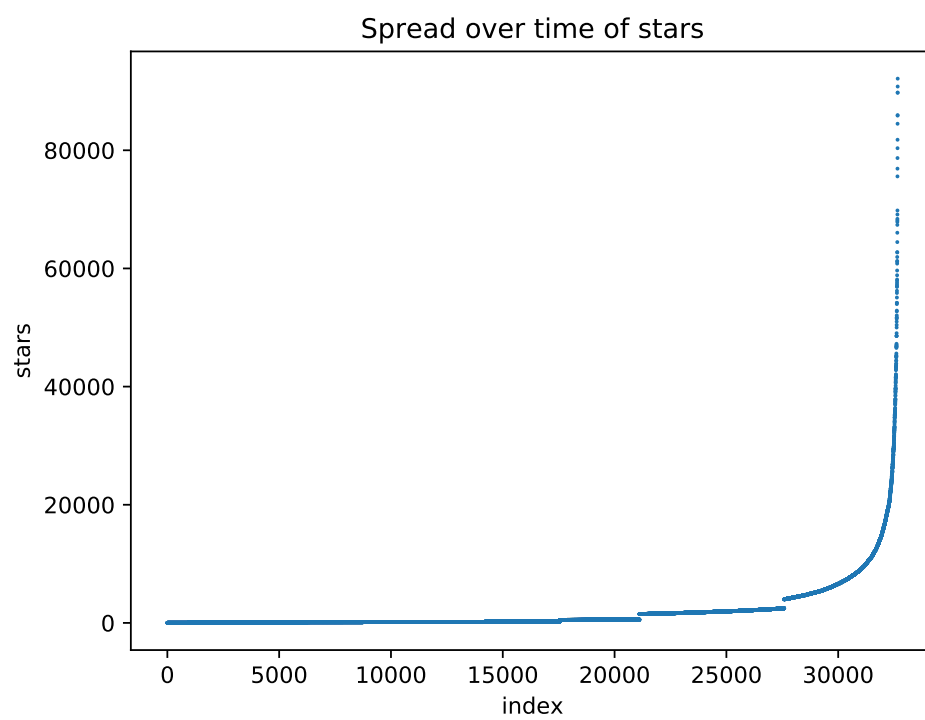


Figure 3: Stars graph

157 4.2 What CI systems are projects using?

158 Like all other research travis is the most popular CI system in use. However
159 over the last 4 years since the github (2017) Circleci has lost out on it's rough
160 quarter that it owned. In particular the rise of github actions seems to have
161 taken second place even though it is still very young in comparison (DATES).
162 However this might not be down to the Circleci loosing out on their existing
163 share. But potentially as the rise in CI usage goes up on github. Projects
are more likely to pick in the built in solutions to github.

Table 2: Configuration types spread

	config	percentage
travis	10607	74%
github	2301	16%
circleci	1109	8%
jenkinsPipeline	161	1%
drone	84	1%
buildkite	32	0%
teamcity	4	0%
semaphore	2	0%
azure	1	0%

4.3 Do certain types of projects use CI more than others?

Below shows all the CI projects sorted then grouped together per 540 projects. Then in this case we choose to categories via star count for each project.

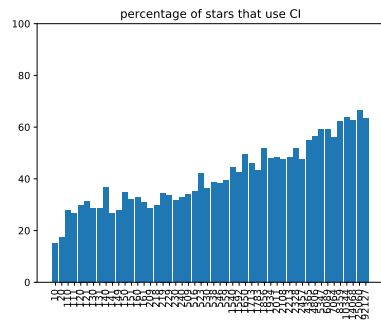


Figure 4: 2020 dataset

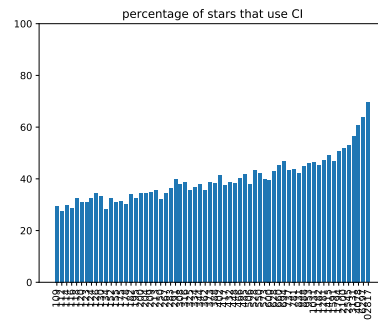


Figure 5: 2016 dataset

In Figure 4 is the results from this research and in Figure 5 is the results from Michael Hilton, Marinov and Dig (2016).

Here we are comparing whether or not in the last 4 years the number of stars increases the CI being used. Their seems to a steeper gradient in the more recent datasets. However as 4 starts at zero stars and 5 starts at 100 stars their is signifacant dip at the start of the first graph.

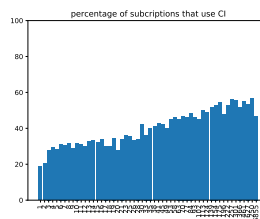


Figure 7: Subs graph

Figure 7 uses the same method as Figure TODO SORT except is does it based the number of subscribers. Subscribers are used on github to keep

175 update on the changes on the project. This ranges from core team memem-
176 bers working on the project to people that want to be notified about a new
177 release. In looking at this metric the hypothesis was that it would have a
178 sharper rise in percentage of projects using CI per subscriber. However that
179 was not the case overall the gradient is not as strong. There is no compar-
180 isson to Michael Hilton, Marinov and Dig (2016) because their final corpus
181 does not contain subscriber count for each project.

182 5 Config file results

183 5.1 configuration errors when loading the config (just 184 yaml parsing errors atm)

Composer error In the example it has two steps that are using an yaml anchor. This allows for the yaml below it to be referenced somewhere else. However if you define the anchor twice it causes a composer error. As you have two references for the samething so it won't know which one to use.

```
definitions:  
steps:  
- step: &build-test  
name: Build and test  
script:  
- mvn package  
- step: &build-test  
name: deploy  
script:  
- ./deploy.sh target/my-app.jar
```

Scanner error The first step of loading the yaml is to scan it to create the tokens. However invalid characters such as "\t" are invalid.

```
definitions: \t
```

185 As can be seen in the table their our configuration files with yaml errors
186 meaning that the CI for that project will not load. Yet it seems that a
187 very small percentage of projects that have them. For example the two

Parse error In this example it has scanned the file and created tokens for the syntax. Now it parses the syntax and works out if each token is valid given it's current context. In this case a closing `]` without an opening `[` is invalid.

```
definitions: ]
```

Table 3: yaml configuration errors

config	composer error	constructor error	parse error	scanner error	no. config
circleci	1	0	0	1	1109
drone	31	0	0	0	84
github	0	1	0	3	2301
travis	6	0	10	21	10607
buildkite	0	0	0	0	32
semaphore	0	0	0	0	2
azure	0	0	0	0	1

highest configuration types with errors are drone (36.90%) followed by travis (0.348%).

In the case for drone all the errors are for the same type of error. Potentially this could be because of how anchors are a lot more common in drone.

For travis it is the most common form of CI found therefore it is more likely to contain more errors. Yet with such a small amount it seems like yaml errors aren't a major problem in CI. Although as they are required to be fixed in order for the CI to run the chances of it working are higher and a more detailed study would need to be done.... ah

5.2 How are comments used in configuration?

198 The assumption was the as continuous integration setups can be compli-
199 cated and have edge cases. Therefore comments would be used to describe
200 and handle that complexity.

201 An example configuration file below for Github actions using the default
202 template slightly altered. Shows two examples of comment usage, the first
203 being including useful information about why a particular version of the
204 programming language was chosen. The second is that the tests have been
205 disabled by commenting them out.

In order to pick up on all these different types of comments. All the CI files were parsed and then regular expressions were used to pick on up key factors such as "note:". Along with multiple single line comments which made up a block/multi-line comment.

For example in to the left there is an example Github Action yaml file. If were it would be parsed we would get: one multi line comment, 15 lines of code, 1 single line comment, a total of 5 comments and 20 lines in the file. Therefore their is a their is a raito of 4:1 for code in this config file.

```
name: Python package
on: [push]
jobs:
  build:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v2
      - name: Set up Python
        uses: actions/setup-python@v1
        # note: only works with python 3
        with:
          python-version: 3.8
      - name: Install dependencies
        run: |
          python -m pip install --upgrade pip
          pip install -r requirements.txt
      # - name: Test with pytest
      #   run: |
      #     pip install pytest
      #     pytest ./src
```

206 Initially before we look at the comments it is important to understand
207 how the rest of the file is made up. In the graph below (Figure 8) it shows
208 how each configuration type is made up by mean of each part of the file. For
209 all the yaml based configurations lines of code and number of lines in total
210 are very close. Then for the number of commmets being very very small on

211 average.

212 In the case for Jenkins pipelines and teamcity there is a much higher
213 usage of having code with comments.

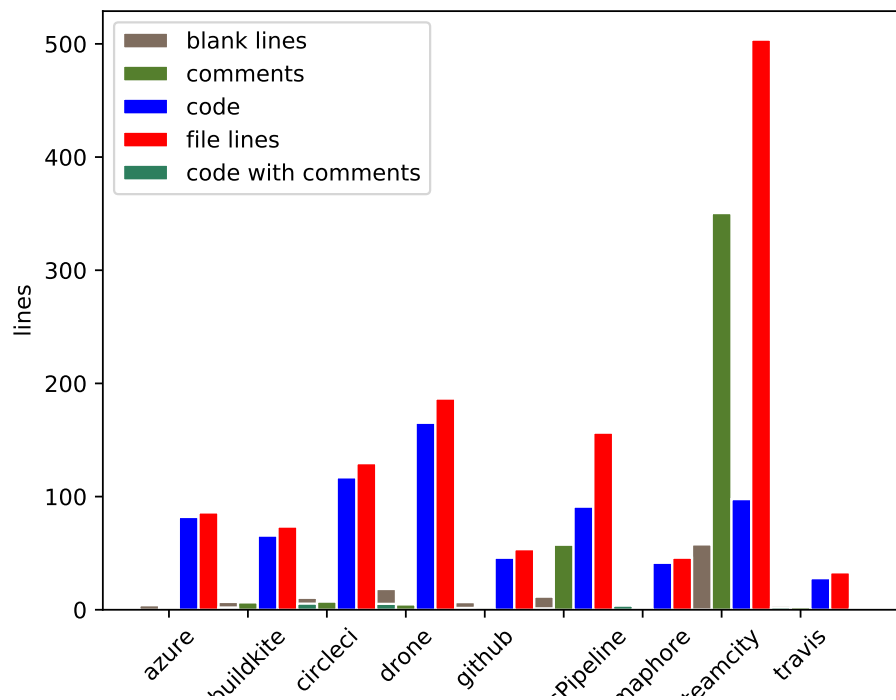


Figure 8: Mean of line counts

214 Raitos:

215 • code: comments

216 • code: line total

217 • code: blank lines

218 • single line comment: multiline comment

219 • single line comment: code with comment

In Figure 9 a regular expression was used to label the comments. There were key different types of comment that we wanted to find. The first being the commented out code which we did by searching for version numbers in comments. The second being useful information about the structure of the CI file such todo, note, important comments (e.g. `//todo`). In order to increase the search for this we included searching for urls and separation comments (e.g. `//===`).

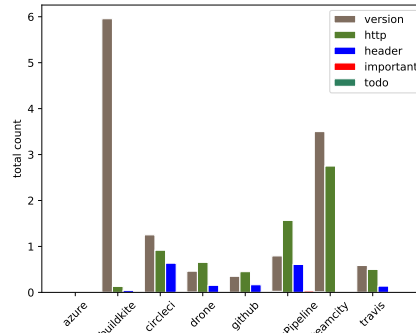


Figure 9: Comment types

220 From labelling the comments in Figure 9 we can see that having com-
 221 mments with versions in and urls is most common. This could indicate
 222 comments from templates or how they are commented. Although yet again
 223 the amount of labels found on average is still very low.

224 Overall we have found that comments are not used a lot. In the cases
 225 that they are used it's more likely to be from a configuration template or
 226 commenting out configuration.

227 5.3 Are external scripts used within the configuration?

228 An external script is a bash or powershell script typically depending on the
229 operating system. It can be used to build, deploy or do any step that CI
230 takes. The key difference between it and the CI configuration is that it be
231 executed on a users machine. Therefore you do get some setups where you
232 have scripts defined for building and deploying the code that the users and CI
233 both use. Most CI systems allow for "script" tags to be used which could be
234 descibed as an internal script. Therefore external scripts are defined outside
235 the CI configuration in the directory.

236 The methodology we used to handle this was too look at how many bash
237 or powershell scripts where used in CI. Using the code the parsed the yaml
238 files for comments we were able to check do a using a regular expression for
239 either of those files.

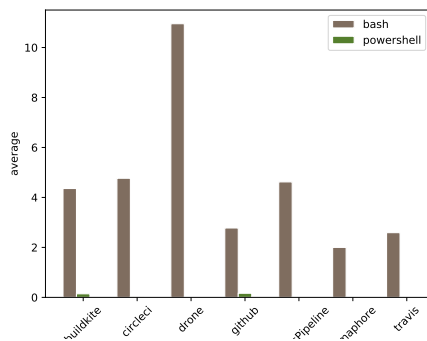


Figure 10: Comment types

	bash	powershell
buildkite	61	2
circleci	1497	8
drone	230	0
github	1097	65
jenkinsPipeline	171	0
semaphore	2	0
travis	5937	3

Figure 11: sum of scripts used

240 In Figure 10 we have the average number of times a script is used for a
241 configuration file that already has a script being used.

242 As some of the necessary actions are being done in the scripts and not in
243 the CI file. Potentailly there could be less lines of code in the configuration
244 for files that use scripts. However in Figure 12 we can see that the data is
245 very spikey with outliers. Then in Figure 13 we can see the same affect when

246 trying to see if the more popular a project is affects the chances of it using
247 CI.

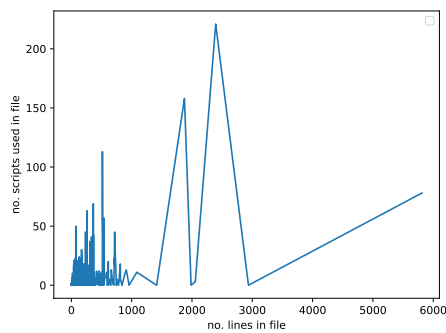


Figure 12: no. scripts to no. lines

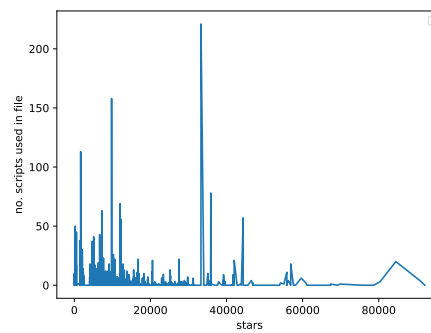


Figure 13: no. scripts to stars

248 To conclude we can see that scripts are used but that much. And their
249 no correlation between lines of code and usage of external scripts.

250 6 Threats to validity

251 The major and most obvious threat is the sample gathered from scraping
252 the data from Github. This has already been touched on in the 3 section but
253 now we are going to look at it in more detail.

254 Firstly if we assume that the scraping works perfectly then it's only at
255 maximum a 1000 open source projects per star. That is excluding closed
256 source projects which would range from personal projects to companies. As
257 well as it is only data from Github not from gitlab, bitbucket or other version
258 control hosting services. This leads to bias in the data for example if gitlab
259 was also scraped then we would get a lot more gitlab ci files. However in
260 order to get best spread of data Github has the best api and most services do
261 not tie you down to use only their service. As well although we could get a
262 1000 projects per star we were still able to get around 30,000 projects and
263 a wide spread across Github. The key aspect being that because it was a
264 sample we focused on getting a good spread of data.

265 Secondly the scraping script is not perfect in how it finds configuration
266 files. As it only looks in the top level directory for the file name pattern
267 described in their docs or unique folder. Therefore if the systems allowed
268 many different names or different names in past it wouldn't have picked it
269 CI system. Additionally we only decided to scrape for certain CI files. Yet
270 we chose a good scope based on previous research into the top CI files. As
271 well the scraping script has been tested worked on to try and minimise any
272 bugs. In the case that we did not pick up a CI file we ran a regexp against
273 the ReadMe file to get a better understanding of the error bounds.

274 Thirdly identifying which projects are programming projects or would
275 have a need for CI. Based on the research Kalliamvakou et al. (2014) it is
276 important to filter out repositories that aren't part of the question being
277 asked. Therefore we could have looked to try and filter out github static
278 sites and other non software based projects. However if assume a certain
279 type of project won't be using CI then we would be introducing bias when
280 trying to answer how CI is used. For further research better labelling of what

281 kind of projects are which would potentially be beneficial though.

282 7 Summary

283 We got a sample of XXXX open source projects from Github and were able
284 to compare that to a previous study 4 years ago. In doing so we found that
285 usage of CI projects was similar and that more popular a project the higher
286 chance it would be using CI. This lined with the research from 4 years ago.
287 The major change was the increase in popularity of Github Actions taking
288 over second place from Circleci. Additionally we look at whether or not the
289 number of people watching the project had the same effect. It did but to a
290 lesser extent.

291 In terms of structure of CI configuration we looked each line of was used
292 in context of comments. We found that a very few projects use comments in
293 their CI. In terms of how they used scripts, we found the majority of projects
294 do not use external scripts.

295 7.1 Discussion and further research

296 In the process of writing this paper we kept on considering more research
297 questions. As there is a lot of meta data that you can get for a single
298 project, in addition to what was used for this paper.

299 Further research into usage that we would like to do is look into how
300 the size of the project affects the chance that it uses CI. Then looking at
301 the usage of scripts within CI configuration, for example using a script tag
302 to run a shell script. As while doing the research we found some projects
303 use scripts a lot while others just used the CI config. This would lead to
304 questions around which CI system has a higher amount of scripts used. But
305 also looking at how much they enable them to be used and what is the size
306 of those scripts. The data for the programming language and version(s) is in
307 the config. Therefore it would be possible to work out how much usage each
308 version is getting of a particular programming language.

309 Further research into structure could look into the naming of each part
310 of the build process that is used. This would be interesting as it would
311 provided insight into what terms are commonly used. As well an idea into

312 how people plan or don't plan out their configuration files. Additionally CI
313 systems can be designed to run on every commit to version control or only
314 commits to certain branches. Therefore by looking at the branching regexp
315 that are being used an better understanding of how branches are actually
316 used in software development where CI is also used could be found out.

317 In addition working on pruning our dataset using methods outlined in
318 Kalliamvakou et al. (2014).

319 8 Acknowledgement

320 We wish to thank Michael Hilton in particular for providing the corpus for
321 their research Michael Hilton, Marinov and Dig (2016).

322 [section ommitted]

323 References

324 Cito, J., Schermann, G., Wittern, J. E., Leitner, P., Zumberi, S. and Gall,
325 H. C. (2017). An Empirical Analysis of the Docker Container Ecosystem
326 on GitHub. In *2017 IEEE/ACM 14th International Conference on Mining*
327 *Software Repositories (MSR)*, pp. 323–333, iSSN: null.

328 Copeland, P. (2010). Google's Innovation Factory: Testing, Culture, and
329 Infrastructure. In *Proceedings of the 2010 Third International Conference*
330 *on Software Testing, Verification and Validation*, Washington, DC, USA:
331 IEEE Computer Society, ICST '10, pp. 11–14.

332 Fowler, M. (2010). Continuous integration. In
333 <https://www.martinfowler.com/articles/continuousIntegration.html>.

334 Gallaba, K. and McIntosh, S. (2018). Use and Misuse of Continuous Inte-
335 gration Features: An Empirical Study of Projects that (mis)use Travis CI.
336 *IEEE Transactions on Software Engineering*, pp. 1–1.

337 github (2017). <https://github.blog/2017-11-07-github-welcomes-all-ci-tools/>.
338 In github.com, ed., *github welcomes all ci tools*.

339 GitHub (2020). github filename search for wrecker.yml files. In *github file-*
340 *name search for wrecker.yml files*.

341 Gitlab (2020). <https://docs.gitlab.com/ee/api/lint.html>. In *Gitlab docs*.

342 Jenkins (2020). <https://jenkins.io/doc/book/pipeline/development/>. In
343 *Jenkins documentation*.

344 Kalliamvakou, E., Gousios, G., Blincoe, K., Singer, L., German, D. M. and
345 Damian, D. (2014). The promises and perils of mining GitHub. Hyderabad,
346 India: Association for Computing Machinery, MSR 2014, pp. 92–101.

347 Michael Hilton, K. H., Timothy Tunnell, Marinov, D. and Dig, D. (2016).
348 Usage, costs, and benefits of continuous integration in open-source projects
349 | Proceedings of the 31st IEEE/ACM International Conference on Auto-
350 mated Software Engineering.

351 Rahman, A., Mahdavi-Hezaveh, R. and Williams, L. (2019). A systematic
352 mapping study of infrastructure as code research. *Information and Soft-*
353 *ware Technology*, 108, pp. 65–77.

354 Shahin, M., Ali Babar, M. and Zhu, L. (2017). Continuous Integration, De-
355 livery and Deployment: A Systematic Review on Approaches, Tools, Chal-
356 lenges and Practices. *IEEE Access*, 5, pp. 3909–3943.

357 Sharma, T., Fragkoulis, M. and Spinellis, D. (2016). Does Your Configuration
358 Code Smell? In *2016 IEEE/ACM 13th Working Conference on Mining*
359 *Software Repositories (MSR)*, pp. 189–200, iSSN: null.

360 travis (2017). travis yaml (old repository). In [https://github.com/travis-](https://github.com/travis-ci/travis-yaml/)
361 [ci/travis-yaml/](https://github.com/travis-ci/travis-yaml/).

362 travis (2020). travis yaml new implementation. In [https://github.com/travis-](https://github.com/travis-ci/travis-yaml/)
363 [ci/travis-yaml/](https://github.com/travis-ci/travis-yaml/).

- 364 Vasilescu, B., Yu, Y., Wang, H., Devanbu, P. and Filkov, V. (2015). Quality
365 and productivity outcomes relating to continuous integration in GitHub.
366 Bergamo, Italy: Association for Computing Machinery, ESEC/FSE 2015,
367 pp. 805–816.
- 368 Wrecker and Oracle (2018). Wrecker ci development blog. In *Wrecker CI*
369 *development blog*.