# Usage and structure of continuous integration as configuration?

Joseph Ling

`jl653@kent.ac.uk`

School of Computing

University of Kent

United Kingdom

Word Count: around 4,400

March 10, 2020

## Abstract

Continous integeration (CI) is becoming more popular as software development moves to an Agile fast paced development life cycle. Most CI is done automatically using a service which run based of configuration. Our major questions is how much is CI acutally being used? As well as how are these files being structured? We got XXXX open source projects from Github to answer these questions. In doing so compared our results against (10) work to see if their has been a increase in usage. We found a shift in CI services being used and were able to get similar results to their study. In terms of structure we found that configuration files are written with no comments normally. We suggest at the end further research is needed to get a better understanding of this growing field.

similar is a bad word to use to describe the comparison

# 1 Introduction

Continous integeration (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 (12) and continuous integeration (2) with (13) demonstrating breadth of the research.

need to be clearer at what the bread of the research means or why it is useful

Continous integeration is a process of automatically running compiling, running tests and checking that the product works. This is can be combined with Continous 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 followed by xml and Java based scripting formats.

In terms of looking at usage we are going answer the same questions as (10) in order to do a comparison. The importanat 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 percentange of open-source projects use CI?

- what is the breakdown of usage of different services?

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

mutliple CI in one project comparison and data is possible

In being able to understand how projects from the sample are being used. From there we will look at the structure of the configuration files.

- What are the common errors when loading yaml configuration?

- How are comments used in the configuration?

- Are external scripts used within the configuration?

2

A key aspect is that these questions do not look too deeply into the individual implementation of each CI system. This is because there are already some good papers looking (4) at this but in order to be able to compare the different configuration types it is important to compare similar attributes.

## 2 Previous Works

### 2.1 Continous integeration

Continous integeration is the frequent submission of work normally tied into a feedback loop. For example using version control daily committing changes. That then a server builds and tests the changes informing you of status of those cahnges. The generally agree upon detailed definition is by Martin Fowler in how blog post on it: (3).

### 2.2 Usage of continous integeration

The actual usage of continous integeration as configuration was looked at by (10). In this they use three source of information github repositories, travis builds and a survery. In order to be do a more systematic study of CI usage than (17). In analysing that data they found that "The trends that we discovered point to an expected growth of CI. In the future, CI will have an even greater influence than it has today.". As we are looking at the same question we will use four of the same research questions out of the fourteen. In order to see what difference four years has made to the growth of usage of CI.

### 2.3 Config as code

Configuration as code or infrastructure as code has been an increasing area of research over the last few years. There seems to be slightly more research in infrastructure as code (12). The has been a focus on Puppet and Chef, for example in (14) looks at code quality by the measure of "code smell" of

Puppet code. This tackles the problem by defining by best practices and analyzing the code against that. In the case of (1) it uses the docker linter in order to be able to analyse the files. For the continous integeration systems we pick we will look into the tooling around that to aid the analysis.

# 3   Methodology

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

After we have got a repository we need to get the CI/CD files from it. This is fairly easy as the CI/CD systems normally require a strict naming convention and location within the repository. However as most of them are yaml based you can have ".yml" and ".yaml" and users can use all sorts of mixtures of upper and lower case. We try to account for this but won't get every scenario. This combined with the fact that we are only looking for top configuration files based on (5) along with github actions and azure pipelines. Is why we also check repositories for their ReadMe.md file to check if it has a build tag.

where did this image come from?? reference it man

In doing so it should give a wider net when sampling and help to understand when a CI system is either not using configuration as code or using a different CI system.

Tooling for the configuration files, we looked into Travis, Github Actions and Jenkins to work out whether or not it could aid in the research or not. As a key part of understanding the first relies on knowing whether or not
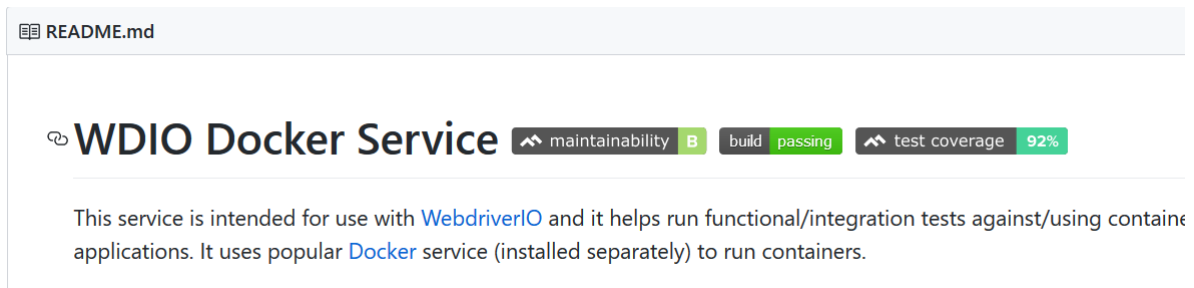
Figure 1: Example of CI tag for Github ReadMe

it is valid. In terms for travis there is currently two parsers to validate the configuration. One which is depracted since 2017 (15) the other which is currently in development (16). Both didn't provided the necessary results with the most recent one not being able to handle default fields. For Github Actions as it's still a new tooling for it hasn't been developed outside of the Github editor web page (11). For Jenkins which is older solution allows validation through http/ssh request to the Jenkins server (Gitlab follows this style as well) (8) (7). This could work well although would require setting up a server for each configuration type and might not validate if varaibles from the config aren't defined on the server. As well as it would be best to be able to validate them all or none of them in terms of being able to compare results easily.

# 4 Usage of CI

## 4.1 What percentange of open-source projects use CI?

Based a search for configuration as configuration files for the following CI systems: Travis, Gitlab, Azure, App Veyor, Drone, Jenkins, Github, Circleci, Semaphore, Teamcity and buildkite. Wrecker got bought by Oracle and from doing a search on Github for what we think based on the docs (18) and (6) for their config file naming convention. We were only able to find 20 results so did not include in the scraping script to speed up the process of searching for the other configuration file 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

Our sample of repositories is 31,494 in comparison to (10) which had a sample of 34,544. The percentage of CI projects they had was 40.27%. As if you combined the "config file(s)" and "found in ReadMe". However in order to work out if a project might be using CI but the config file wasn't picked a search string is used. Therefore it is not as accurate as finding a config file as their could be false postives.

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

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

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

## 4.2 What CI systems are projects using?

In Table 2 we find like all other research travis is the most popular CI system in use. However over the last 4 years since the (5) Circleci has lost out on it's rough quarter that it owned. In particular the rise of github actions seems to have taken second place even though it is still very young in comparison (DATES). However this might not be down to the Circleci loosing out on their existing 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.
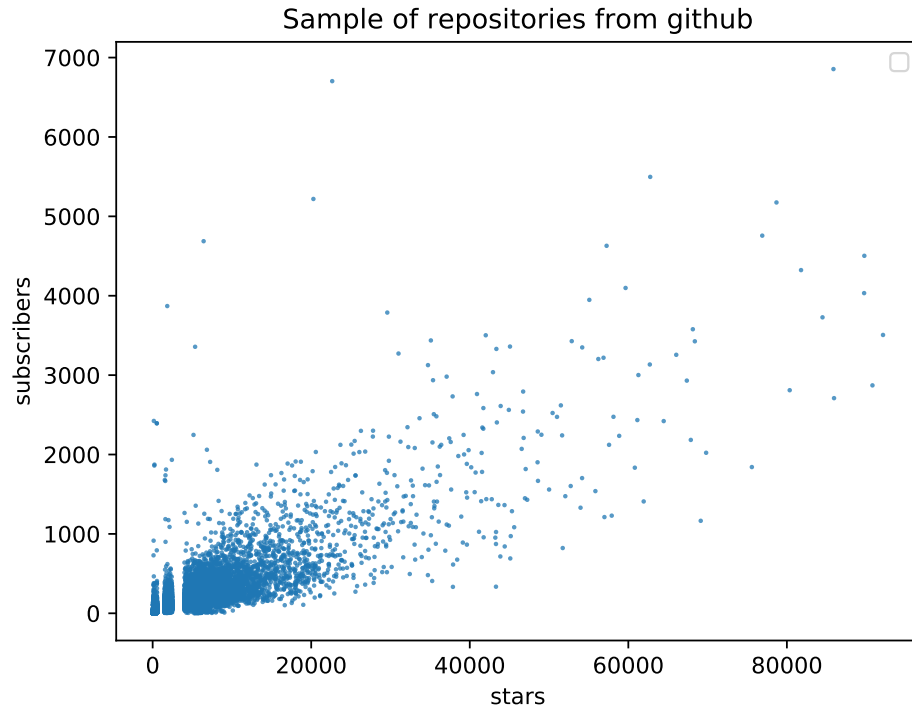
Figure 2: Scatter graph of Github stars against subscribers

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% |

Figure 3: Stars graph

## 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.

Here in Figure 4 and 5 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 uses the same method as Figure 4 except is does it based the number of subscribers. Subscribers are used on github to keep update on the changes on the project. This could range from core team memembers

when the writing is good and nearly polished make this the proper size
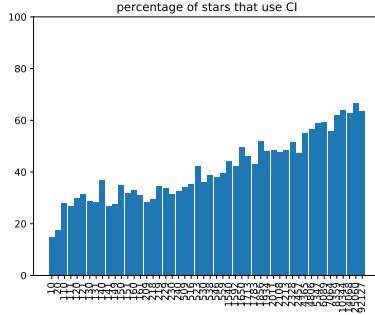
8

Figure 4: 2020 dataset



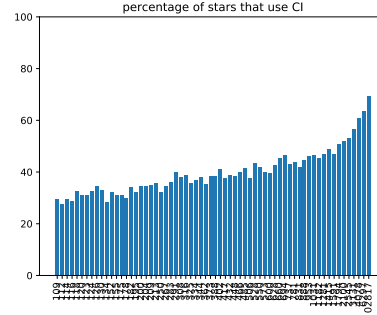Figure 5: 2016 dataset

Figure 6: In Figure 4 is the results from this research and in Figure 5 is the results from (10).
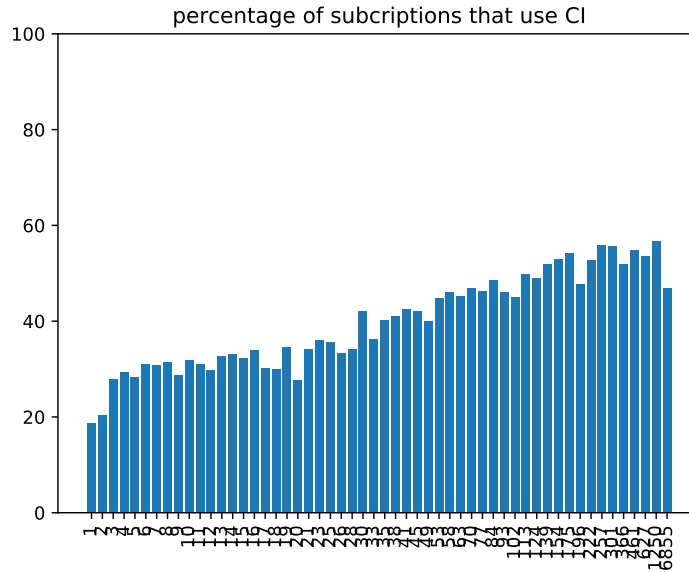


Figure 7: Subs graph

working on the project to people that want to be notified about a new release. In looking at this metric the hypothesis was that it would have a sharper rise in percentage of projects using CI per subscriber. However that was not the case overall the gradient is not as strong. There is no comparisson to (10) because their final corpus does not contain subscriber count for each project.

9

# 5 Config file results

## 5.1 What are the common errors when loading yaml configuration?

**Composer error** In the example it has two steps that are using an yaml anchor. This allows for the yaml to be referenced somewhere else. However if you define the anchor twice with the same name it causes an composer error. As you have two references using the same name 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
```

**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: ]
```

As can be seen in the Table 3 their our configuration files with yaml errors meaning that the CI for that project will not load. Yet it seems that a very small percentage of projects that have them. For example the two highest configuration types with errors are Drone (36.90%) followed by Travis (0.348%).

10

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 |

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

For Travis as it is the largest config type out of the sample by a signifacant amount 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 are the ones with errors ones that are being changed when the scraping was being done. Meaning that as the CI has been set up correctly for the other 99.632% as they are not needing to change because their our no yaml errors in it and presumbely it is doing what they intend for it to do.

11

## 5.2   How are comments used in configuration?

181   The assumption was the as continuous integeration setups can be compli-
182   cated and have edge cases. Therefore comments would be used to describe
183   and handle that complexity.

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

In order to pick up on all these dif-
ferent types of comments. All the
CI files were parsed and then reg-
ular expressions were used to pick
on up key factors such as "note:".
Along with multiple single line com-
ments 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 to-
tal 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
```

189   Initally before we look at the comments it is important to understand
190   how the rest of the file is made up. In the graph below (Figure 8) it shows
191   how each configuration type is made up by mean of each part of the file. For
192   all the yaml based configurations lines of code and number of lines in total
193   are very close. Then for the number of commments being very very small on

average.

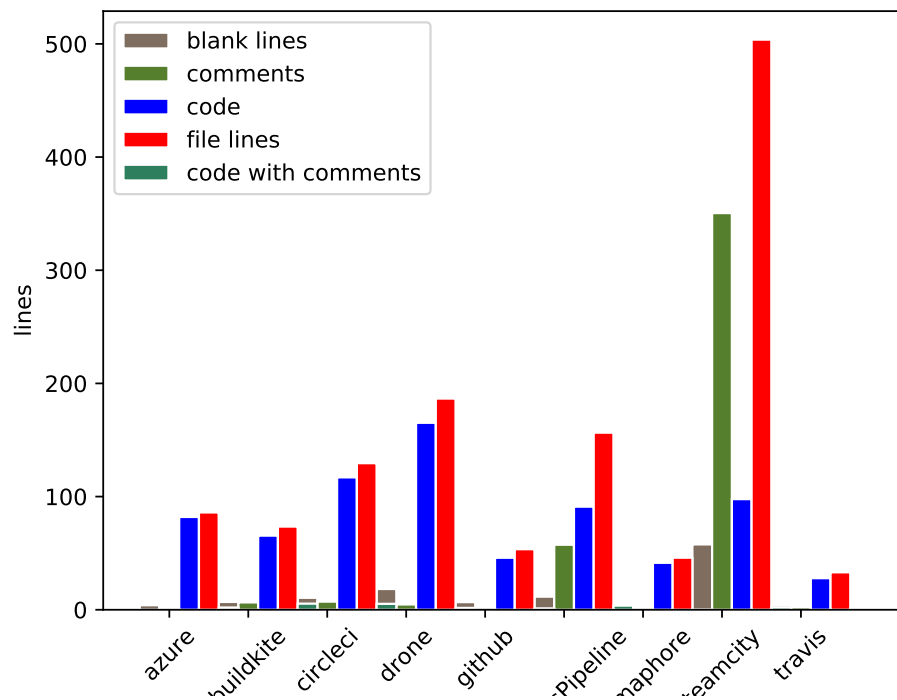In the case for Jenkins pipelines and teamcity there is a much higher usage of having code with commments.



Figure 8: Mean of line counts

Raitos:

- code: comments

- code: line total

- code: blank lines

- single line comment: multiline comment

- 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 commments. The second being useful information about the structure of the CI file such todo, note, importanat comments (e.g. //todo). In order to increase the search for this we included searching for urls and seperation comments (e.g. //===).
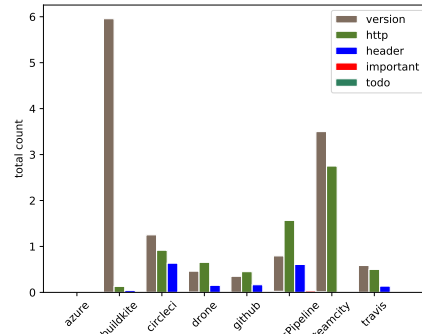


Figure 9: Comment types

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

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

14

## 5.3 Are external scripts used within the configuration?

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

The methodology we used to handle this was too look at how many bash or powershell scripts where used in CI. Using the code the parsed the yaml files for comments we were able to check do a using a regular expression for either of those files.
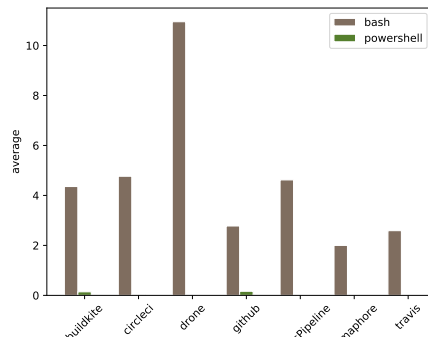
Figure 10: Comment types

|  | 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

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

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

15

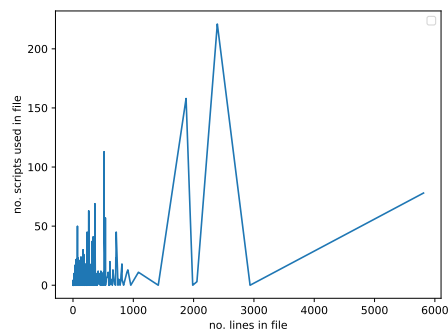trying to see if the more popular a project is affects the chances of it using
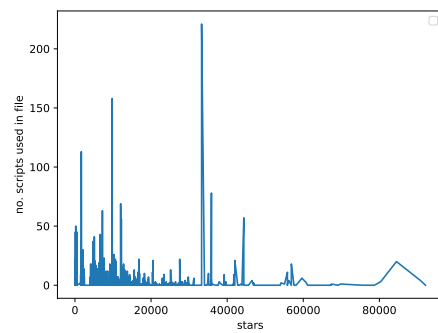CI.



Figure 12: no. scripts to no. lines



Figure 13: no. scripts to stars

percentage of
usage needed
like we had for
comments

Overall we can see that scripts are not used that much. And their no
correlation between lines of code and usage of external scripts.

# 6 Threats to validatity

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

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

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

Thirdly identifying which projects are programming projects or would have a need for CI. Based on the research (9) it is important to filter out repositories that aren't part of the question being asked. Therefore we could have looked to try and filter out github static sites and other none software based projects. However if assume a certain type of project won't be using CI then we would be introducing bais when trying to answer how CI is used. For further research better labelling of what kind of projects are which would

17

265    potentailly benefical though.

# 7  Summary

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

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

From this a better understanding of this topic could be gathered by looking into the data gathered more. As we found we were faced with a lot more questions while doing this research as we go into below.

## 7.1  Discussion and further research

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

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

19

Further research into structure could look into the naming of each part of the build process that is used. This would be interesting as it would provided insight into what terms are commonly used. As well an idea into how people plan or don't plan out their configuration files. Additionally CI systems can be designed to run on every commit to version control or only commits to certain branches. Therefore by looking at the branching regexp that are being used an better understanding of how branches are actually used in software development where CI is also used could be found out.

In addition working on pruning our dataset using methods outlined in (9).

# 8    Acknowledgement

[section ommitted]

# References

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

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

[3] Fowler,        M.        (2010).        Continous        integration.        In *https://www.martinfowler.com/articles/continuousIntegration.html*.

20

321 [4] Gallaba, K. and McIntosh, S. (2018). Use and Misuse of Continuous In-
322 tegration Features: An Empirical Study of Projects that (mis)use Travis
323 CI. *IEEE Transactions on Software Engineering*, pp. 1–1.

324 [5] github (2017). https://github.blog/2017-11-07-github-welcomes-all-ci-
325 tools/. In github.com, ed., *github welcomes all ci tools*.

326 [6] GitHub (2020). github filename search for wrecker.yml files. In *github*
327 *filename search for wrecker.yml files*.

328 [7] Gitlab (2020). https://docs.gitlab.com/ee/api/lint.html. In *Gitlab docs*.

329 [8] Jenkins (2020). https://jenkins.io/doc/book/pipeline/development/. In
330 *Jenkins documentation*.

331 [9] Kalliamvakou, E., Gousios, G., Blincoe, K., Singer, L., German, D. M.
332 and Damian, D. (2014). The promises and perils of mining GitHub.
333 Hyderabad, India: Association for Computing Machinery, MSR 2014,
334 pp. 92–101.

335 [10] Michael Hilton, K. H., Timothy Tunnell, Marinov, D. and Dig, D.
336 (2016). Usage, costs, and benefits of continuous integration in open-
337 source projects | Proceedings of the 31st IEEE/ACM International Con-
338 ference on Automated Software Engineering.

339 [11] mscoutermarsh Github memeber of staff (2019). Yaml validator for
340 github actions possible expansion of variables. In *YAML validator for*
341 *Github Actions possible expansion of variables*.

342 [12] Rahman, A., Mahdavi-Hezaveh, R. and Williams, L. (2019). A system-
343 atic mapping study of infrastructure as code research. *Information and*
344 *Software Technology*, 108, pp. 65–77.

345 [13] Shahin, M., Ali Babar, M. and Zhu, L. (2017). Continuous Integration,
346 Delivery and Deployment: A Systematic Review on Approaches, Tools,
347 Challenges and Practices. *IEEE Access*, 5, pp. 3909–3943.

348 [14] Sharma, T., Fragkoulis, M. and Spinellis, D. (2016). Does Your Config-
349      uration Code Smell? In *2016 IEEE/ACM 13th Working Conference on*
350      *Mining Software Repositories (MSR)*, pp. 189–200, iSSN: null.

351 [15] travis (2017). travis yaml (old repository). In *https://github.com/travis-*
352      *ci/travis-yaml/*.

353 [16] travis    (2020).    travis    yaml    new    implementation.    In
354      *https://github.com/travis-ci/travis-yml/*.

355 [17] Vasilescu, B., Yu, Y., Wang, H., Devanbu, P. and Filkov, V. (2015).
356      Quality and productivity outcomes relating to continuous integration
357      in GitHub. Bergamo, Italy:  Association for Computing Machinery,
358      ESEC/FSE 2015, pp. 805–816.

359 [18] Wrecker and Oracle (2018). Wrecker ci development blog. In *Wrecker*
360      *CI development blog*.