# The Hidden Costs of Automation: An Empirical Study on GitHub Actions Workflow Maintenance

Pablo Valenzuela-Toledo[1,4], Alexandre Bergel[2], Timo Kehrer[1], Oscar Nierstrasz[3]

[1]*Software Engineering Group, University of Bern*, Bern, Switzerland
[2]*RelationalAI*, Bern, Switzerland
[3]*feenk GmbH*, Wabern, Switzerland
[4]*Universidad de La Frontera*, Temuco, Chile

*Abstract*—**GitHub Actions (GA) is an orchestration platform that streamlines the automatic execution of software engineering tasks such as building, testing, and deployment. Although GA workflows are the primary means for automation, according to our experience and observations, human intervention is necessary to correct defects, update dependencies, or refactor existing workflow files. In fact, previous research has shown that software artifacts similar to workflows, such as build files and bots, can introduce additional maintenance tasks in software projects. This suggests that workflow files, which are also used to automate repetitive tasks in professional software production, may generate extra workload for developers. However, the nature of such effort has not been well studied. This paper presents a large-scale empirical investigation towards characterizing the maintenance of GA workflows by studying the evolution of workflow files in almost 200 mature GitHub projects across ten programming languages. Our findings largely confirm the results of previous studies on the maintenance of similar artifacts, while also revealing GA-specific insights such as bug fixing and CI/CD improvement being among the major drivers of GA maintenance. A direct implication is that practitioners should be aware of proper resource planning and allocation for maintaining GA workflows, thus exposing the "hidden costs of automation." Our findings also call for identifying and documenting best practices for such maintenance, and for enhanced tool features supporting dependency tracking and better error reporting of workflow specifications.**

*Index Terms*—**Continuous Integration, GitHub Actions, Workflow Maintenance, Empirical Study**

## I. INTRODUCTION

The software industry has widely adopted Continuous Practices (CP) such as Continuous Integration and Delivery (CI/CD) to automate software engineering tasks [1]. The premise driving these practices suggests that they can help minimize integration issues, enable frequent integration, automatically deploy changes, and speed up feedback loops to software developers [2]–[4].

Within the GitHub ecosystem, developers can automate software engineering tasks through GitHub Actions (GA), a widely adopted tool for implementing CP [1]. GA uses `YAML` workflow files as the building blocks for automating software engineering tasks. Developers can specify these files to perform jobs when a specific trigger event occurs throughout the definition of executable processes. The workflows can execute jobs concurrently, sequentially, or following a specific order defined by the developer. Each job comprises sequentially executed steps that embody commands, similar to bash scripts, to perform a common CI/CD task. Events that trigger workflows include pushing changes to a repository, generating new pull requests, and regularly scheduled events.

Although workflow files are the primary means for specifying automation, human intervention is necessary to correct defects, update dependencies, or refactor existing workflow files. For instance, the pull request titled "use github.ref rather than github.event.ref in deploy.yml" of the `tldraw` GitHub repository[1] illustrates an instance of CI/CD enhancement. In this case, developers include Dependabot[2], a tool to automatically update dependencies, as part of the CI/CD process [5]. This means that automating the dependency update process comes with a cost: an additional workload, as developers need to set up the new tool in conjunction with the project's workflows. This example aligns with previous studies showing that integrating bots to automate dependency updates generates maintenance overhead for developers [6].

More generally, previous research has shown that software artifacts similar to workflows, such as build files and bots, can introduce additional maintenance tasks in software projects [6], [7]. This suggests that workflow files, which are also used to automate repetitive tasks in professional software production, may generate extra workload for developers. However, the nature of such effort has not been well studied.

We present a large-scale empirical study towards characterizing the effort needed to maintain GA workflow files. We review the change history of 183 projects hosted on GitHub using GA across 10 programming languages. Inspired by previous work on similar kinds of development artifacts [7], [8], our study focuses on three primary aspects: (1) the number and size of workflow files distributed across multiple projects as well as the frequency of changes (*i.e.*, file churn rate), (2) workflow coupling, which denotes how frequently source code changes require workflow changes, and (3) workflow ownership, which involves identifying the developers responsible for these changes.

On the one hand, as expected, we find that workflow files constitute a rather small proportion of a project's total amount of files, and that they are changed less frequently than other kinds of source code files. On the other hand, however, the average size of workflow files is comparable

---

[1]https://github.com/tldraw/tldraw/pull/2495
[2]https://github.com/dependabot/dependabot-core

to those of production and test code files, and they are far from being stable in terms of evolution. Developers tend to update production and test code when workflow files are modified, primarily due to bug fixing and CI/CD improvement. Changes are often conducted by only a few developers, but these developers also take the roles of production and test code developers.

The implications of our study are manifold. Practitioners should be aware that while automation through GA workflows helps to streamline the CI/CD process [9], it also incurs significant maintenance efforts expressed in workflow modifications. This awareness can lead to better planning and resource allocation for maintaining GA automation workflows. Moreover, like any other artifact in professional software production, GA workflows should be designed with maintainability in mind. From a researchers' perspective, our findings highlight the need to identify and document best practices for maintaining GA workflows. In addition, supporting tools should be equipped with features that facilitate easier updates and maintenance, particularly regarding dependency tracking and improved error reporting.

In summary, this paper provides the following contributions:

- A large-scale empirical study providing quantitative results on the volume, evolution and ownership of GA workflow files as well as their logical coupling with other source code files.
- A detailed qualitative investigation that yields a taxonomy comprising four major reasons for the logical coupling of workflow files and other source code files.
- A dataset curated from 183 mature GitHub projects facilitating further research on studying the evolution and maintenance of GA workflows.

The data and tools used in our study are publicly available for the sake of replication and reproduction [10].

## II. RESEARCH QUESTIONS

The goal of this study is to characterize the maintenance of GA workflow files. We propose five concrete research questions for guiding our study, inspired by prior empirical research on the maintenance of similar artifacts [7], [8]. With these RQs, we aim to better understand (i) the *volume and evolution* of workflow files within a software project ($RQ_1$ and $RQ_2$), (ii) the *logical coupling* of workflow files and other source code files including the contributing factors ($RQ_3$ and $RQ_4$), and (iii) *workflow ownership* in terms of developers being responsible for changing workflow files ($RQ_5$):

**$RQ_1$: How many workflow files are in a project, and what are their sizes?** We aim to quantify the volume of workflow files to assess their significance across multiple projects. Determining their size serves as first high-level indicator for assessing their complexity.

**$RQ_2$: What is the rate of changes in workflow files?** We seek to analyze the rate of change in workflow files (*i.e.*, file churn rate [11]) by measuring the fraction of workflow files that have undergone changes over a specified period. This is motivated by previous research by Shin *et al.*, and Jiang and

Adams [8], [12], which underscores that a high churn rate may indicate the presence of flawed or potentially vulnerable files, requiring increased maintenance effort.

**$RQ_3$: To what extent are workflow files logically coupled with other source code files?** We want to examine the degree of interdependence between workflow files and other file categories using logical coupling [13]. Previous studies indicate that artifacts similar to workflows, such as build and Infrastructure as Code files, exhibit low coupling with source code files [7], [8]. In this context, we aim to determine whether this is different for workflow files or holds true for them as well.

**$RQ_4$: What are the reasons that contribute to workflow files being logically coupled with other source code files?** Given that some degree of logical coupling can be observed, we aim to uncover the causes of logical coupling between workflow files and other file categories. By knowing the reasons behind such coupling across multiple projects, we can assess whether it is an inherent aspect of the system architecture or if there is a common pattern of developers' behavior. We identify these reasons through a qualitative analysis of logical workflow coupling and develop a comprehensive taxonomy of the contributing factors.

**$RQ_5$: How is workflow ownership distributed among different roles of developers?** We want to determine how many developers are involved in maintaining workflow files, and how such ownership is distributed among different roles of developers. This will help us to understand how responsibilities are distributed within developer teams and whether there are potential bottlenecks or knowledge silos.

## III. STUDY DESIGN

In this section, we give an overview of our study design comprising three major phases, namely project selection, data curation, and data analysis. We structure Sections III-A and III-B according to the individual steps of each of the first two phases. For the sake of avoiding redundant descriptions, Section III-C provides a more general introduction into the analysis techniques used in our study, as they are partially reused across our five RQs. We will describe their applications in terms of the concrete analyses in more detail in Section IV.

### A. Project Selection

We collected a dataset from GitHub projects using GA. Our records include source code change logs and details about the developers responsible for the changes. To ensure the quality and relevance of our dataset, we followed the methodology proposed by Kalliamvakou *et al.* [14], splitting the selection process into four inclusion criteria: (a) initial inclusion criteria, (b) projects that use GA and specific programming languages, (c) active, and large projects, and (d) software projects that are not duplicates of each other.

*(a) Initial inclusion criteria:* We leverage the `seart-ghs` website [15] to identify software projects on GitHub [15]. To find projects with a rich change history, we selected only those with at least 500 commits. Furthermore, using the number

of stars as a proxy for a project's popularity [16], we only included projects with at least 500 stars to avoid irrelevant or toy projects. We did not include forked projects because they largely contain duplicated project histories, which would bias our analysis. To include a large number of projects, we selected those created before December 31, 2023.

*(b) Projects using popular programming languages and GA:* We considered projects from the top programming languages from 2020 to 2023 (namely JavaScript, Python, TypeScript, Java, C#, C++, PHP, C, Shell, and Ruby) due to their relevance in the GitHub community during the study period [17]. We identified projects using GA through the GitHub API because they are the focus of our study.

*(c) Active and large projects:* We aimed to include those projects demonstrating a collaborative, long-term software development process. To identify such projects, we defined thresholds using the "knee method", following the strategy outlined by Weeraddana *et al.* [18]. This strategy involved evaluating: (1) the number of commits, (2) the number of contributors, and (3) the number of stars. To evaluate the number of commits, we established a threshold of 8.151, chosen due to its proximity to the "knee" of the curve, identifying 6,867 projects. For the number of contributors, a threshold of 135 was similarly selected, also close to the curve's knee, identifying 2,417 projects. Likewise, for the number of stars, we set a threshold of 28,691, identifying 211 projects with significant popularity. From these projects, we only selected those providing a "readme" file written in English, resulting in 187 projects.

*(d) Deduplication of software projects:* Finally, we manually identified and removed duplicate and non-software projects before cloning the remaining ones. We found that the projects `AutoGPT` and `Auto-GPT` are the same non-forked repository from the `Significant-Gravitas` owner, so we excluded the second. After excluding three more projects due to cloning errors, our dataset consists of 183 diverse software projects from well-known organizations such as Google, Microsoft, and others.

### B. Data Curation

The data curation process consists of four steps: (a) Classifying projects into single-workflow and multi-workflow groups, (b) collecting commit data for each project, (c) classifying a project's files into production code, test code, and workflow files, and (d) unifying contributor identities.

*(a) Classification into single- and multi-workflow projects:* To account for different strategies of workflow modularization, we classify the projects into two groups: single-workflow, consisting of 27 projects with only one workflow file, and multi-workflow, containing 156 projects with more than one workflow file. Our study compares these groups to assess differences in their maintenance and evolution characteristics.

*(b) Collecting commit data for each project:* We collect commit data for each project by iterating through the commit history of all branches using the command `git log --topo-order`. We extract the unique commit ID (SHA),

timestamp, contributor's name, and the files modified or created for each commit. We filter the commits to include only those made from the point of introduction of the first GA-related commit up to the end of 2023. This ensures that the timeframes will be associated with the use of GA workflows.

*(c) Unifying contributor identities:* Changes in a committer's name or email address can result in incorrect attribution of commits on GitHub. To unify contributor identities, we use the GitHub-alias-merging script by Vasilescu *et al.* [19], which employs heuristics to link different aliases and email addresses of the same contributor.

*(d) Classifying a project's files into production code, test code, and workflow files:* We classify the files in each project as production, test, or workflow files, following the method of Nejati *et al.* [20]. We disregard any files that do not fall into these categories. We use file extensions, naming, and location conventions to identify file types. We focus on source code files written in the project's primary programming language, categorizing them as production or test code files. For example, in Java projects like `spring-boot`[3] and `dbeaver`[4], we consider files with the `.java` extension. Using project-specific naming and location conventions, we create regular expressions to identify test files, classifying non-matching source files as production code. We identify workflow files by their location in the `.github/workflows` directory and their `.yml` or `.yaml` extension. Unlike previous studies, we exclude build files and similar artifacts due to the complexity of identifying these across multiple programming languages (e.g., due to the heterogeneity of build systems and the various files involved as build scripts).

### C. Analysis Techniques

We use three main analysis techniques: (a) statistical tests, (b) association rules, and (c) open coding & card sorting.

*(a) Statistical tests:* We primarily used the Kruskal-Wallis and Mann-Whitney tests [22] for our statistical analyses ($RQ_1$-$RQ_5$). We applied the non-parametric Kruskal-Wallis test to determine if the distribution of a measure varied between the three file categories. If we rejected the null hypothesis—stating *"there is no significant difference between the means of the three categories"*—it indicated that at least one category had a distinct distribution of the metric under examination. For post-hoc testing, we used Mann-Whitney tests to identify which specific categories had distinct distributions. We performed these tests between every pair of file categories, applying the Bonferroni correction to the alpha value (0.05 by default in all our tests).

*(b) Association rules:* We investigate logical coupling in $RQ_3$ and $RQ_5$ through association rules [21]. Specifically, to assess the coupling relationship between different file categories, we analyze each pair $\langle A, B \rangle$ of file categories. As for quantitative assessment, we use the usual metrics known as *Support (Supp)*, *Confidence (Conf)* and *Lift*, as defined in

---

[3]https://github.com/spring-projects/spring-boot
[4]https://github.com/dbeaver/dbeaver

TABLE I: Association Rule metrics: Support, Confidence and Lift [21] (here: to assess the coupling relationship between changes in different file categories).

| Metric | Formula | Description |
|---|---|---|
| Supp(A) | $P(A)$ | Determines the frequency of changes in file category A. A value of 1 (or 100%) means that every commit changes at least one file in category A. |
| Supp(A & B) | $P(A \cap B)$ | Determines how frequently changes in file categories A and B occur together. A maximum value of 1 means that every commit changes files in categories A and B. |
| Conf(A⇒B) | $\frac{P(A \cap B)}{P(A)}$ | Determines how often changes in file category B occur along with changes in file category A. A maximum value of 1 means that whenever a file of category A is changed, at least one file of category B is changed in the same commit. |
| Lift(A⇒B) | $\frac{P(A \cap B)}{P(A)P(B)}$ | Determines how much more (or less) likely changes in file category B are to occur along with changes in file category A, compared to when no file of category A has changed. Lift = 1 means that changes in file categories A and B occur together as expected if they were independent, while Lift >1 (Lift <1) indicates a strong (weak) association A ⇒ B. |



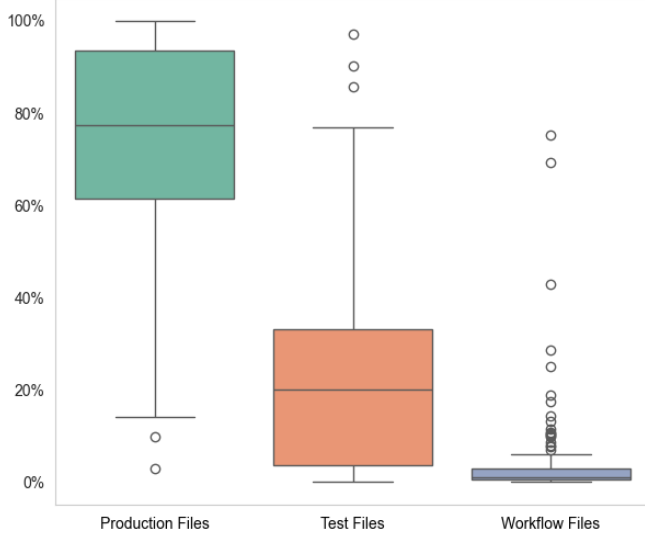Fig. 1: Mean relative frequencies of the studied file categories in the multi-workflow group.



Fig. 2: Mean file sizes in NLOC (without outliers) of the studied file categories in the multi-workflow group.

Table I for assessing frequencies of and relationships between file category changes per commit (RQ$_3$). The metrics used for assessing frequencies of and relationships between developer roles per project history (RQ$_5$) are defined analogously. Afterwards, similar to Jiang and Adams [8], we performed chi-square statistical tests to test whether the obtained confidence values were significant or not higher than expected due to chance.

*(c) Open coding & card sorting:* For RQ$_4$, we conduct open coding [23] using a randomly selected set of logically coupled commits that have a lift greater than one. Then, we apply open card sorting [24] to codify the reasons for coupled changes. Open Coding and Card Sorting are widely used techniques in software engineering useful to derive taxonomies from data [24].

## IV. ANALYSIS AND RESULTS

### A. *Relative Frequency and Size of Workflow Files (RQ$_1$)*

To answer RQ$_1$, we study the relative frequency and size of each file category (i.e., production, test, and workflow) across all projects, separated by single-workflow and multi-workflow groups. We use the Kruskal-Wallis test to evaluate the statistical significance of our findings.
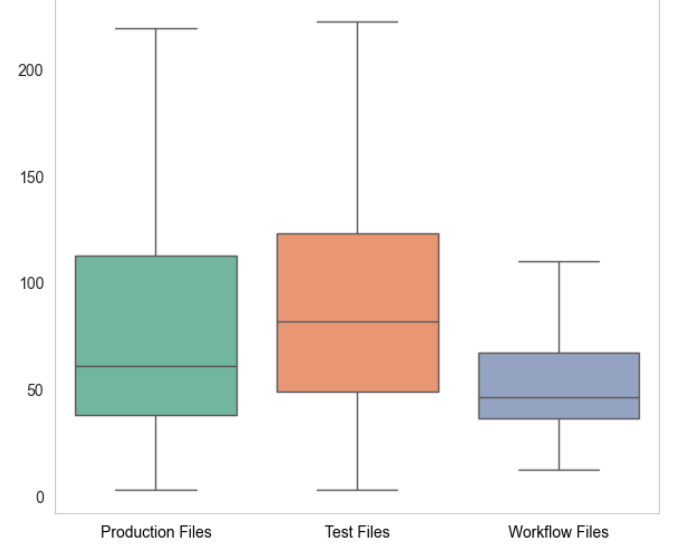
Overall, the 183 studied projects contain between 1 and 70 workflow files. The multi-workflow group has a mean relative frequency of 3.74% workflow files, while the single-workflow group has 3.05%. Production files make up the most significant proportion, with a mean relative frequency of 74.59% in the multi-workflow group and 83.31% in the single-workflow group. Test files are the second-largest category, with a mean relative frequency of 21.65% in the multi-workflow group and 13.63% in the multi-workflow group. Figure 1 shows a boxplot of the proportion of the three file categories in the multi-workflow group, reflecting these trends. We do not include a boxplot of the single-workflow group since it follows the same trends. More detailed statistics on both distributions may be found in our replication package.

As we are dealing with textual source code artifacts in our study, we measure their size in Number of Lines of Code (NLOC). Production and test files are generally larger than workflow files. In the multi-workflow group, workflow files have a mean NLOC of 60.99, while production and test files expose mean NLOC values of 103.26 and 101.37, respectively. In the single-workflow group, the mean NLOC increases to 98.93 for workflow files, 128.19 for production files, and

113.32 for test files. Figure 2 shows a boxplot of the mean file sizes in NLOC (without outliers) of the studied file categories in the multi-workflow group. We do not include a boxplot for the single-workflow group since it follows the same trends. Again, more detailed statistics on both distributions may be found in our replication package.

> **Finding 1:** Workflow files account for a small proportion of files, with a mean relative frequency of 3.74% in the multi-workflow group and 3.05% in the single-workflow group. However, their mean size of NLOC, at 60.99 in the multi-workflow group and 98.92 in the single-workflow group, is close to that of production and test files, which are of the same order of magnitude.

### B. Rate of Changes in Workflow Files (RQ2)

To assess the rate of changes in workflow files, we first determined the amount of changed files per month of each project. Subsequently, to enable comparisons over time, we normalize the number of changed files by dividing it by the total number of files in the corresponding category for that month, yielding the proportion of changed files in terms of file churn rate [11]. For each project, we then calculate the average proportion of changed files per month, and we study the distribution of this average across all projects. Furthermore, we conducted a Kruskal-Wallis test and post-hoc tests to examine the statistical significance of our results.

The analysis of both single- and multi-workflow groups reveals that the monthly proportion of changed workflow files is significantly lower than for production and test code files. In the single-workflow group, workflow files have a mean value of 6.8%, while in the multi-workflow group, the mean value is 8.2%. Production files have mean values of 84.53% and 74.12% for the single-workflow and multi-workflow groups, respectively. Test code files have mean values of 10.75% and 18.90%, correspondingly.

Figure 3 illustrates the distribution of the rate of changes for the workflow, test, and production file categories. The solid lines represent the median percentages, and the dotted lines represent the first and third quartiles. The figure shows that workflow files have a relatively narrow distribution with a concentration around lower change rates, indicating they experience fewer changes overall. Additionally, the plot for workflow files reveals a single peak, suggesting a consistent rate of changes, mainly concentrated at lower proportions. In contrast, production files show a wide distribution with a higher concentration around the upper change rates and a pronounced peak, reflecting their more frequent updates. Test files exhibit a distribution between the workflow and production files, indicating moderate change frequencies.

A Kruskal-Wallis test confirmed significant differences in the distributions of the monthly change percentages among the different file categories in both groups, with extremely low p-values of 5.06e-33 for the multi-workflow group and

0.0 for the multi-workflow group. Further Mann-Whitney post-hoc tests showed statistically significant differences between Production and Workflow files, with p-values of 1.77e-14 and 0.0, respectively. The tests also showed significant differences between Production and Test files, with p-values of 2.04e-29 and 0.0 each. Additionally, there were significant differences between Test and Workflow files, with p-values of 4.76e-11 for the multi-workflow group and 1.32e-255 for the multi-workflow group.

> **Finding 2:** Workflow files change less frequently than production and test code files. In the multi-workflow group, a mean of 8.2% of the workflow files undergo at least one monthly change, while repositories with a single-workflow account for a mean 6.9%. In the same time period, production files have a mean change proportion of 74.12% in the multi-workflow group and 84.53% in the single-workflow group. Similarly, test files change 18.90% of the time in multi-workflow and 10.75% of the time in repositories with a single workflow.

### C. Quantitative Analysis of Logical Workflow Coupling (RQ3)

The results of our quantitative analysis indicate that workflow files are the least frequently changed, as evidenced by the low *Support* metric values (1.40% in the multi-workflow group and 3.13% in the single-workflow group)(See Table II). While modified more often than workflow files, test files still show relatively low-frequency changes (*Support* values of 19.07% in the multi-workflow group and 32.57% in the single-workflow group). On the other hand, production files exhibit the highest frequency of changes, with *Support* values of 92.57% in the multi-workflow group and 88.55% in the single-workflow group.

Regarding the relationship between workflow and test files (*Supp(Workflow & Test)*), the *Support* is 0.06% in the single-workflow group and 0.23% in the multi-workflow group, indicating that commits involving changes in both categories are extremely rare. For the relationship between workflow and production files (*Supp(Workflow & Production)*), the *Support* is slightly higher, with values of 0.17% in the single-workflow group and 0.53% in the multi-workflow group.

Changes in workflow files are rarely associated with changes in other categories, and when they are, they are more likely to coincide with changes in production files rather than test files. Regarding the relationship between workflow and test files (*Conf(Workflow ⇒ Test)*), the *Confidence* is 4.60% in the multi-workflow group and 7.40% in the single-workflow group, indicating that if a commit involves changes in workflow files, there is a 4.60% chance in the multi-workflow group and a 7.40% chance in the single-workflow group that it also involves changes in test files. For the relationship between workflow and production files (*Conf(Workflow ⇒ Production)*), the *Confidence* values are higher, with 12.26% in the single-workflow group and 17.00% in the multi-workflow
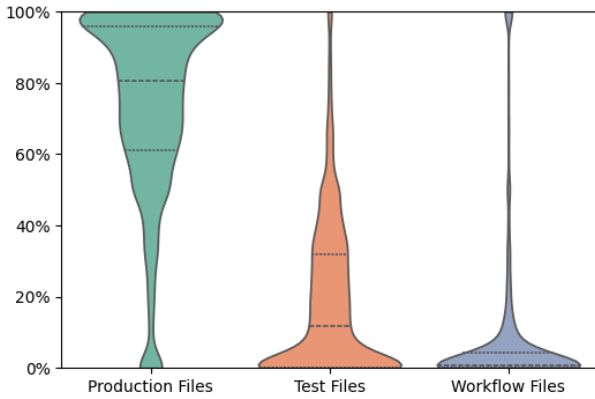
Fig. 3: Distribution of the rate of changes for the workflow, test, and production file categories. We calculated the monthly rate of change by normalizing the number of changed files against the total files in each category, then averaging this proportion per month for each project.

group. If a commit involves changes in workflow files, there is a 12.26% chance in the single-workflow group and a 17.00% chance in the multi-workflow group that it also involves changes in production files.

Changes in workflow files tend to occur independently of changes in both test and production files, as indicated by the *Lift* metrics. For the relationship between workflow and test files (*Lift(Workflow ⇒ Test)*), the *Lift* is 24.12% in the single-workflow group and 22.71% in the multi-workflow group. These values indicate that changes in workflow files reduce the likelihood of changes in test files to 24.12% and 22.71% of what would be expected by random chance, respectively. This suggests a negative association, meaning that changes in workflow files rarely coincide with changes in test files.

For the relationship between workflow and production files (*Lift(Workflow ⇒ Production)*), the *Lift* values are 13.25% in the single-workflow group and 19.20% in the multi-workflow group. These values show that when workflow files change, production files rarely change at the same time: only 13.25% of the time in the single-workflow group and 19.20% of the time in the multi-workflow group. Although the *Lift* values for workflow and production files are slightly higher than those for workflow and test files, they still indicate a negative association.

The combination of workflow and test yields a chi-square value of 4502.59 with a p-value of $< 0.0001$, suggesting a strong correlation between these two categories. Furthermore, the combination of workflow and production reveals an even stronger association, with a chi-square value of 78,799.33 and a p-value of $< 0.0001$. Similarly, the combination of test and production also shows a highly significant correlation, with a chi-square value of 58,511.84 and a p-value of $< 0.0001$. These results imply that there are significant relationships between these categories and that they are not independent of each other.

TABLE II: File category changes per commit: Support, Confidence, and Lift metrics for single- and multi-workflow groups.

| Metric | Category Combination | Single | Multi |
|---|---|---|---|
| Support | Workflow | 0.0140 | 0.0313 |
| | Test | 0.1907 | 0.3257 |
| | Production | 0.9257 | 0.8855 |
| | Workflow & Test | 0.0006 | 0.0023 |
| | Workflow & Production | 0.0017 | 0.0053 |
| Confidence | Workflow ⇒ Test | 0.0460 | 0.0740 |
| | Workflow ⇒ Production | 0.1226 | 0.1700 |
| Lift | Workflow ⇒ Test | 0.2412 | 0.2271 |
| | Workflow ⇒ Production | 0.1325 | 0.1920 |

**Finding 3:** The coupling between workflow files and other file categories is minimal and weak. For the relationship between workflow and test files, there is a 4.60% chance in the single-workflow group and a 7.40% chance in the multi-workflow group that a commit involving workflow files also changes test files. For workflow and production files, there is a 12.26% chance in the single-workflow group and a 17.00% chance in the multi-workflow group. In the multi-workflow group, there is a slightly higher but still weak interdependence between workflow and production files.

### D. Qualitative Analysis of Logical Workflow Coupling (RQ₄)

In addition to RQ₃, we also conducted a qualitative analysis of commits that contribute to the highest levels of logical coupling, as indicated by their Lift values (see Table II). We applied open coding and card sorting for this analysis. We selected 200 commits with the highest Lift scores (100 from the single-workflow group and 100 from the multi-workflow group) to investigate why workflow changes were tightly coupled, similar to the previous work by Jiang et al. [8].

We followed a structured approach to card sorting, carried out in three distinct phases as recommended by Zimmermann [24]: First, in the *preparation phase* we use descriptive coding [25]. We set up a collection of cards to provide a detailed understanding of the workflow coupling. Each commit was meticulously described based on the GitHub single commit view, including the reasons for and locations of the changes. Second, in the *execution phase*, we labelled and organized each card into relevant groups with clear, descriptive categories. Finally, in the *analysis phase*, we created abstract hierarchies without predefined categories to identify broader trends and general categories.

To ensure accuracy, the description criteria were tested by three participants: two undergraduate students, one graduate student and the first author. Each of them was given 100 commits to describe the reason for coupling. These descriptions were then reviewed and refined collaboratively to establish a final structure. During the sorting process, we used spreadsheets to facilitate easier handling and better understanding of each commit. A hierarchy was created manually to organize a set of categories.

TABLE III: Main Reasons for Logical Coupling.

| Main Reasons for Logical Coupling | Specific reason | Workflow & Test | Workflow & Production |
|---|---|---|---|
| Management and Configuration | Configuration and Infrastructure | 1% | 3% |
| | Dependency and Build Management | 7% | 6% |
| | Merging and Synchronization | 1% | 5% |
| Development and Enhancement | CI/CD Improvement | 61% | 33% |
| | Code Quality and Maintenance | 7% | 9% |
| | Features and Enhancements | 6% | 11% |
| Bug Fixes and Compatibility | — | 12% | 18% |
| Testing and Documentation | Testing | 4% | 12% |
| | Documentation | 1% | 3% |
| Total | | 100% | 100% |

As a result, we built a taxonomy consisting of nine reasons for logical coupling including a total of 200 labelled commits (Table III). We grouped them into four main categories: (1) Management and Configuration, (2) Development and Enhancement, (3) Bug Fixes and Compatibility, and (4) Testing and Documentation.

*1) Management and Configuration*

This category involves tasks related to setting up and managing the project's environment, including configuring infrastructure, handling dependencies, and synchronizing code across different branches (9% of workflow & test, 13% of workflow & production). These activities ensure that the foundational aspects of the project are well-organized and functioning correctly. Three reasons for coupling are related to this category: (i) Configuration and Infrastructure, (ii) Dependency and Build Management, and (iii) Merging and Synchronization.

*1.a) Configuration and Infrastructure* encompasses modifications related to setting up and maintaining the infrastructure and configuration necessary for the project. Proper configuration management is essential for ensuring a stable and efficient development environment. The commit "fix(dev): Add .python-version back" illustrates this by addressing issues related to the Python version management within the development environment.[5]

*1.b) Dependency and Build Management* involves managing the project's dependencies, orchestrating the build process, and ensuring that all necessary components are correctly integrated. The commit "Revert Switch to Debian 11 (bullseye) as base for our dockerfiles" involves reverting a change to Debian 11 (bullseye) that caused all pull requests to fail.[6] The revert restored stability, ensuring dependencies and configurations returned to a reliable state.

*1.c) Merging and Synchronization* cover tasks related to merging code from different branches and synchronizing changes across the codebase. The commit "Merge branch '6.4' into 7.0" exemplifies this category by integrating changes from branch 6.4 into 7.0, which involves combining contributions from different parents and resolving potential conflicts.[7] It ad-dresses specific issues, such as fixing the silencing of the wait command for `sigchild-enabled` binaries, and updates multiple files to ensure consistent behaviour and compatibility.

*2) Development and Enhancement*

This category encompasses tasks related to improving the project's functionality and performance. It involves introducing new features, enhancing existing ones, and ensuring the software efficiently meets user needs and standards (84% of workflow & test, 53% of workflow & production). Three factors contributing to coupling are associated with this category: (i) CI/CD Improvement, (ii) Code Quality and Maintenance, and (iii) Features and Enhancements.

*2.a) CI/CD Improvement* relates to optimizing and automating the continuous CI/CD Improvement. It includes setting up and refining pipelines and ensuring smooth and efficient code integration. The commit "Use npm v7 with workspaces for local development and testing" explains how developers upgrade the project to use `npm` v7 with `workspaces`, enhancing continuous integration and deployment processess.[8]

*2.b) Code Quality and Maintenance* involve activities that ensure the software remains functional, efficient, and easily understood over time. This includes practices such as refactoring code to improve readability and performance and writing and maintaining tests to catch bugs early. The commit "tools: automate histogram update" enhances the CI/CD workflow by automating the update process for the project's histogram dependency.[9] It adds a new script, update-`histogram.sh`, to the `tools` directory, which fetches and updates the histogram to its latest version from GitHub releases.

*2.c) Features and Enhancements* involve activities to add new functionality or improve existing features in a software project. This includes designing and implementing new features to meet user needs, optimizing performance, enhancing user interfaces, and extending the software's capabilities. For example, the commit "Add –os and –arch flags to readall" exemplifies this category by adding the `--os` and `--arch` flags to the `brew readall` command, allowing users to read using specified operating systems and CPU architectures.[10]

---

[5]https://bit.ly/3RDUhR8
[6]https://bit.ly/3VEq99g
[7]https://bit.ly/45AHXqG

[8]https://bit.ly/3RBzSMw
[9]https://bit.ly/3VBIaFl
[10]https://bit.ly/4eyYX4E

TABLE IV: Developer roles per project history: Support and Confidence metrics for single- and multi-workflow groups.

| Metric | Category Combination | Single | Multi |
|---|---|---|---|
| Support | Workflow | 0.0435 | 0.0486 |
| | Test | 0.2150 | 0.4664 |
| | Production | 0.9746 | 0.9636 |
| | Workflow & Test | 0.0165 | 0.0224 |
| | Workflow & Production | 0.0155 | 0.0225 |
| | Test & Production | 0.1751 | 0.4156 |
| Confidence | Workflow $\Rightarrow$ Test | 0.3809 | 0.4609 |
| | Workflow $\Rightarrow$ Production | 0.3571 | 0.4637 |
| | Test $\Rightarrow$ Workflow | 0.0771 | 0.0480 |
| | Production $\Rightarrow$ Workflow | 0.0159 | 0.0233 |

*3) Bug Fixes and Compatibility*

This category refers to activities dedicated to identifying and fixing bugs, as well as ensuring that the project remains compatible with different environments and dependencies. These tasks are crucial for maintaining the stability and reliability of the project (12% of Workflow & Test, 18% of Workflow & Production). The commit "alt-svc: enable by default", illustrates this category: `altsvc` support by default in curl and removes the unused `CURLALTSV CIMMEDIATELY` option.[11] The update involves modifications across 27 files, including changes to configuration files for CI pipelines, such as `.azure-pipelines.yml`, and `.github/workflows/macos.yml`.

*4) Testing and Documentation*

This category focuses on tasks related to ensuring the reliability of the code through testing and improving project documentation (5% of Workflow & Test, 15% of Workflow & Production). Two factors contributing to coupling are associated with this category: (i) Testing, and (ii) Documentation.

*4.a) Testing* refers to modifications that enhance the testing framework or add new tests to ensure more comprehensive and effective software validation. This includes updates to testing scripts, configurations, or the addition of new test cases. The commit titled "fix: nschematics install Windows" addresses issues with installing the schematics on Windows systems.[12] The changes include adding smoke tests for Angular Schematics on multiple operating systems (Ubuntu, Windows, macOS), updating `Node.js` setup in GA, and modifying various files to support the Windows installation process.

*4.b) Documentation* refers to the changes to the documentation within the software project. This includes updates, additions, or deletions to ensure that the documentation is accurate and up-to-date. The commit "Change supported PyPy versions to 3.9 and 3.10", illustrates the updates the the ReadTheDocs (RTD) URL for the coverage documentation is updated to point to the latest version.[13] There are also adjustments to the GA workflow, specifically for testing on `PyPy` versions and ensuring the correct system libraries are installed based on the Python version being used.

[11]https://bit.ly/4cuZ4fK
[12]https://bit.ly/4exe8LJ
[13]https://bit.ly/3xsxYHi

**Finding 4:** Given that some degree of logical coupling of workflow files and other file categories can be observed, our results identified key reasons for this. The main contributors are *Management and Configuration* tasks, with significant activities including *Dependency and Build Management*, and *Merging and Synchronization*. The primary drivers are *Development and Enhancement* activities, particularly *CI/CD Improvement* processes. *Bug Fixes and Compatibility* also play a critical role, along with *Testing and Documentation* efforts.

*E. Workflow Ownership Distribution (RQ5)*

Identifying ownership for each category involves examining the commits made to the repository. We look at the author of each commit to determine their role in the development process. If a commit modifies a workflow file, then the author of that commit is considered a workflow developer. An author can have multiple roles, such as being both a workflow developer and a production code developer, even for the same commit.

We compute the same metrics as for $RQ_3$, this time for the ownership. For example, *Supp(Workflow)* indicates the percentage of developers changing workflow files from the total number of developers. *Supp(Workflow & Production)* is the percentage of developers changing workflow and production files from the total number of developers. The *Conf(Workflow ⇒ Production)* represents the percentage of workflow developers who also modify production files, in relation to the total number of developers who have made at least one change to a workflow file.

Based on the *Support* results, workflow developers have the lowest proportion among all developers, with only 4.35% in the single-workflow group and 4.86% in the multi-workflow group being involved in changes to workflow files (Table IV). These results could indicate that workflow changes are more specialized tasks handled by a smaller subset of developers. In contrast, production code developers are the most common among all developers, with 97.46% in the single-workflow group and 96.36% in the multi-workflow group making changes to production files.

Regarding the combination of workflow and test file changes, 1.65% of developers in the single-workflow group and 2.24% in the multi-workflow group are involved in both. These low results indicate that it is relatively uncommon for developers to work on both workflow and test files. The infrequency of these combined changes suggests that workflow and test tasks are often handled separately, with few developers integrating their efforts across these areas.

Similarly, 1.55% of developers in the multi-workflow group and 2.25% in the single-workflow group are involved in the combination of workflow and production file changes. These low percentages hint that it is uncommon for developers to work on both workflow and production files. The rarity of these combined changes suggests that workflow and produc-

tion tasks are typically managed independently, with limited overlap in developer responsibilities.

In terms of *Confidence*, 38.09% of developers in the single-workflow group who changed workflow files also changed test files, compared to 46.09% in the multi-workflow group *(Conf(Workflow ⇒ Test))*. Although it is relatively uncommon for workflow developers to also be involved in test file changes, these confidence values are moderately high in this context. The higher confidence in the multi-workflow group indicates that developers in this group are more likely to integrate their workflow changes with test changes, highlighting a more collaborative and thorough development process.

Most developers do not modify workflow files, suggesting a specialization in these tasks. The considerable values of *Conf(Production ⇒ Workflow)* and *Conf(Test ⇒ Workflow)* indicate that while workflow developers are relatively few, those involved in production and test activities frequently engage with workflow changes. Specifically, 15.90% of production code developers in the single-workflow group and 23.30% in the multi-workflow group also change workflow files, while 7.71% of test developers in the single-workflow group and 4.80% in the multi-workflow group also change workflow files.

> **Finding 5:** Few developers work on workflow files, with only 4.35% in the single-workflow group and 4.86% in the multi-workflow group. However, a notable proportion of developers who work on workflow files also work on test files (38.09% in the multi-workflow group and 46.09% in the multi-workflow group) and production files (35.71% in the single-workflow group and 46.37% in the multi-workflow group).

## V. DISCUSSION

In this section, we discuss the findings presented in the previous section and provide a set of practical implications.

*(a) RQ₁ & RQ₂.* Compared to production and test code files, workflow files represent only a minor fraction of the overall project files, underscoring their specialized role in managing CI/CD pipelines and automation tasks. Additionally, workflow files change significantly less frequently than production and test files. This characteristic is likely due to their specific and well-defined roles in automating CI/CD processes, which do not change as dynamically as production and test code. Nonetheless, they are far from being stable, which confirms our hypothesis that manual effort is needed for maintaining workflow files over time. Although workflow files are typically smaller than production code files, their size is not drastically different from test files, especially in the multi-workflow group where the means are relatively close. That is, workflow files being designed to automate specific tasks still contain a substantial number of lines of code.

*Implication 1:* Given workflow files' substantial size and critical role in CI/CD processes, practitioners should allocate resources accordingly. Despite their smaller proportion, these files require careful maintenance. Practitioners should dedicate sufficient time and personnel to ensure the stability and efficiency of workflows, as this directly impacts the overall development and deployment process.

*(b) RQ₃ & RQ₄.* Although the confidence of a workflow file change coinciding with a production file change is relatively low, and even lower with test files, workflow files exhibit minimal coupling with these file categories. A thorough investigation of commits that contribute to the logical coupling reveals that they are primarily updated for CI/CD improvements and bug fixes.

*Implication 2:* By categorizing the logical coupling between workflow files and other file categories, our taxonomy reveals that current tools lack robust error reporting and support features. This classification identifies areas where tools fail to meet users' needs, revealing the need for researchers to develop advanced tools to address these gaps.

*(c) RQ₅.* Our analysis finds that few developers work on workflow files. However, confidence values reveal that a notable proportion of developers who work on workflow files also work on test and production files, suggesting that workflow maintenance requires cross-functional knowledge.

*Implication 3:* Practitioners should be aware of these dependencies and plan their development and maintenance activities to account for the interconnected nature of changes in workflow files. Organizations should ensure that their teams include developers skilled in both workflow management and general software development.

## VI. THREATS TO VALIDITY AND LIMITATIONS

### A. Internal Validity

*Selection Bias:* The selection of projects based on criteria like the number of commits and stars may introduce selection bias, potentially not representing the entire population of GA workflow projects. We mitigated this by using a systematic selection methodology and including a diverse set of projects across multiple programming languages.

*Measurement Error and Categorization Accuracy:* Errors or inconsistencies in the repositories' metadata, commit messages, or file histories could affect our results. Additionally, misclassification of file changes into production, test, and workflow categories could occur due to naming conventions and file locations. We used established tools and methods for data curation and analysis, and a rigorous methodology including regular expressions and manual checks, to ensure consistency and accuracy.

*Confounding Variables and Temporal Effects:* Factors such as developer experience, project size, and complexity could influence our findings. Additionally, development practices and technologies may have evolved over the four-year study period. To address these issues, we generalized our results across multiple projects and languages, reducing the impact of any single variable. We also included a broad time frame to capture trends and changes over time.

*Developer Attribution and File Classification:* Unifying contributor identities can be challenging due to aliasing and

changes in email addresses or usernames, and errors in classifying files into production, test, and workflow categories can lead to incorrect conclusions. We used heuristics to merge aliases, manual verification for accurate developer attribution, followed a standardized classification scheme, and performed manual checks to verify accuracy.

### B. External Validity

*Platform and Ecosystem Specificity:* Our study focuses exclusively on GA workflows, which may not apply to other CI/CD platforms such as GitLab CI, Travis CI, or Jenkins. Additionally, projects hosted on GitHub may exhibit unique characteristics compared to those on other platforms. To mitigate this, we suggest future research replicate the study across different CI/CD platforms and hosting environments to verify the generalizability of our results.

*Project and Domain Diversity:* Although we included a diverse set of projects across various programming languages, certain project types or domains may still be underrepresented. We addressed this by systematically selecting projects from a variety of domains and ensuring broad representation.

### C. Construct Validity

We assume that studying changes is a good proxy to assess the maintenance of a software artifact. This assumption may not always hold true, as even simple and small changes in a workflow file could result from significant cognitive effort. This simplification was necessary because we aimed to measure across multiple repositories, which is feasible when considering changes rather than the broader concept of maintenance. To mitigate, we performed a qualitative analysis of a sample of commits to better understand the context and complexity of the changes. By examining the reasons for changes we aim to provide a more nuanced understanding of the maintenance activities.

### VII. Related Work

Prior empirical studies have focused on various aspects of maintenance and evolution related to Continuous Integration and Continuous Delivery (CI/CD) practices, essential for modern development using tools like Jenkins, Travis, CircleCI, and GA. Jiang and Adams as well as Adams and McIntosh explored CI configurations and build maintenance efforts, highlighting challenges in maintaining scripts and configurations [7], [8]. Gallaba *et al.* provided a comprehensive overview of CI/CD practices based on eight years of data [26]. Studies on the `Greenkeeper` dependency bot also emphasize the need to evaluate automation tools' benefits and costs [6]. Zhao et al. [27] and Cassee et al. [28] examined how the introduction of the Travis CI tool affected development practices. However, none of the above mentioned studies specifically focused on GA.

Research shows that GA, introduced in 2019, quickly replaced Travis due to its integration and ease of use [29]. Rostami Mazrae *et al.* found modifications to be the most common GA workflow changes, revealing hidden maintenance

costs and broader implications [30]. Delicheh et al. [31] provide preliminary insights into the dependencies of GA actions. Decan et al. [32] have studied the extent to which automation workflows are outdated with respect to updated GA actions, concluding that better policies and tooling are needed to keep workflows up-to-date. Valenzuela-Toledo and Bergel also identified the need for better tooling for GA workflows [33]. The research by Kinsman *et al.* examined how developers use GA and how activity indicators change after its adoption [34]. Saroar and Nayebi surveyed developers to understand their perceptions of GA, revealing motivations, decision criteria, and challenges [35]. Zhang *et al.* conducted a large-scale empirical study to create a taxonomy of GA-related problems, while Bouzenia and Pradel looked into resource usage and optimization opportunities in GA workflows [36]. Other studies, such as Koishybayev *et al.* [37], explored security risks associated with excessive privileges in GA workflows. The EGAD project [38], [39] provides a tool environment to explore and analyze GA workflows, contributing to a deeper understanding of their usage and maintenance practices. Cardoen et al. [40] developed an open-source tool designed to extract the commit histories of changes made to workflow files in GitHub repositories, along with a raw dataset that has been collected using their tool. Despite the relevance of all ot these studies, none of them specifically addressed the characterization of the maintenance of GA workflow files, which is the primary focus of our research.

### VIII. Conclusion

Automating tasks with GA has become standard practice in the software industry, enabling developers to automate building, testing, and deployment processes. However, our research reveals that implementing and maintaining these workflow files is not without additional costs. A large-scale empirical study spanning 183 GitHub projects in ten different programming languages identified several key features and challenges in maintaining GA workflow files.

Our findings confirm results from previous studies on the maintenance of similar artifacts such as Build files and Infrastructure as Code. While beneficial for efficiency, automation entails hidden costs that must be adequately managed. Practitioners must plan and allocate sufficient resources for maintaining these workflows, including identifying and documenting best practices.

Despite the strengths of our study, there are threats to construct validity that need to be considered. The scope of our analysis may not cover all aspects of workflow maintenance, and certain contextual factors unique to individual projects may influence the results. To mitigate these threats in future work, we plan to conduct more focused studies on specific aspects of workflow maintenance, exploring different perspectives and contexts. Additionally, we aim to collaborate with practitioners to gather more detailed insights and develop tailored strategies for efficient workflow management.

## REFERENCES

[1] M. Golzadeh, A. Decan, and T. Mens, "On the rise and fall of CI services in GitHub," in *2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2022, pp. 662–672.

[2] M. Fowler and M. Foemmel, "Continuous integration," 2006.

[3] J. Humble and D. Farley, *Continuous delivery: reliable software releases through build, test, and deployment automation*. Pearson Education, 2010.

[4] B. Fitzgerald and K.-J. Stol, "Continuous software engineering: A roadmap and agenda," *Journal of Systems and Software*, vol. 123, pp. 176–189, 2017.

[5] tldraw GitHub Repository, "tldraw/tldraw Pull Request #2495," https://github.com/tldraw/tldraw/pull/2980, 2024, accessed on March 20, 2024.

[6] B. Rombaut, F. R. Cogo, B. Adams, and A. E. Hassan, "There's no such thing as a free lunch: Lessons learned from exploring the overhead introduced by the Greenkeeper dependency bot in npm," *ACM Transactions on Software Engineering and Methodology*, vol. 32, no. 1, pp. 1–40, 2023.

[7] S. McIntosh, B. Adams, T. H. Nguyen, Y. Kamei, and A. E. Hassan, "An empirical study of build maintenance effort," in *Proceedings of the 33rd international conference on software engineering*, 2011, pp. 141–150.

[8] Y. Jiang and B. Adams, "Co-evolution of infrastructure and source code-an empirical study," in *2015 IEEE/ACM 12th Working Conference on Mining Software Repositories*. IEEE, 2015, pp. 45–55.

[9] M. Wessel, J. Vargovich, M. A. Gerosa, and C. Treude, "GitHub actions: the impact on the pull request process," *Empirical Software Engineering*, vol. 28, no. 6, p. 131, 2023.

[10] P. Valenzuela-Toledo, A. Bergel, T. Kehrer, and O. Nierstrasz, "The hidden costs of automation: An empirical study on GitHub actions workflow maintenance replication package," https://zenodo.org/record/12485092, May 2024, version 1.0, Accessed: 2024-06-23.

[11] N. Nagappan and T. Ball, "Using software dependencies and churn metrics to predict field failures: An empirical case study," in *First International Symposium on Empirical Software Engineering and Measurement (ESEM 2007)*. IEEE, 2007, pp. 364–373.

[12] Y. Shin, A. Meneely, L. Williams, and J. A. Osborne, "Evaluating complexity, code churn, and developer activity metrics as indicators of software vulnerabilities," *IEEE transactions on software engineering*, vol. 37, no. 6, pp. 772–787, 2010.

[13] H. Gall, K. Hajek, and M. Jazayeri, "Detection of logical coupling based on product release history," in *Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272)*. IEEE, 1998, pp. 190–198.

[14] E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. M. German, and D. Damian, "The promises and perils of mining GitHub," in *Proceedings of the 11th working conference on mining software repositories*, 2014, pp. 92–101.

[15] O. Dabic, E. Aghajani, and G. Bavota, "Sampling projects in github for MSR studies," in *18th IEEE/ACM International Conference on Mining Software Repositories, MSR 2021*. IEEE, 2021, pp. 560–564.

[16] I. GitHub, "Saving repositories with stars," 2024, accessed: 2024-06-07. [Online]. Available: https://docs.github.com/en/get-started/exploring-projects-on-github/saving-repositories-with-stars

[17] GitHub, "The state of open source and ai," https://github.blog/2023-11-08-the-state-of-open-source-and-ai/, 2023, [Online; accessed 7-May-2024].

[18] N. R. Weeraddana, X. Xu, M. Alfadel, S. McIntosh, and M. Nagappan, "An empirical comparison of ethnic and gender diversity of DevOps and non-DevOps contributions to open-source projects," *Empirical Software Engineering*, vol. 28, no. 6, p. 150, 2023.

[19] B. Vasilescu, A. Serebrenik, and V. Filkov, "A data set for social diversity studies of GitHub teams," in *2015 IEEE/ACM 12th working conference on mining software repositories*. IEEE, 2015, pp. 514–517.

[20] M. Nejati, M. Alfadel, and S. McIntosh, "Code review of build system specifications: Prevalence, purposes, patterns, and perceptions," in *Proc. of the International Conference on Software Engineering (ICSE)*, 2023, p. 1213–1224.

[21] R. Agrawal, T. Imieliński, and A. Swami, "Mining association rules between sets of items in large databases," in *Proceedings of the 1993 ACM SIGMOD international conference on Management of data*, 1993, pp. 207–216.

[22] S. Developers, "scipy.stats.kruskal," https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.kruskal.html, 2024, accessed: 2024-06-08.

[23] K. Charmaz, "Grounded theory in global perspective: Reviews by international researchers," *Qualitative inquiry*, vol. 20, no. 9, pp. 1074–1084, 2014.

[24] T. Zimmermann, "Card-sorting: From text to themes," in *Perspectives on data science for software engineering*. Elsevier, 2016, pp. 137–141.

[25] S. B. Lo, A. D. Svensson, C. J. Presley, and B. L. Andersen, "A cognitive–behavioral model of dyspnea: qualitative interviews with individuals with advanced lung cancer," *Palliative & Supportive Care*, vol. 21, no. 6, pp. 1070–1077, 2023.

[26] K. Gallaba, M. Lamothe, and S. McIntosh, "Lessons from eight years of operational data from a continuous integration service: An exploratory case study of circleci," in *Proceedings of the 44th international conference on software engineering*, 2022, pp. 1330–1342.

[27] Y. Zhao, A. Serebrenik, Y. Zhou, V. Filkov, and B. Vasilescu, "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)*. IEEE, 2017, pp. 60–71.

[28] N. Cassee, B. Vasilescu, and A. Serebrenik, "The silent helper: the impact of continuous integration on code reviews," in *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2020, pp. 423–434.

[29] P. Rostami Mazrae, T. Mens, M. Golzadeh, and A. Decan, "On the usage, co-usage and migration of CI/CD tools: A qualitative analysis," *Empirical Software Engineering*, vol. 28, no. 2, p. 52, 2023.

[30] P. R. Mazrae, A. Decan, T. Mens, and M. Wessel, "A preliminary study of github actions workflow changes," 2023.

[31] H. O. Delicheh, A. Decan, and T. Mens, "A preliminary study of GitHub actions dependencies," in *CEUR Workshop Proceedings*, vol. 3483, 2023, pp. 66–77.

[32] A. Decan, T. Mens, and H. O. Delicheh, "On the outdatedness of workflows in the GitHub actions ecosystem," *Journal of Systems and Software*, vol. 206, p. 111827, 2023.

[33] P. Valenzuela-Toledo and A. Bergel, "Evolution of GitHub action workflows," in *2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2022, pp. 123–127.

[34] T. Kinsman, M. Wessel, M. A. Gerosa, and C. Treude, "How do software developers use GitHub actions to automate their workflows?" in *2021 IEEE/ACM 18th International Conference on Mining Software Repositories (MSR)*. IEEE, 2021, pp. 420–431.

[35] S. G. Saroar and M. Nayebi, "Developers' perception of github actions: A survey analysis," in *Proceedings of the 27th International Conference on Evaluation and Assessment in Software Engineering*, 2023, pp. 121–130.

[36] I. Bouzenia and M. Pradel, "Resource usage and optimization opportunities in workflows of GitHub actions," in *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, 2024, pp. 1–12.

[37] I. Koishybayev, A. Nahapetyan, R. Zachariah, S. Muralee, B. Reaves, A. Kapravelos, and A. Machiry, "Characterizing the security of GitHub CI workflows," in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 2747–2763.

[38] P. Valenzuela-Toledo, A. Bergel, T. Kehrer, and O. Nierstrasz, "EGAD: A moldable tool for GitHub action analysis," in *2023 IEEE/ACM 20th International Conference on Mining Software Repositories (MSR)*. IEEE, 2023, pp. 260–264.

[39] P. Valenzuela-Toledo, A. Bergel, T. Kehrer, and O. Nierstrasz, "Exploring GitHub actions through EGAD: an experience report," in *Proceedings of the International Workshop on Smalltalk Technologies, Lyon, France; August 29th-31st, 2023*, ser. CEUR Workshop Proceedings, vol. 3627. CEUR-WS.org, 2023.

[40] G. Cardoen, T. Mens, and A. Decan, "A dataset of GitHub actions workflow histories," in *2024 IEEE/ACM 21st International Conference on Mining Software Repositories (MSR)*. IEEE, 2024, pp. 677–681.