

tool-recommender-bot

Anonymous Author(s)

ABSTRACT

To increase software engineer productivity, toolsmiths create tools and features to automatically complete software development tasks. However, these useful tools are often undiscovered or ignored by developers, which is problematic for software applications that rely on programmer efficiency and correctness. This paper introduces a new approach to making software engineering tool recommendations that integrates characteristics from user-to-user suggestions and industry practices for researchers to increase awareness of their products among developers. To help improve tool adoption among software engineers, we implemented this approach in *tool-recommender-bot*, an automated recommendation system, and found our design is more effective for increasing tool discovery compared to other styles of tool recommendations.

KEYWORDS

Tool Recommendations; Tool Discovery;

ACM Reference format:

Anonymous Author(s). 2018. *tool-recommender-bot*. In *Proceedings of The 26th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Lake Buena Vista, Florida, United States, 4âc9 November, 2018 (ESEC/FSE 2018)*, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

To maintain and meet rising demands for technology, software engineers emphasize *software quality* throughout the development process, monitoring metrics that impact software producers and consumers [19]. However, despite increased attention to quality, buggy code remains a problem as the number of software errors increases [23]. The Software Fail Watch by Tricentis suggests software failures impacted 3.7 billion users and caused \$1.7 trillion in financial losses in 2017.¹ Additionally, the process of finding and fixing bugs, or debugging, is a time-consuming and costly activity. The National Institute of Standards and Technology reported software engineers spend 70-80% of work time debugging and on average one error takes 17.4 hours to debug [37]. Finding and fixing defects early during development is also important since studies show debugging costs increase the longer a bug remains in code [5, 9].

To improve code quality, researchers and toolsmiths have created software engineering tools to aid developers in their work. Research

shows that tools for static analysis [2], refactoring [30], security, and more are beneficial for improving code and preventing bugs. These tools can automatically perform a wide variety of software development tasks to save time and effort for developers. Additionally, the Software Engineering Body of Knowledge recommends using development tools because they can be used to achieve “desirable characteristics of software products” [43]. Software engineering tools are useful for reducing software errors and debugging costs while increasing developer productivity.

Although software quality is important, software engineering tools created to improve code quality are often ignored by developers [16]. Tilley and colleagues suggest tool adoption should be one of the most important goals of software engineering research, and present challenges preventing professional developers from adopting programming tools created by researchers, or research-off-the-shelf (ROTS) software [44]. Furthermore, researchers have explored barriers preventing software engineers from utilizing useful tools for debugging [8], static analysis [17], security [47], refactoring [29], documentation [14], and more. One barrier to tool adoption is the discoverability barrier, where users are unaware of a tool’s existence [28]. This lack of tool awareness in the software engineering industry can lead to poor-quality applications, inconvenienced users, and wasted time and money fixing errors. Automated approaches have been developed to increase awareness of software tools and features, however Murphy-Hill and colleagues found that developers prefer learning about tools from colleagues during normal work activities, or *peer interactions* [33].

To help solve the tool discovery problem among software engineers, we developed an approach for making automated tool recommendations. Our approach integrates characteristics of peer interactions and software engineering practices to suggest useful tools to developers. We base our technique on three design pillars: *commend*, *apprehend*, and *recommend*. These principles allow our approach to be applied to a variety software engineering tools and to make customized recommendations to developers. Our goal is for researchers and toolsmiths to implement this approach to increase awareness and adoption of their tools.

To evaluate our approach, we implemented *tool-recommender-bot*. The initial implementation of *tool-recommender-bot* for this study recommends ERROR PRONE², an open source static analysis tool for Java code, to developers on GitHub³, a popular code hosting and collaboration website. We measure the effectiveness of our approach by observing the frequency of recommendations and how developers react to receiving suggestions for tools from our system compared to other recommendation methods. This research makes the following contributions:

- a novel approach for recommender systems to increase awareness of software engineering tools, and
- an implementation and evaluation of *tool-recommender-bot* recommending a static analysis tool to developers.

²<http://errorprone.info>

³<https://github.com>

¹<https://www.tricentis.com/software-fail-watch/>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ESEC/FSE 2018, 4âc9 November, 2018, Lake Buena Vista, Florida, United States

© 2018 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

2 RELATED WORK

This project builds on previous work by researchers studying recommendations in software engineering.

Many researchers have developed approaches to try to solve the tool adoption problem for software users. Fischer and colleagues suggest passive help systems, which require software users to explicitly seek information, are ineffective and inconvenient, and *active help systems* designed to automatically guide and support users in their work are more effective [12]. Passive recommendation approaches such as *organizational memory* implemented in ANSWER GARDEN [1], *group memory* in TeamInfo [4], and *idea gardening* [7] by Cao and colleagues use social and collaborative means to improve tool adoption and help software users overcome barriers. Other active solutions like Maltzahn’s prototype recommender system TOOLBOX [26], CommunityCommands by Matejka and colleagues [27], and Linton and colleagues’ *organization-wide learning* (OWL) technique [24] monitor software actions by users by tracking logs on networked computers or stored in a database to suggest helpful commands and strategies to other users performing similar tasks. This research also aims to improve tool adoption among software users. To do improve the effectiveness of software recommendations, we introduce a new approach with three unique design principles and implement this approach in a novel active help system, *tool-recommender-bot*, that focuses on suggesting software engineering tools to developers.

More specifically, researchers have created techniques to improve the tool adoption problem among software engineers. Dubois and Tamburrelli argue that *gamification* can motivate software engineers to adopt new tools and better software development practices [10], and this has been implemented in systems such as Blaze [42], HALO [40], and Teamfeed [41]. Latoza and colleagues suggest *crowdsourcing* is also beneficial for software engineering [22]. This has been used by researchers to improve developer learning and feature adoption in systems like Codepourri for Python tutorials [15], CrowdStudy for web applications [34], and platforms such as TopCoder [21]. Our work builds on research in this area by introducing a distinctive recommendation approach with three original design principles to make tool suggestions to developers more effective.

Furthermore, researchers and toolsmiths have developed automated recommender systems to increase awareness of software engineering tools. SpyGlass automatically monitors developer actions, infers navigation tasks, and recommends tools to help users navigate code more efficiently [45]. Systems such as Review Bot by Balachandran [3] and Tricorder by Sadowski and colleagues [38] seek to improve adoption of program analysis in industry by automatically running static analysis tools and presenting the output to software engineers during code reviews and project builds. To improve the effectiveness of automated system-to-user recommendations, we implemented *tool-recommender-bot* to evaluate our recommendation approach based on design principles derived from user-to-user recommendations. Additionally, this is the first approach to our knowledge that can be extended to recommend different types of software engineering tools to developers and that can be utilized by other researchers to increase awareness of their work.

Finally, existing research has also analyzed different types and styles of software engineering tool recommendations. Murphy-Hill and colleagues outlined existing recommendation algorithms to propose new algorithms based on the previous methods [31]. Moreover, Murphy-Hill explored seven different methods of tool recommendations: peer observations, peer recommendations, tool encounters, tutorials, written descriptions, Twitter or RSS Feeds, and discussion threads- and found that peer interactions (peer observation and recommendation) were the most effective mode of tool discovery among software engineers [32, 33]. Piorkowski and colleagues evaluated different recommendation algorithms to determine if *information foraging theory* [36] can be applied to improve recommendations using structure and wording based on how software engineers seek information [35]. Additionally, Winters and colleagues analyzed different styles of tool recommendations and found that suggestion style impacts tool adoption for software developers [Mike]. This study compares the effectiveness of our approach to other recommender styles by modifying our design principles to investigate different wording and methods of communication in tool recommendations to developers.

3 APPROACH

Our approach uses three design pillars to increase awareness of software engineering tools: *commend* developers on their work, *apprehend* additional opportunities for improvement, and *recommend* useful development tools to enhance the quality of the code. Figure 1 provides a general overview of how we implement these design pillars as an active help system.

3.1 Commend

The first design principle of our approach finds places where a tool would be useful. To do this, our technique commends developers for their updates making improvements to code. We analyze code changes made by developers and then praise them for their work when they improve the quality of code. Additionally, psychology research also suggests people respond better to praise than blame [18]. Our approach applauds developers when they produce good work rather than rebuke them for making mistakes when they introduce new issues.

To determine improvements in code quality, our approach monitors projects for code changes made by software engineers (Figure 1A, 1.B). Software engineering tools are used to determine if programming tasks are performed by developers when they make code changes by automatically runs tools on the original version of the code and the modified version (Figure 1.C). Then, we compare the tool output from on both versions of code to determine if any development tasks were performed by the developer (Figure 1.D). For example, if we want to recommend a refactoring tool to a developer in an Integrated Development Environment (IDE), then we would wait for changes to the code base and use a tool to detect if a refactoring activity occurred in the modifications. If a refactoring activity is detected, our approach commends the programmer for their code change to improve the software quality when making a recommendation for the refactoring tool.

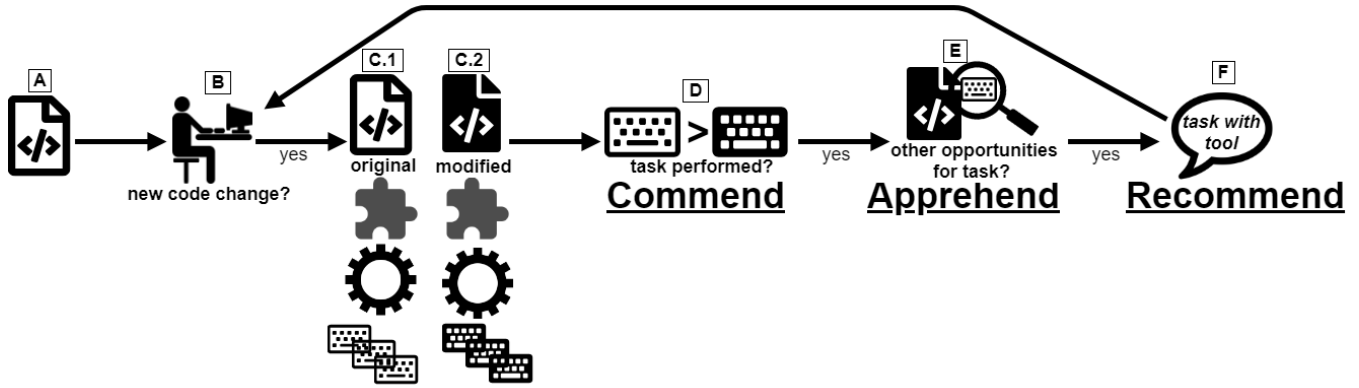


Figure 1: Flowchart presenting a general implementation of our approach: A) Store project source code; B) Check for code changes; C.1) Add tool plugin, run tool on original code, and collect output; C.2) Add tool plugin, run tool on modified code, and collect output; D) Check for development tasks performed between versions of code; E) Search for other places where task can be applied; F) Recommend tool to help perform task.

3.2 Apprehend

Next, our approach finds additional places to recommend useful software engineering tools. To find other opportunities for recommendations, we apprehend instances where developers can make source code improvements similar to their previous code changes and based on development tool analysis from the commend phase. This design principle helps reduce effort and provides information-related needs to software engineers, which Fischer suggests is a goal for help systems [12].

After determining whether a development activity occurred, our approach searches for similar instances in the code where the same activity can be applied. To apprehend opportunities for similar programming tasks, our approach involves automatically analyzing the tool output from the updated code to automatically find new opportunities where developers can perform the same task (Figure 1.E). For example, after the programmer completes the refactoring task in the IDE our approach analyzes the entire updated version of the code base to find other places where the same refactoring activity is applicable.

3.3 Recommend

The final step of our approach is to recommend software engineering tools to developers (Figure 1.F). The recommendation to developers includes commending them for improving code quality based on their changes and apprehending new opportunities to make the same improvement. In the continuing example, after determining refactoring activity was performed and finding other locations in the code where the same task can be performed, then our approach will make a recommendation to the developer in the IDE. The recommendation will praise the developer for their refactoring change, provide additional places in the code where they can make the same refactoring change, and suggest using a refactoring tool for future refactoring activities. Our goal is that through commending and apprehending, we can improve the effectiveness of recommendations for software engineering tools and increase adoption among developers.

4 MOTIVATION

Our tool recommendation approach is motivated by previous research in tool discovery and peer interactions.

4.1 Peer Interactions

Research suggests recommendations between colleagues is the most effective mode of tool discovery. Murphy-Hill and colleagues interviewed software engineers and found that, even though they appear less frequently in the workplace, developers prefer peer interactions over system-to-user recommendations. [33]. Similarly, Welty discovered that software users sought help from colleagues more often than search engines and help menus [46].

To better understand what makes peer interactions effective, Brown and colleagues observed how software users recommend tools to each other while completing tasks. Their results suggest *receptiveness* is a significant factor in determining the outcome of tool recommendations, while other characteristics such as politeness and persuasiveness are not as important [6]. Fogg argues a receptive audience is vital for designing persuasive technology, and describes receptiveness as: 1) demonstrating a desire and 2) familiarity with acquiring a target behavior [13]. In this research, our target behavior is the adoption of useful development tools and these two criteria for receptiveness influenced the design principles for our approach for recommending software engineering tools.

4.1.1 Demonstrate Desire. Brown and colleagues defined demonstrating a desire to adopt a particular behavior as users expressing interest in “discovering, using, or learning more information” about a new practice. In software engineering, developers desire to write high-quality programs without errors. This desire is demonstrated by the fact software engineers spend the majority their time testing and debugging programs [37]. Our approach capitalizes on developers’ desires to produce mistake-free code by commending developers for their work. To help improve tool adoption among software engineers, our approach analyzes code changes made by

developers to determine their intentions, then compliments the developer on their attempt to improve the code quality.

4.1.2 Familiarity. Familiarity is also a key part of receptiveness and important in increasing adoption of a target behavior. Fogg suggests users are more likely to adopt a target behavior if they are familiar with it [13]. To improve usage of software engineering tools, our approach integrates familiarity. First, we make recommendations based on developers' contributions to recent projects. Scalabrino and colleagues claim code understandability is one of the most important factors for software development, maintenance, debugging, and testing [39]. Developers should be familiar with the project before making changes to the code base. More specifically, our approach apprehends code the developer may also want to change based on their previous update.

5 IMPLEMENTATION WALKTHROUGH

To implement our approach, we developed *tool-recommender-bot* to recommend the ERROR PRONE static analysis tool to developers on GitHub. This section provides a general overview of how our automated recommendation system works.

5.1 Commend

First, we clone repositories locally to detect when GitHub users make changes to a project as a commit or pull request (Figure 3.A.3, 3.B). *tool-recommender-bot* automatically analyzes code changes to find opportunities to commend developers when bugs reported by ERROR PRONE are removed. ERROR PRONE is a static analysis tool that uses a suite of defined bug patterns to detect errors in Java code. Static analysis tools like ERROR PRONE are useful for debugging and preventing errors in source code for applications, however they are often underutilized by software engineers [17]. More details on how our implementation determines a bug fix can be found in Section 6.2.

5.2 Apprehend

Next, *tool-recommender-bot* apprehends places where similar opportunities to the code can be made (Figure 1.E). GitHub developers should be knowledgeable and familiar with the changes they propose, as well as the code base to which they are contributing. *tool-recommender-bot* suggests ERROR PRONE when a reported bug is fixed by a developer in a commit or pull request, but the same error still exists elsewhere in the code.

5.3 Recommend

Finally, *tool-recommender-bot* recommends ERROR PRONE to the GitHub user at the location of their fix in the code. We commend developers by telling them "Good job!" and presenting the error they fixed. *tool-recommender-bot* provides the apprehended opportunities for similar code changes as links to buggy lines with the same error in other parts of the code. Finally, we provide a link to the ERROR PRONE website for users to learn more about the tool if they are interested. Figure 2 provides a sample recommendation from *tool-recommender-bot* to a developer for ERROR PRONE on a GitHub pull request.

6 IMPLEMENTATION DESCRIPTION

This sections provides technical details on the implementation of our approach in *tool-recommender-bot*. Figure 3 shows a specific flowchart displaying the implementation of our design principles for the evaluation. The design of our system is motivated by current practices in the software engineering industry to make *tool-recommender-bot* more similar to a peer for developers.

6.1 Continuous Integration

To analyze developer changes, our system utilizes continuous integration concepts and tools to observe code modifications to GitHub repositories. *tool-recommender-bot* is implemented as a plugin for Jenkins, "the leading open source automation server" for source code deployment and delivery.⁴ We use Jenkins to periodically check for new modifications, commits and pull requests, to GitHub projects every 15 minutes. Our system ignores any commits or pull requests from developers that do not modify a Java file. When a new code change is found, Jenkins to automatically analyze the patch and run our approach.

To analyze the source code, we target projects that use the Maven⁵ build automation and software dependency management tool for Java applications. We clone the GitHub repository, automatically inject ERROR PRONE as a Maven plugin to a repository's project object model file (*pom.xml*), and then run the build process with the static analysis tool. *tool-recommender-bot* builds both the original version of the code before the proposed changes (base) and the updated version of the repository with the modifications from the developer (head) to inspect differences in the ERROR PRONE output. The JGit Java API⁶ is used to clone projects and checkout the base and head versions of the code for our analysis.

6.2 Debugging

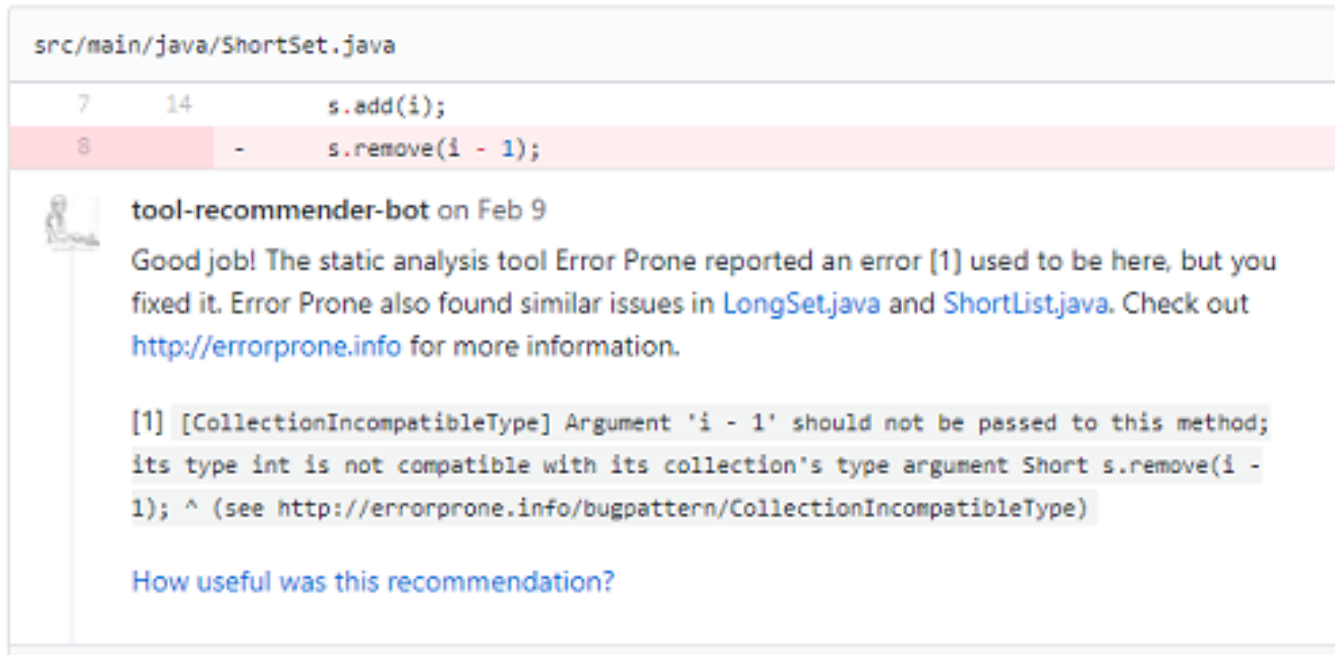
Before commending developers on their work, *tool-recommender-bot* must debug to find errors and determine if proposed code changes are a fix. After building the base and head versions with ERROR PRONE, our system parses the output of each build to determine if any faults reported were removed between versions of code. To determine if a change fixes a defect, we developed an algorithm using the code differencing tool GumTree [11]. GumTree allows us to identify actions (addition, delete, insert, move, and update) performed between the altered versions of the project.

To determine if an error was fixed, we take several things into consideration: First, we ignore errors that are removed but were not located in a file modified by the developer. This ensures that the GitHub user will be familiar with the code changed and potential error fixed. Second, we ignore changes where only delete actions were detected between the base and head versions of a file. This avoids making recommendations in situations where defects were only removed by developers. For example, deleting a class will remove errors reported by ERROR PRONE in the source code, however the intention was not to fix the bugs. Thus, for a change to be considered a fix there must be new code added or existing code modified by the developer. Similarly, we also ignore classes that

⁴<https://jenkins.io/>

⁵<https://maven.apache.org/>

⁶<https://eclipse.org/jgit/>

Figure 2: Screenshot of recommendation from *tool-recommender-bot*

are deprecated by developers. These conditions were put in place minimize false positives and prevent errant recommendations to software engineers in our approach.

When a fix is identified in the changed version of code, *tool-recommender-bot* finds the location of the fix in the head version. To find the modified line of code that fixed a bug, we use GumTree to parse the source code and convert it to abstract syntax trees. We look for the action closest to the offset of the error node determined from the line number reported by ERROR PRONE. If the closest action is not a delete, then our approach uses the location of that action. Otherwise, our algorithm iteratively searches for the closest sibling node or parent nodes that is not a delete action. To apprehend different opportunities for similar changes, we iterate through the list of errors reported by ERROR PRONE in the head version of code and look for instances of the same error that was fixed by the developer.

6.3 Code Review

Code reviews between developers are a standard procedure of software engineering teams to improve code quality [25]. This practice also applies to GitHub projects, with many repositories requiring approval from another developer before changes can be merged into the main code base [48]. Our approach simulates peer code reviews by making recommendations for static analysis tools as a comment on pull requests and commits. Github allows users to make comments on specific lines of code in situ with the code changes. *tool-recommender-bot* determines where to automatically make recommendation comments by converting the line number of the identified fix to the equivalent position in the diff file, or textual representation of code changes made in a commit or pull request,

represented by the number of lines below the “@@” symbol in the header⁷. Our system uses the jcabi Object Oriented GitHub API⁸ to automatically comment on GitHub commits and pull requests. An example of a comment from *tool-recommender-bot* is displayed in Figure 2.

To increase the likelihood of tool adoption, *tool-recommender-bot* implements our approach by commending developers on their changes, apprehending chances for similar modifications, and recommending software engineering tools to find even more related errors. In the comment, *tool-recommender-bot* uses language similar to a peer to compliment the author’s code contribution. For instance, our system uses “Good job!” to commend developers for fixing an error. Additionally, our system presents similar instances of the fixed error found elsewhere in the code. *tool-recommender-bot* automatically adds direct links to at most two different locations of the same defect where a similar fix can be applied. Finally, we recommend ERROR PRONE and provide information about the tool to encourage developers to use software engineering tools in their future work to fix more related errors and more.

7 METHODOLOGY

Our study evaluates the effectiveness of *tool-recommender-bot* by analyzing how often our system recommends software engineering tools and how developers respond to recommendations from our system.

⁷<https://developer.github.com/v3/pulls/comments/#create-a-comment>

⁸<http://github.jcabi.com/>

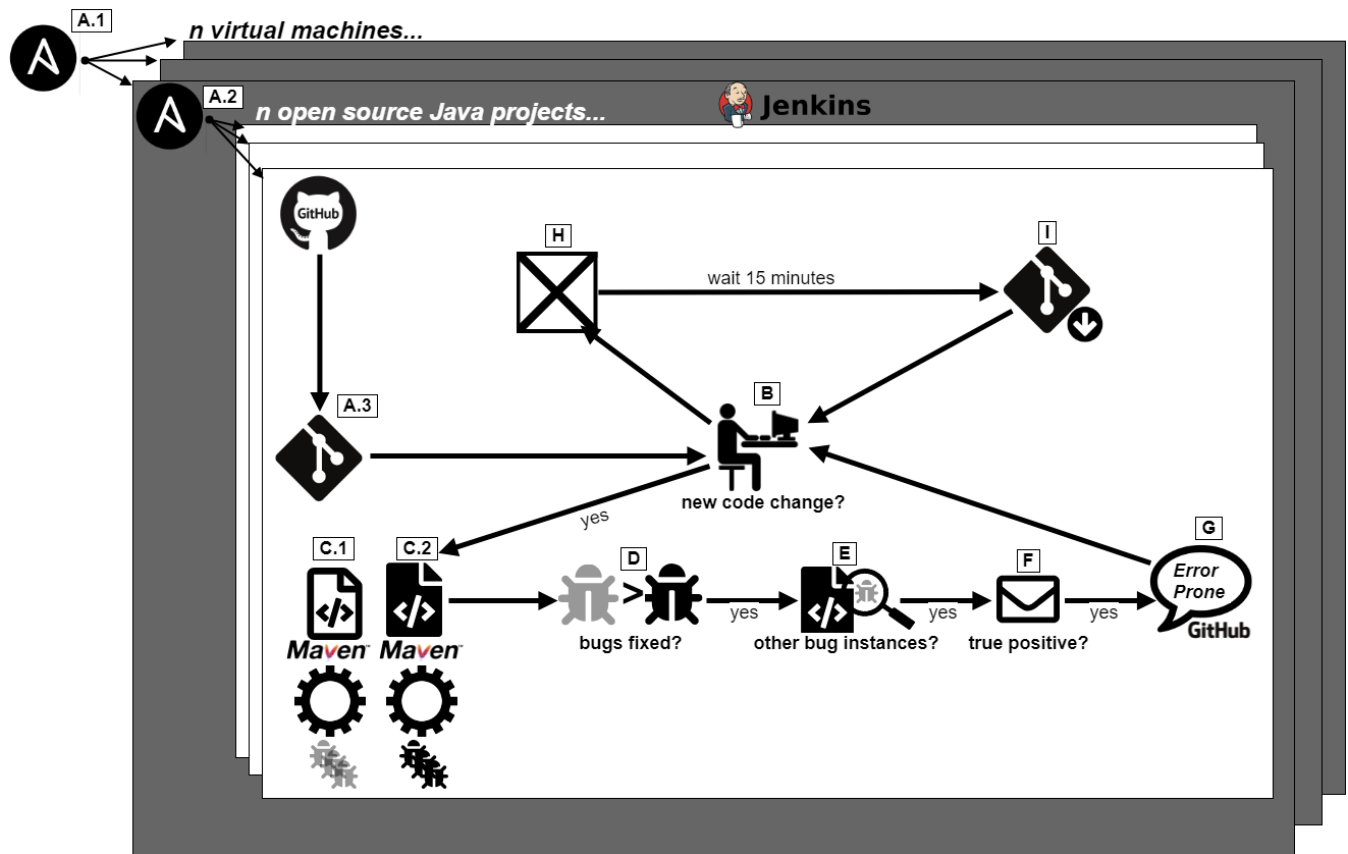


Figure 3: Flowchart presenting our implementation of tool-recommender-bot for this evaluation: A) Store project source code; B) Check for code changes; C.1) Add tool plugin, run tool on original code, and collect output; C.2) Add tool plugin, run tool on modified code, and collect output; D) Check for development tasks performed between versions of code; E) Search for other places where task can be applied; F) Recommend tool to help perform task.

7.1 Projects

We used real-world open source software applications to evaluate *tool-recommender-bot*. To choose projects for this study from the millions of GitHub repositories online, we used the following criteria:

- primarily written in Java,
- build using Maven,
- do not already use ERROR PRONE,
- ranked among the most popular and most recently updated repositories

Our evaluation was limited to Java projects since ERROR PRONE can only analyze Java source code. To determine if a repository used Maven as a build system, we automatically checked if a Project Object Model (*pom.xml*) configuration file was located in the home directory. We also checked to make sure that the *pom.xml* did not already contain the ERROR PRONE plugin to avoid projects that already use ERROR PRONE. We selected projects that don't use ERROR PRONE to increase awareness of the tool in recommendations. Developers are less likely to know about the tool if the projects they contribute to do not implement it in their build.

To get the most popular repositories, we filtered GitHub projects by the amount of stars based on numbers in the Fibonacci sequence. We chose the most starred repositories to study the most popular projects on GitHub. Stars are a social aspect of GitHub where users can indicate their favorite projects and repositories of interest⁹. Using Fibonacci numbers allowed us to get a higher concentration of projects with a lower amount of stars, while fewer projects will have a very large number of stars. To filter of repositories, we grouped projects with 1 or 2 stars, 2 or 3 stars, between 3 and 5 stars, between 5 and 8 stars, etc., and sorted the top 100 projects in each group by when they were most recently updated. After using a GitHub search API to find projects that met the above criteria, we compiled a list of 789 code repositories. Out of those projects, one repository failed due to broken Unicode text. The projects used for our evaluation include a wide range of software applications providing a variety of services from large software companies such as Google and Apache to individual developers. A list of projects used for this study is publicly available online.¹⁰

⁹<https://help.github.com/articles/about-stars/>

¹⁰<https://gist.github.com/tool-recommender-bot/1769ccd148508beabcd273a731723860>

7.2 Study Design

To evaluate the effectiveness of our approach, we compared *tool-recommender-bot* to different styles and mediums of tool recommendations.

7.2.1 Setup. To analyze 700+ GitHub repositories simultaneously, we set up *tool-recommender-bot* on virtual machines. We used Ansible¹¹ to automatically setup and install dependencies on the all the virtual machines, and then on each VM we use it to setup the Jenkins server and generate Jenkins jobs for each project. Each Jenkins job clones the GitHub repository, periodically checks for new commits and pull requests every 15 minutes, and runs *tool-recommender-bot* to analyze project code changes. Figure 3.A.1-A.3 displays the setup process for our evaluation of *tool-recommender-bot*. Additionally, the scripts and instructions to setup the study environment are publicly available online.¹²

7.2.2 Recommendation Styles. To test the effectiveness of our approach, we used four different recommendation styles based on our design principles of commending developers for their work and apprehending opportunities to make similar changes. The four recommendation styles include commending and apprehending, commending only, apprehending only, and neither commending nor apprehending. In each case, a recommendation was always made. We adapted *tool-recommender-bot* to make recommendations without commending or without apprehending. Without apprehending, our system simply commends the developer for their fix and does not link to similar errors found in other parts of the code. For the final recommendation without commending or apprehending, we decided to send recommendations to developers by email without commending them for their patch and not providing other instances to make a similar fix. An overview of each style is presented in Table 1.

7.2.3 Recommendation Review. For our evaluation, we manually reviewed each instance where *tool-recommender-bot* reported a recommendation to a developer. This step can be seen in Figure 3.F. We do this for several reasons: First, we wanted to ensure that only one ERROR PRONE recommendation was made to each project to avoid biasing our results by minimizing the chance other GitHub users working on the same project discover the tool from a comment made to that repository. Additionally, in some cases one commit or pull request provided multiple opportunities for a recommendation. In this case, the authors examined each instance and selected one of errors reported to make a recommendation on.

Next, we inspected cases to avoid making unnecessary recommendations to developers. Johnson and colleagues discovered one of the primary barriers of static analysis tool usage for software engineers is false positive output [17]. To prevent multiple comments on one project and false positives from our system, we modified *tool-recommender-bot* to send an email to the researchers with the proposed recommendation. After checking the code changes, we proceeded to use *tool-recommender-bot* to post the comment recommending ERROR PRONE to the developer on their pull request or commit if the fix was deemed to be a true positive and the repository had not previously received a recommendation. Otherwise,

we noted the instance as a false positive in our implementation and did not make a recommendation.

7.2.4 Debriefing. To gather data on the usefulness of our system, we sent a follow-up survey to developers. Survey participants were users who received a recommendation from *tool-recommender-bot* on their pull request or commit. We asked developers about their awareness of ERROR PRONE and how likely they are to use the tool in the future. The survey also included a free-response section to provide an opportunity for participants to add comments on the usefulness of the recommendation.

Developers voluntarily consented to complete the survey and provide feedback on our system. To ensure developers answered honestly, we notified respondents that their answers will be used for research purposes. Previous research has shown that survey participants are more motivated to answer truthfully if they know they are contributing to research [20].

To further examine the effectiveness of *tool-recommender-bot*, we compared our approach to sending email recommendations to developers. To study this, we found similar instances of code fixes by GitHub users where our system would recommend a tool and, instead of making the recommendation with *tool-recommender-bot* on the pull request or commit, send an email suggesting ERROR PRONE to the developers. The emails also contained the recommendation feedback survey for developers, and we compared results to see how software engineers responded to receiving a recommendation by email vs. on GitHub from *tool-recommender-bot*. To send a recommendation via email, the developer must have an email address publicly available on the GitHub user profile.

7.3 Data Analysis

We analyzed the data collected in our study to determine effectiveness of our automated tool recommendation system.

7.3.1 Quantitative. To determine the effectiveness of our approach, we observed how often *tool-recommender-bot* makes recommendations on commits and pull requests. In addition to the frequency of recommendations, we also tracked instances where ERROR PRONE defects were removed but not reported as fixed according to our fix identification algorithm in Section III.B.2, the number of occurrences where errors were fixed but no other instances of that bug were found in the code, and the false positive rate.

Effective tool recommendation systems should have ample opportunities to regularly make recommendations to users. To determine how often *tool-recommender-bot* automatically recommends ERROR PRONE, we observed the total number of new pull requests and commits, and compared it to the amount recommendations made by our tool. We calculated the rate of true positive recommendations during the span of our study to measure the recommendation rate for each GitHub repository used in our evaluation. To calculate the false positive rate, we compared the number of unnecessary instances where our tool proposed a recommendation found by the authors to the total number of instances where a recommendation reported by *tool-recommender-bot*.

7.3.2 Qualitative. For our second research question, we accumulated responses from developers in our follow-up survey presented

¹¹<https://www.ansible.com/>

¹²https://github.com/chbrown13/tool-recommender-bot/tree/master/eval_scripts/

Table 1: Study Recommendation Styles

Apprehend	Commend	
	Yes	No
	No	Yes
Yes	<i>tool-recommender-bot</i>	Modified <i>tool-recommender-bot</i> (without commending)
No	Modified <i>tool-recommender-bot</i> (without apprehending)	Email

in recommendations from *tool-recommender-bot* and by email to determine the usefulness of our system. We utilized a five-point Likert scale for participants to rank how knowledgeable they were about the existence of ERROR PRONE before the recommendation and how likely they are to use ERROR PRONE for future development tasks. An optional free response section was provided at the end for respondents to describe explain why or why not they found the recommendation useful. These responses were used to analyze developers' reactions to our automated recommendation. Finally, researchers analyzed and independently coded open-ended responses from participants to further analyze the effectiveness of our approach based on feedback from software developers. We measured the response rate by observing the total number of recommendations made, the total number of survey responses, and the positive and negative responses from developers who received a recommendation via *tool-recommender-bot* and via email.

8 RESULTS

8.1 Quantitative

Tons of recommendations...

No false positives...

8.2 Qualitative

8.2.1 *Recommend.*

8.2.2 *Commend, Recommend.*

8.2.3 *Apprehend, Recommend.*

8.2.4 *Commend, Apprehend, Recommend.* Excellent responses from recommendees...

Something statistically significant...

9 DISCUSSION

9.1 Observations

9.1.1 *Why Were There So Few Recommendations?* Non-java changes, number of errors removed but not fixed, number of errors fixed without another instance in the code, manual inspection of pull requests and commits...

9.2 Implications

Here's what our results say about improving tool recommendation systems...

10 LIMITATIONS

An internal threat to the validity of this work is our use of code differencing to determine if developers intended to fix a bug in a commit or pull request. We cannot definitively determine the intentions of GitHub developers making changes to a repository, however two authors analyzed the code changes and came to an agreement on if the modified patch was a fix before making a recommendation to the developer. Additionally, although we used a Likert scale to measure if GitHub users who received a recommendation were likely to use ERROR PRONE in the future, we did not measure if the tool was actually adopted by the developers for future tasks.

Our evaluation has limited generalizability due to the fact we only assessed recommendations for the ERROR PRONE static analysis tool. This restricted our study to evaluate Java projects and a specific set of errors that can be reported by the tool. We selected ERROR PRONE because it is able to report a wide variety of errors based on bug patterns for Java code. Another external threat to the validity of this study is the projects selected for our evaluation. We only examined open source repositories on GitHub, and these results may not generalize to developers of closed source software or projects on other code hosting sites, such as BitBucket.¹³ To minimize this threat, we evaluated *tool-recommender-bot* on a large number of popular software applications on GitHub that provide many different services from a wide variety of software companies and developers.

11 FUTURE WORK

The main goal of our future work is to increase the practicality of *tool-recommender-bot* for researchers and toolsmiths to implement our approach with their tools. A next step is to implement our approach with multiple static analysis tools, such as Checkstyle¹⁴ for Java. This will allow us to have more opportunities to make recommendations to developers based on the output from different tools. Additionally, we plan to update our system to work with software engineering tools that integrate as plugins for different build systems such as Gradle¹⁵ and Travis CI¹⁶, as opposed to just Maven projects. Future work will also extend *tool-recommender-bot* to work with different types of tools to increase adoption and usage, for example Find Security Bugs¹⁷ which is useful for finding security vulnerabilities in Java web applications. Finally, future work will consist of expanding *tool-recommender-bot* to work with software engineering tools for different programming languages,

¹³<http://bitbucket.org>

¹⁴<http://checkstyle.sourceforge.net/>

¹⁵<https://gradle.org/>

¹⁶<https://travis-ci.org/>

¹⁷<https://find-sec-bugs.github.io/>

such as the Pylint¹⁸ static analysis tool for Python, clang¹⁹ compiler and static analyzer for C and C++, and more.

12 CONCLUSION

tool-recommender-bot is awesome

REFERENCES

- [1] M. S. Ackerman and T. W. Malone. Answer garden: A tool for growing organizational memory. In *Proceedings of the ACM SIGOIS and IEEE CS TC-OA Conference on Office Information Systems*, COCS '90, pages 31–39, New York, NY, USA, 1990. ACM.
- [2] N. Ayewah, D. Hovemeyer, J. D. Morgenthaler, J. Penix, and W. Pugh. Using static analysis to find bugs. *IEEE software*, 25(5), 2008.
- [3] V. Balachandran. Reducing human effort and improving quality in peer code reviews using automatic static analysis and reviewer recommendation. In *Software Engineering (ICSE), 2013 35th International Conference on*, pages 931–940. IEEE, 2013.
- [4] L. M. Berlin, R. Jeffries, V. L. O'Day, A. Paepcke, and C. Wharton. Where did you put it? issues in the design and use of a group memory. In *Proceedings of the INTERACT'93 and CHI'93 conference on Human factors in computing systems*, pages 23–30. ACM, 1993.
- [5] B. W. Boehm. Software engineering economics. *IEEE transactions on Software Engineering*, (1):4–21, 1984.
- [6] C. Brown, J. Middleton, E. Sharma, and E. Murphy-Hill. How software users recommend tools to each other. In *Visual Languages and Human-Centric Computing*, 2017.
- [7] J. Cao, S. D. Fleming, and M. Burnett. An exploration of design opportunities for end-user programmers' ideas. In *Visual Languages and Human-Centric Computing (VL/HCC), 2011 IEEE Symposium on*, pages 35–42. IEEE, 2011.
- [8] J. Cao, K. Rector, T. H. Park, S. D. Fleming, M. Burnett, and S. Wiedenbeck. A debugging perspective on end-user mashup programming. In *Visual Languages and Human-Centric Computing (VL/HCC), 2010 IEEE Symposium on*, pages 149–156. IEEE, 2010.
- [9] M. Dawson, D. N. Burrell, E. Rahim, and S. Brewster. Integrating software assurance into the software development life cycle (sdlc). *Journal of Information Systems Technology & Planning*, 3(6):49–53, 2010.
- [10] D. J. Dubois and G. Tamburrelli. Understanding gamification mechanisms for software development. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, pages 659–662. ACM, 2013.
- [11] J. Falleri, F. Morandat, X. Blanc, M. Martinez, and M. Monperrus. Fine-grained and accurate source code differencing. In *ACM/IEEE International Conference on Automated Software Engineering, ASE '14, Vasteras, Sweden - September 15 - 19, 2014*, pages 313–324, 2014.
- [12] G. Fischer, A. Lemke, and T. Schwab. *Active help systems*, pages 115–131. Springer Berlin Heidelberg, Berlin, Heidelberg, 1984.
- [13] B. Fogg. Creating persuasive technologies: An eight-step design process. In *Proceedings of the 4th International Conference on Persuasive Technology*, Persuasive '09, pages 44:1–44:6, New York, NY, USA, 2009. ACM.
- [14] A. Forward and T. C. Lethbridge. The relevance of software documentation, tools and technologies: a survey. In *Proceedings of the 2002 ACM symposium on Document engineering*, pages 26–33. ACM, 2002.
- [15] M. Gordon and P. J. Guo. Codepourri: Creating visual coding tutorials using a volunteer crowd of learners. In *2015 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*, pages 13–21, Oct 2015.
- [16] V. Ivanov, A. Rogers, G. Succi, J. Yi, and V. Zorin. What do software engineers care about? gaps between research and practice. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*, pages 890–895. ACM, 2017.
- [17] B. Johnson, Y. Song, E. Murphy-Hill, and R. Bowdidge. Why Don't Software Developers Use Static Analysis Tools to Find Bugs? In *Proceedings of the 2013 International Conference on Software Engineering (ICSE)*, ICSE '13, pages 672–681, Piscataway, NJ, USA, 2013. IEEE Press.
- [18] W. A. Kennedy and H. C. Willcutt. Praise and blame as incentives. *Psychological Bulletin*, 62(5):323, 1964.
- [19] B. Kitchenham and S. L. Pfleeger. Software quality: the elusive target [special issues section]. *IEEE software*, 13(1):12–21, 1996.
- [20] J. A. Krosnick. Response strategies for coping with the cognitive demands of attitude measures in surveys. *Applied cognitive psychology*, 5(3):213–236, 1991.
- [21] K. Lakhani, D. Garvin, and E. Lonstein. Topcoder (a): Developing software through crowdsourcing. 2010.
- [22] T. D. LaToza and A. van der Hoek. Crowdsourcing in software engineering: Models, motivations, and challenges. *IEEE software*, 33(1):74–80, 2016.
- [23] Z. Li, L. Tan, X. Wang, S. Lu, Y. Zhou, and C. Zhai. Have things changed now?: An empirical study of bug characteristics in modern open source software. In *Proceedings of the 1st Workshop on Architectural and System Support for Improving Software Dependability*, ASID '06, pages 25–33, New York, NY, USA, 2006. ACM.
- [24] F. Linton, D. Joy, H. peter Schaefer, and A. Charron. Owl: A recommender system for organization-wide learning. 2000.
- [25] L. MacLeod, M. Greiler, M. A. Storey, C. Bird, and J. Czerwinka. Code reviewing in the trenches: Understanding challenges and best practices. *IEEE Software*, PP(99):1–1, 2017.
- [26] C. Maltzahn and D. Vollmar. Toolbox: a living directory for unix tools owned by the community. Technical report, Citeseer, 1994.
- [27] J. Matejka, W. Li, T. Grossman, and G. Fitzmaurice. Communitycommands: command recommendations for software applications. In *Proceedings of the 22nd annual ACM symposium on User interface software and technology*, pages 193–202. ACM, 2009.
- [28] E. Murphy-Hill. Continuous social screencasting to facilitate software tool discovery. In *Proceedings of the 34th International Conference on Software Engineering*, ICSE '12, pages 1317–1320, Piscataway, NJ, USA, 2012. IEEE Press.
- [29] E. Murphy-Hill and A. Black. Breaking the barriers to successful refactoring. In *2008 ACM/IEEE 30th International Conference on Software Engineering*, pages 421–430, May 2008.
- [30] E. Murphy-Hill and A. P. Black. Refactoring Tools: Fitness for Purpose. *IEEE Software*, 25(5):38–44, 2008.
- [31] E. Murphy-Hill, R. Jiresal, and G. C. Murphy. Improving software developers' fluency by recommending development environment commands. In *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*, FSE '12, pages 42:1–42:11, New York, NY, USA, 2012. ACM.
- [32] E. Murphy-Hill, D. Y. Lee, G. C. Murphy, and J. McGrenere. How do users discover new tools in software development and beyond? *Computer Supported Cooperative Work (CSCW)*, 24(5):389–422, 2015.
- [33] E. Murphy-Hill and G. C. Murphy. Peer interaction effectively, yet infrequently, enables programmers to discover new tools. In *Proceedings of the ACM 2011 Conference on Computer Supported Cooperative Work, CSCW '11*, pages 405–414, New York, NY, USA, 2011. ACM.
- [34] M. Nebeling, M. Speicher, and M. C. Norrie. Crowdstudy: General toolkit for crowdsourced evaluation of web interfaces. In *Proceedings of the 5th ACM SIGCHI symposium on Engineering interactive computing systems*, pages 255–264. ACM, 2013.
- [35] D. Piorkowski, S. Fleming, C. Scaffidi, C. Bogart, M. Burnett, B. John, R. Bellamy, and C. Swart. Reactive information foraging: An empirical investigation of theory-based recommender systems for programmers. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '12, pages 1471–1480, New York, NY, USA, 2012. ACM.
- [36] P. Pirolli and S. Card. Information foraging. *Psychological review*, 106(4):643, 1999.
- [37] S. Planning. The economic impacts of inadequate infrastructure for software testing. 2002.
- [38] C. Sadowski, J. Van Gogh, C. Jaspan, E. Söderberg, and C. Winter. Tricorder: Building a program analysis ecosystem. In *Proceedings of the 37th International Conference on Software Engineering-Volume 1*, pages 598–608. IEEE Press, 2015.
- [39] S. Scalabrino, G. Bavota, C. Vendome, M. Linares-Vásquez, D. Poshyvanyk, and R. Oliveto. Automatically assessing code understandability: how far are we? In *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering*, pages 417–427. IEEE Press, 2017.
- [40] S. Sheth, J. Bell, and G. Kaiser. Halo (highly addictive, socially optimized) software engineering. In *Proceedings of the 1st International Workshop on Games and Software Engineering*, pages 29–32. ACM, 2011.
- [41] L. Singer and K. Schneider. It was a bit of a race: Gamification of version control. In *Games and Software Engineering (GAS), 2012 2nd International Workshop on*, pages 5–8. IEEE, 2012.
- [42] W. Snipes, A. R. Nair, and E. Murphy-Hill. Experiences gamifying developer adoption of practices and tools. In *Companion Proceedings of the 36th International Conference on Software Engineering*, pages 105–114. ACM, 2014.
- [43] I. C. Society, P. Bourque, and R. E. Fairley. *Guide to the Software Engineering Body of Knowledge (SWEBOK(R)) - Version 3.0*. IEEE Computer Society Press, Los Alamitos, CA, USA, 3rd edition, 2014.
- [44] S. Tilley, S. Huang, and T. Payne. On the challenges of adopting rots software. In *Proceedings of the 3rd International Workshop on Adoption-Centric Software Engineering*, pages 3–6, 2003.
- [45] P. Viriyakattiyaporn and G. C. Murphy. Improving program navigation with an active help system. In *Proceedings of the 2010 Conference of the Center for Advanced Studies on Collaborative Research, CASCON '10*, pages 27–41, Riverton, NJ, USA, 2010. IBM Corp.
- [46] C. J. Welty. Usage of and satisfaction with online help vs. search engines for aid in software use. In *Proceedings of the 29th ACM International Conference on Design of Communication, SIGDOC '11*, pages 203–210, New York, NY, USA, 2011. ACM.

¹⁸<https://www.pylint.org/>

¹⁹<https://clang.llvm.org/>

- [47] S. Xiao, J. Witschey, and E. Murphy-Hill. Social influences on secure development tool adoption: Why security tools spread. In *Proceedings of the 17th ACM Conference on Computer Supported Cooperative Work & Social Computing, CSCW '14*, pages 1095–1106, New York, NY, USA, 2014. ACM.
- [48] Y. Yu, H. Wang, G. Yin, and C. X. Ling. Reviewer recommender of pull-requests in github. In *2014 IEEE International Conference on Software Maintenance and Evolution*, pages 609–612, Sept 2014.