

ACWRECOMMENDER: A Tool for Validating Actionable Warnings with Weak Supervision

Zhipeng Xue[†], Zhipeng Gao^{*†}, Xing Hu[†], Shanping Li[†]

[†]Zhejiang University, Hangzhou, China
{ zhipengxue, zhipeng.gao, xinghu, shan }@zju.edu.cn

Abstract—Static analysis tools have gained popularity among developers for finding potential bugs, but their widespread adoption is hindered by the accompanying high false alarm rates (up to 90%). To address this challenge, previous studies proposed the concept of actionable warnings, and apply machine-learning methods to distinguish actionable warnings from false alarms. Despite these efforts, our preliminary study suggests that the current methods used to collect actionable warnings are rather shaky and unreliable, resulting in a large proportion of invalid actionable warnings. In this work, we mined 68,274 reversions from Top-500 Github C repositories to create a substantial actionable warning dataset and assigned weak labels to each warning’s likelihood of being a real bug. To automatically identify actionable warnings and recommend those with a high probability of being real bugs (AWHB), we propose a two-stage framework called ACWRECOMMENDER. In the first stage, our tool use a pre-trained model, i.e., UniXcoder, to identify actionable warnings from a huge number of SA tool’s reported warnings. In the second stage, we rerank valid actionable warnings to the top by using weakly supervised learning. Experimental results showed that our tool outperformed several baselines for actionable warning detection (in terms of F1-score) and performed better for AWHB recommendation (in terms of nDCG and MRR). Additionally, we also performed an in-the-wild evaluation, we manually validated 24 warnings out of 2,197 reported warnings on 10 randomly selected projects, 22 of which were confirmed by developers as real bugs, demonstrating the practical usage of our tool.

Index Terms—Actionable warning recommendation, Static analysis, Weak supervision, Data mining

I. INTRODUCTION

The static analysis tool has been widely used by software developers and companies to detect potential bugs and report warnings in recent years [1]–[3]. For example, Facebook has developed *infer*, a static code analysis tool for checking generic bug patterns (e.g, null pointer exceptions, memory leaks, race conditions) in their Android and iOS apps (including the main Facebook, Whatsapp, Instagram app and many others). Due to the lightweight analysis and low computational cost, these static analysis tools have gained popularity among developers. However, static analysis tools face two major challenges: Firstly, they have a high false alarm rate, with warnings often having a false-positive rate of reaching up to 90% [4]. Secondly, developers often become overwhelmed with information overload while using these tools, which can

cause them to overlook real bugs and getting lost in irrelevant information.

To fill the gap and help developers to better make use of static analysis tools, previous researchers have introduced the concept of **actionable warning**, namely the warnings need to be acted on by developers. [5], [6]. Particularly, if a warning presents in one revision and disappears in a subsequent revision, then this warning is referred as an actionable warning, otherwise, it is regarded as a false alarm. Different methods have been proposed to identify actionable warnings [7]–[9]. However, two limitations still exist: (i) Most of these methods rely on hand-crafted features, which heavily depend on manual design and expert domain knowledge. (ii) Prior research has focused solely on detecting actionable warnings without addressing the crucial concern that not all actionable warnings are valid indicators of real bugs.

In this paper, we aim to automate the task of identifying actionable warnings produced by static analysis tool (*infer* in this study) and further validatin actionable warnings by recommending AWHB (Actionable Warning with High probability to be real Bug). To achieve this, we build a large dataset of actionable warnings and propose a two-stage framework called ACWRECOMMENDER which includes a coarse-grained detection stage and a fine-grained reranking stage. The dataset is collected from the Top 500 popular C projects on Github, consisting of 538 actionable warnings and 30,590 false alarms. Then each actionable warning is assigned a weak label using semantic and structural matching rules to estimate the likelihood of it being a real bug. In the coarse-grained detection stage, we train a detector to predict whether a warning is actionable or not using the dataset. In the fine-grained reranking stage, we fine-tune the model to prioritize AWHB using weakly supervised learning.

We conducted extensive experiments on two tasks, actionable warnings detection and AWHB recommendation, to evaluate the effectiveness of our ACWRECOMMENDER. Our proposed model demonstrated superiority over several baselines. Results show that our first-stage detector significantly outperforms other baselines by 91.7% in terms of F1-score for the actionable warning detection task, and our second-stage reranker performs better than its three baselines in terms of nDCG and MRR for the AWHB recommendation task. An in-the-wild evaluation was also conducted, where our tool recommended 24 actionable warnings to Github developers,

* corresponding author.

22 of which were confirmed as real bugs, further justifying the practical usage of our approach.

II. APPROACH

We first introduce how to build an actionable warning dataset under weak supervision. We then present the details of our proposed two-stage model. The Overall framework is illustrated in Fig. 1

A. Actionable Warning Collection and Labeling

The aim of this stage is to gather all actionable warnings and assign a label for each actionable warning to represent its likelihood of being a real bug. Regarding the actionable warning collection, we followed the process of previous studies [7], for each revision a given project, we run static analysis tool *infer* to generate a list of warnings. Then we automatically check if the warning disappears in later revisions, if yes then the warning is labeled as an actionable warning, otherwise the warning is treated as a false alarm. Any warnings that have not been resolved for over two years are also deemed false alarms. Finally, we have collected 538 actionable warnings and 30,590 false alarms from top-500 GitHub C projects.

However, we find that **actionable warnings collected by the current pipeline are largely invalid and may not necessarily represent real bugs**, this observation is also consistent with the latest empirical findings [10]. The main reason is that the current method regards all disappeared warnings as actionable warnings, while this assumption is rather shaky because the disappearance of such warning(s) can be caused by a non-relevant fix. In this study, we aim to take one more step further by assigning actionable warnings with different probability scores under weak supervision. The higher probability scores indicate actionable warning(s) are more likely to be real bug(s) (referred to **AWHB** in this study), which should be inspected at the beginning. To gather more accurate actionable warnings, we estimate the “matching degree” of each actionable warning and its bug-fix commit in terms of two perspectives, i.e., semantical matching rule (using commit message) and structural matching rule (using code change context).

1) *Semantic Matching Rule.*: The commit message of a bug-fix reversion summarizes the commit and can be used to estimate a semantic matching score. As shown in Table I, a score of 3 is assigned if the commit message contains keywords that exactly match the warning type (column 3), indicating a high likelihood of a real bug. A score of 2 is assigned if warning context keywords are mentioned in the commit message (column 4). A score of 1 is assigned if the commit message only contains fix-related common keywords (column 5). If the commit message does not match any keyword, a score of 0 is assigned, indicating an unlikely relation to a real bug.

2) *Structural Matching Rule.*: Besides semantic information, we use code change context to assist determination if an actionable warning is fixed by the corresponding bug-fix reversion. As shown in Table II, a structural matching score

is assigned based on code change matching rules. A score of 3 is assigned if the code change matches the fix pattern of the warning type (column 2). A score of 1 is assigned if the code change falls within expected bug-fix scope (column 3). For warnings that do not match any rule, a score of 0 is assigned.

3) *Aggregation of Weal Labels.*: To obtain a more reliable and robust label for each actionable warning, we use majority voting to combine the above semantic matching score and structural matching score, as demonstrated in Equ. 1.

$$Label(x) = \begin{cases} \text{VTB}, & CM(x) + CC(x) > 3 \\ \text{LTB}, & 2 \leq CM(x) + CC(x) \leq 3 \\ \text{UTB}, & 0 \leq CM(x) + CC(x) < 2 \end{cases} \quad (1)$$

Given an actionable warning x , $CM(x)$ refers to the semantic score for x using the commit message matching rule, and $CC(x)$ refers to the structural score for x using the code change matching rule. If the sum of $CM(x)$ and $CC(x)$ is greater than 3, it means the actionable warning is very likely to be a real bug from both semantic and structural aspects and we label x as **VTB (Very Likely To be Bugs)**. Similarly, if the sum of $CM(x)$ and $CC(x)$ falls between 2 and 3, it means the warning x matches the bug-fix reversion from either the semantic or structural aspect and is likely to be a real bug, we label x as **LTB (Likely To be Bugs)**. Lastly, if the sum of $CM(x)$ and $CC(x)$ is less than 2, it means the actionable warning x mismatches the bug-fix reversion and is unlikely to be a real bug, we label such instances as **UTB (Unlikely To be Bugs)**. We then define the warning x whose $Label(x)$ is **VTB** or **LTB** as **AWHB (Actionable Warning with High probability to be real Bug)**.

B. Two-stage Model

In order to identify actionable warnings and recommend AWHB, we propose a two-stage modal which includes a detector and a reranker. In the coarse-grained detection stage, actionable warnings dataset is used to warm-up and let the model learn how to distinguish actionable warnings from false alarms. In the fine-grained reranking stage, we further fine-tune the model to rerank the AWHB to the top by weakly supervised learning.

1) *Warn-up the detector model.*: The goal of this stage is to develop a model that can differentiate between actionable warnings and false alarms. To achieve this, UniXcoder [11] is used to process two types of key information: text-related input (such as bug type) and code-related input (such as AST). The use of UniXcoder allows for unified multi-modal data training, which is ideal for the encoding task. The text-related input is obtained by extracting bug type, qualifier, procedure, and filename from the warning report generated by *infer*. The code-related input is generated by identifying the bug location and associated buggy statement, locating the parent node of the statement in the AST, and extracting control flow information from the AST as code context. Both types of input are fed into UniXcoder, and the resulting embeddings are obtained by concatenating the pre-trained model outputs. The model

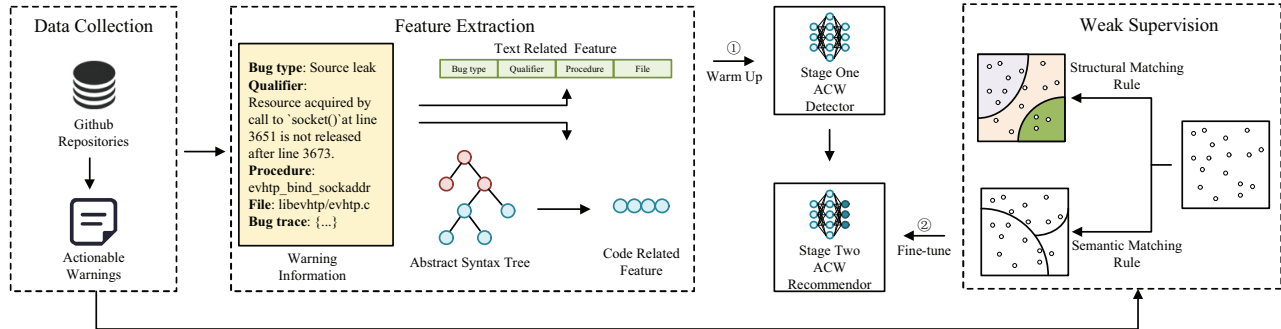


Fig. 1: Overview of Our Tool

TABLE I: Semantic Matching Rule

Warning Type	Warning Qualifier Template	Warning Type Keyword	Warning Context Keyword	Common Keyword
Uninitialized Variable	The value read from <i>variable</i> was never initialized.	initial, define, assign, etc.	<i>variable</i>	fix, repair, bug, warning, solve, problem, etc.
Null Dereference	<i>pointer</i> last assigned on line # could be null and is dereferenced at line #	null dereference, null pointer, etc.	<i>pointer</i>	
Resource Leak	Resource acquired to <i>variable</i> by call to <i>function</i> at line # is not released after line #	resource, leak, etc.	<i>variable, function</i>	
Dead Store	The value written to <i>variable</i> is never used.	dead store, unused, never, etc.	<i>variable</i>	

TABLE II: Structural Matching Rule

Warning Type	Fix Pattern	Scope Pattern
Uninitialized Variable	assign value by assignment or reference	before warning
Null Dereference	add a null-check	before warning
Resource Leak	invoke resource-free-related function	after warning
Dead Store	use the variable, remove assignment	after warning

is trained to identify actionable warnings by warming up UniXcoder. The actionable warning identification task can be viewed as a binary classification problem. That is, for a given reported warning x , we use the model $f(x; \theta)$ to determine whether x is actionable or not. The actionable warning dataset without weak supervision is used to warm up $f(x; \theta)$ and the optimization goal is defined as follows:

$$\min_{\theta} \frac{1}{N} \sum_{x \in \mathcal{X}} \mathcal{L}(f(x; \theta), y_x) \quad (2)$$

where \mathcal{L} denotes a loss function and y_x is the actionable warning label without weak supervision. Any loss function suitable for a classification task can be used in the warm-up process, and in this study we use the Binary Cross Entropy Loss.

2) *fine-tune the reranker model*: Based on the warmed-up actionable warning detector model, we fine-tuned a reranking model by continuing training with the actionable warning dataset under weak supervision. To rank different levels of actionable warnings, we transform the ranking problem into multiclass classification task. That is, for a given reported

warning x , we aim to predict x as V/L/U/False Warning based on our actionable warnings dataset under weak supervision. The optimization problem is defined as follows:

$$\min_{\theta} \frac{1}{N} \sum_{x \in \mathcal{X}} \mathcal{J}(g(f(x; \theta)), \tilde{y}_x) \quad (3)$$

The \tilde{y}_x is the label aggregated from weak supervision. $g(x)$ is the softmax function to compute the probability of each class for x , and \mathcal{J} is Cross Entropy Loss, which is suitable for multiclass classification task. When the optimization is done, the final ranking score for each warning x can be inferred as follows:

$$\mathcal{S}(x) = \begin{cases} class(x) + g_{\bar{y}}(x), & \bar{y} \in \text{VTB, LTB, UTB} \\ class(x) - g_{\bar{y}}(x), & \bar{y} \in \text{False Warning} \end{cases} \quad (4)$$

where $class(x)$ maps each warning x to a base class score 0/1/2/3 if the predicted class is False Warning/UTB/LTB/VTB, \bar{y} denotes the predicted class of x and $g_{\bar{y}}(x)$ denotes the probability of the predicted class.

C. Implementation

We implement our tool by Python 3.9, leveraging the PyQt5 [12], which is a popular Python GUI module. To use our tool, developers first input a `infer` report (txt file from their local computer) by selecting the local file and clicking the `open` button, then the developers click the `Analyze` button in the bottom of the tool interface to call our pre-trained model by feeding the input `infer` reports. In the backend, `ACWRECOMMENDER` can read each warning from

the report, pre-process the necessary text-related and code-related information automatically, then ACWRECOMMENDER identifies the potential actionable warnings and rerank AWHB to the top. Finally, the top returned warnings will be listed in the right side of our tool, we also highlight the warnings with red and orange colours to help developers to address the AWHB with high priority.

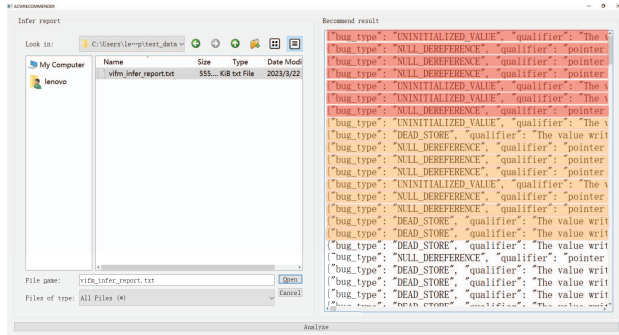


Fig. 2: User Interface of ACWRECOMMENDER

III. EVALUATION

To evaluate the performance of our tool, we first build actionable warning dataset from the top 500 repositories (ordered by the number of stars) in Github for C repositories. In this study, we use *infer* with the default settings to generate warnings, which primarily fall into four categories: uninitialized variable, resource leak, null dereference, and dead store. Finally, we collect both actionable and false warnings from 68,274 revisions of the 394 projects that we reminded, including 538 actionable warnings and 30,590 false warnings. According to our weak supervision rules. Among the 538 actionable warnings, we assign the 57 of them as **VTB**, 59 **LTB** and 422 **UTB**. We only regard the actionable warning whose aggregated labels are **VTB** and **LTB** as **AWHB**.

Then we compare ACWRECOMMENDER with three baselines: Golden Feature [7], Random Forest and Random Selection/Ranking. Golden Feature is a state-of-the-art tool for actionable warning identification, which uses the manually defined golden features (e.g., added lines, number of methods) and leveraged the SVM model. Random Forest uses a random forest model which performed the best in our preliminary experiments, to identify and recommend actionable warning based warning feature embedding. Random Selection/Ranking is a popular baseline, which identifies the actionable warning based on the ratio of actionable warning dataset.

A. Quantitative Analysis

1) *RQ1: The Identification Effectiveness Evaluation:* To evaluate the effectiveness of our proposed tool, i.e., ACWDETECTOR, we evaluate it and the baseline methods on our testing set in terms of Precision, Recall and F1-score. The evaluation results is shown in Table III.

Our proposed tool, ACWDETECTOR, surpasses all the baseline methods by a significant margin in terms of all

TABLE III: Identification Effectiveness Evaluation

Measure	Precision	Recall	F1-score
Random Selecting	0.019	0.019	0.019
Golden Feature	0.294	0.183	0.154
Random Forest	0.433	0.217	0.289
ACWDETECTOR	0.472	0.671	0.554

evaluation metrics. In specific, ACWDETECTOR increase the Precision, Recall and F1-score of Random Forest 9.0%, 184.3% and 91.7%, respectively.

2) *RQ2: The Recommendation Effectiveness Evaluation:* Since the AWHB should be recommended and handled earlier, we evaluate the recommendation effectiveness of our reranker tool. We build 100 queries by randomly selecting 1,000 warnings from testing set. We set K to 1, 3, and 5 for nDCG@K metric. The evaluation result is listed in Table IV.

TABLE IV: Recommendation Effectiveness Evaluation

Measure	nDCG@1	nDCG@3	nDCG@5	MRR
Random Ranking	0.020	0.031	0.052	0.007
Golden Feature	0.074	0.095	0.137	0.064
Random Forest	0.143	0.211	0.239	0.204
ACWRECOMMENDER	0.212	0.381	0.416	0.396

Our proposed tool is effective for AWHB recommendation, and outperforms all the baseline methods. The nDCG@1 value exceeds 0.2, indicating that over one-fifth of the queries can accurately identify AWHB in the topmost position of the recommended warning list with a high degree of certainty. The MRR value is approximately 0.4, suggesting that on average, the first AWHB can be found at the third position in the recommended warning list.

B. In-the-wild Evaluation

Our goal is to help developers find real bugs from static analysis tools while inspecting as few warnings as possible. We conducted an in-the-wild evaluation to assess the practical value of our ACWRECOMMENDER in real-world Github repositories. We randomly selected 10 C repositories from GitHub with more than 500 stars and ran *infer* and obtained 2,197 reported warnings. ACWRECOMMENDER reranked the warnings and generated a list for checking. We manually checked 659 (Top 30%) warnings and found 24 potential bugs, of which 22 were confirmed by developers, with 20 merged and 2 approved by code reviewers. Notably, **12 of the 22 confirmed bugs were identified within the Top 10% of recommended warnings**, demonstrating the effectiveness of our tool in reducing the effort required to find real bugs among a large number of warnings.

IV. CONCLUSION

This research aims to identify actionable warnings produced by static analysis tool and recommend the Actionable Warning with High probability to be real Bug. To address these task, we first collect actionable warnings from top-500 Github

repositories. To the best of our knowledge, this is the first actionable warning dataset collected by mining the histories of popular Github repositories. We propose an approach named ACWRECOMMENDER, which leverages UniXcoder to learn the semantic features of warnings and involves a coarse-grained detection stage and a fine-grained reranking stage. Extensive experiments on the real-world Github repositories have demonstrated its effectiveness and promising performance.

ACKNOWLEDGMENT

This research is supported by the Starry Night Science Fund of Zhejiang University Shanghai Institute for Advanced Study, Grant No. SN-ZJU-SIAS-001. This research is partially supported by the Shanghai Rising-Star Program (23YF1446900) and the National Science Foundation of China (No. 62202341).

REFERENCES

- [1] G. Brat and A. Venet, "Precise and scalable static program analysis of nasa flight software," in *2005 IEEE Aerospace Conference*, 2005, pp. 1–10.
- [2] M. Beller, R. Bholanath, S. McIntosh, and A. Zaidman, "Analyzing the state of static analysis: A large-scale evaluation in open source software," in *2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*, vol. 1, 2016, pp. 470–481.
- [3] M. Christakis and C. Bird, "What developers want and need from program analysis: An empirical study," in *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 332–343. [Online]. Available: <https://doi.org/10.1145/2970276.2970347>
- [4] S. Heckman and L. Williams, "A systematic literature review of actionable alert identification techniques for automated static code analysis," *Information and Software Technology*, vol. 53, no. 4, pp. 363–387, 2011, special section: Software Engineering track of the 24th Annual Symposium on Applied Computing. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584910002235>
- [5] E. A. Alikhashashneh, R. R. Raje, and J. H. Hill, "Using machine learning techniques to classify and predict static code analysis tool warnings," *2018 IEEE/ACS 15th International Conference on Computer Systems and Applications (AICCSA)*, pp. 1–8, 2018.
- [6] K. Heo, H. Oh, and K. Yi, "Machine-learning-guided selectively unsound static analysis," *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*, pp. 519–529, 2017.
- [7] J. Wang, S. Wang, and Q. Wang, "Is there a" golden" feature set for static warning identification? an experimental evaluation," in *Proceedings of the 12th ACM/IEEE international symposium on empirical software engineering and measurement*, 2018, pp. 1–10.
- [8] U. Yuksel and H. Sözer, "Automated classification of static code analysis alerts: A case study," *2013 IEEE International Conference on Software Maintenance*, pp. 532–535, 2013.
- [9] S. Lee, S. Hong, J. Yi, T. Kim, C.-J. Kim, and S. Yoo, "Classifying false positive static checker alarms in continuous integration using convolutional neural networks," *2019 12th IEEE Conference on Software Testing, Validation and Verification (ICST)*, pp. 391–401, 2019.
- [10] H. J. Kang, K. L. Aw, and D. Lo, "Detecting false alarms from automatic static analysis tools: how far are we?" in *Proceedings of the 44th International Conference on Software Engineering*, 2022, pp. 698–709.
- [11] D. Guo, S. Lu, N. Duan, Y. Wang, M. Zhou, and J. Yin, "Unixcoder: Unified cross-modal pre-training for code representation," in *Annual Meeting of the Association for Computational Linguistics*, 2022.
- [12] M. Summerfield, "Rapid gui programming with python and qt: The definitive guide to pyqt programming," 2007. [Online]. Available: <https://api.semanticscholar.org/CorpusID:60150603>