

A DETAILED INFORMATION

A.1 Introduction of BLEU and CodeBLEU Metrics

A.1.1 BLEU. The BLEU (Bilingual Evaluation Understudy) [51] metric is a widely used method for evaluating the quality of text which has been machine-translated from one language to another. Developed to assess the accuracy of machine translation outputs, BLEU compares the machine-produced translations to one or more human reference translations. It quantifies translation quality by calculating the precision of n-grams (sequences of n words) in the machine-generated text relative to the reference texts, while also incorporating a penalty for overly brief translations. This metric provides a score ranging from 0 to 1, where a score closer to 1 indicates a greater similarity between the machine translation and the human reference, suggesting higher translation quality. BLEU is praised for its simplicity and objectivity, making it a standard benchmark in the field of natural language processing for comparing the performance of different translation systems.

The overall BLEU score is calculated as:

$$\text{BLEU} = BP \cdot \exp \left(\sum_{n=1}^N \frac{1}{N} \log p_n \right), \quad (12)$$

where BP is the brevity penalty, N is the maximum n-gram length, and p_n is the precision of n-grams. The brevity penalty is used to penalize overly short translations, and the precision of n-grams is calculated as the ratio of the number of n-grams in the machine translation that appear in the reference translations to the total number of n-grams in the machine translation. The BLEU score is the geometric mean of the precision of n-grams, with the brevity penalty applied to the score. The brevity penalty is calculated as:

$$BP = \begin{cases} 1 & \text{if } c > r, \\ \exp \left(1 - \frac{r}{c} \right) & \text{if } c \leq r, \end{cases} \quad (13)$$

For each n-gram level (e.g., unigram, bigram, trigram, etc.), the precision is calculated as:

$$p_n = \frac{\sum_{\text{clip}_n} c_{\text{clip}_n}}{\sum_{\text{count}_n} c_{\text{count}_n}}, \quad (14)$$

where clip_n is the number of n-grams that appear in the machine translation and the reference translations, count_n is the number of n-grams in the machine translation, and c_{clip_n} and c_{count_n} are the corresponding counts. The BLEU score is the weighted geometric mean of the precision of n-grams, with the weights being the inverse of the number of n-grams. The weights are used to balance the contributions of different n-gram levels to the overall score, with higher weights assigned to longer n-grams. This is done to reflect the fact that longer n-grams are more informative and carry more meaning, and thus should be given more importance in the evaluation.

A.1.2 CodeBLEU. CodeBLEU [57] is an evaluation metric specifically designed for assessing the quality of code generated by machine learning models in programming tasks. It extends the principles of the BLEU metric, traditionally used in natural language processing for evaluating machine translations, to the domain of source code generation. CodeBLEU takes into account not only the syntactic accuracy by comparing n-grams between the generated code and the reference code, but also incorporates semantic and structural aspects unique to programming languages. This includes considering code abstract syntax trees (ASTs), data flow, and logical control structures to better capture the functional correctness of the generated code relative to the reference implementations. By integrating these dimensions, CodeBLEU aims to provide a more comprehensive and meaningful assessment of code generation models, reflecting both the stylistic and functional fidelity of the produced source code. This metric has become increasingly important

as the field of code generation and software engineering assisted by artificial intelligence continues to evolve. CodeBLEU can be defined as the weighted combination of four parts:

$$\text{CodeBLEU} = \alpha \cdot \text{BLEU} + \beta \cdot \text{BLEU}_{\text{weight}} + \gamma \cdot \text{Match}_{\text{ast}} + \delta \cdot \text{Match}_{\text{df}}. \quad (15)$$

where BLEU is calculated by standard BLEU [51], $\text{BLEU}_{\text{weight}}$ is the weighted n-gram match, obtained by comparing the hypothesis code and the reference code tokens with different weights, $\text{Match}_{\text{ast}}$ is the syntactic AST match, exploring the syntactic information of the code, and Match_{df} is the semantic data flow match, capturing the semantic information of the code. The weights α , β , γ and δ are used to balance the contributions of these four parts.

A.2 Details of Datasets and Tasks

In our experiments, we consider the following five widely-used tasks, (1) defect detection, (2) clone detection, (3) authorship attribution, (4) code translation as well as (5) code summarization.

- **Defect detection:** Given a source code, this task is to identify whether it contains defects that may be used to attack software systems, such as resource leaks, use-after-free vulnerabilities and DoS attack. In our experiments, we use the dataset provided by Zhou *et al.*³ [87]. It consists of 27, 318 functions collected from two large C-language open-source projects that are popular among developers and diversified in functionality, i.e., FFmpeg⁴ and Qemu⁵. The defect detection task is treated as binary classification. The positive label indicates that the current project has defects while the negative one represents the opposite case.
- **Clone detection:** Given two code snippets as input, clone detection task aims to check whether they are equivalent in terms of operational semantics. This paper considers the widely used clone detection benchmark BigCloneBench [63] in our experiment. It consists over 6, 000, 000 true clone pairs and 260, 000 false clone pairs from 10 different functionalities. In BigCloneBench, each code fragment is a Java method. Following the settings of Zhou *et al.* [72], we discard those unlabeled data while we use 901, 028 code fragments for training and the other 415, 416 ones for validation and testing purposes.
- **Authorship attribution:** The purpose of this task is to identify the author of a given code snippet. In this paper, we use the Python dataset provided by Alsulami *et al.* [3]. It is collected from the Google Code Jam1 (GCJ)⁶, an annual competition held by Google since 2008. This dataset consists of solutions to 10 problems implemented by 70 authors.
- **Code translation:** This task aims to migrate legacy software from one programming language in a platform to another. The training data for code translation is the code pairs with equivalent functionality in two programming languages. In this paper, we use the dataset provided by Lu *et al.* [41]. It is collected from 4 open-source projects including Lucene⁷, POI⁸, JGit⁹ and Antlr¹⁰. These projects are originally developed in Java and subsequently translated into C#. This dataset consists of 11, 800 pairs of functions or methods, from which 500 pairs have been randomly selected for validation purposes and the other 1, 000 pairs are used for testing.
- **Code summarization:** The objective of code summarization is to generate a natural language comment for a given code snippet. In this paper, we use the CodeSearchNet dataset [31] that

³<https://sites.google.com/view/devign>

⁴<https://github.com/FFmpeg/FFmpeg>

⁵<https://github.com/qemu/qemu>

⁶<https://codingcompetitions.withgoogle.com/codejam>

⁷<https://lucene.apache.org/>

⁸<https://poi.apache.org/>

⁹<https://eclipse.dev/jgit/>

¹⁰<https://github.com/antlr/>

consists of six programming languages, including Python, Java, JavaScript, PHP, Ruby, and Go. The data are collected from publicly available open-source non-fork GitHub repositories, with each documentation representing in the first paragraph of the code. To enhance its overall quality, we employ a filtering process according to the guidelines outlined in [41]. The statistics about the filtered CodeSearchNet dataset is listed in Table 3.

Table 3. Data statistics about the filtered CodeSearchNet dataset.

Language	Training	Validation	Testing
GO	167,288	7,325	8,122
Java	164,923	5,183	10,955
JavaScript	58,025	3,885	3,291
PHP	241,241	12,982	14,014
Python	251,820	13,914	14,918
Ruby	24,927	1,400	1,261

A.3 Details of Evaluation Metrics

To quantitatively evaluate the performance different methods, we consider the following four metrics in our empirical study.

- Attack success rate (ASR): the percentage of the number of successful adversarial attacks examples (denoted as N_{succ}) w.r.t. the total number of examples generated by the corresponding algorithm (denoted as N_{total}):

$$ASR = \frac{N_{succ}}{N_{total}} \times 100\%. \quad (16)$$

More specifically, for classification tasks, adversarial examples that lead to inconsistencies between the victim model's predictions and the original classification results are considered successful. For generation tasks, following the approach in NLP of setting a threshold [17, 60], we consider adversarial examples successful when the BLEU [51] or CodeBLEU [57] scores are below 50% of the original values. The higher the ASR is, the better performance of the algorithm achieves.

- Average adversarial loss (AAL): the average value of adversarial loss.

$$AAL = \frac{1}{N_{succ}} \sum_{i=1}^{N_{succ}} f_1(\mathbf{x}'). \quad (17)$$

- Average semantic similarity (ASS): the average value of semantic similarity.

$$ASS = \frac{1}{N_{succ}} \sum_{i=1}^{N_{succ}} f_2(\mathbf{x}'). \quad (18)$$

- Average modification rate (AMR): the average ratio of the perturbation tokens w.r.t. the total number of tokens:

$$AMR = \frac{1}{N_{succ}} \sum_{i=1}^{N_{succ}} f_3(\mathbf{x}'). \quad (19)$$

- Average query count (AQC): the average number of queries (denoted as N_{query}) w.r.t. the victim model to find a successful adversarial example:

$$AQC = \frac{1}{N_{\text{succ}}} \sum_{i=1}^{N_{\text{succ}}} N_{\text{query}}. \quad (20)$$

It is important to note that the primary goal of MOAA is to generate a diversity trade-off AEs, in contrast to existing methods which only focus on finding an AE. To ensure a fair comparison with AQC, we evaluate the efficiency of algorithms based on the number of query required to generate the first AE.

Recap from Section 2.2 that we have defined our objectives as to minimize AL, SS and MR. Consequently, we would expect algorithms that yield AAL, ASS, AMR and AQC as low as possible.

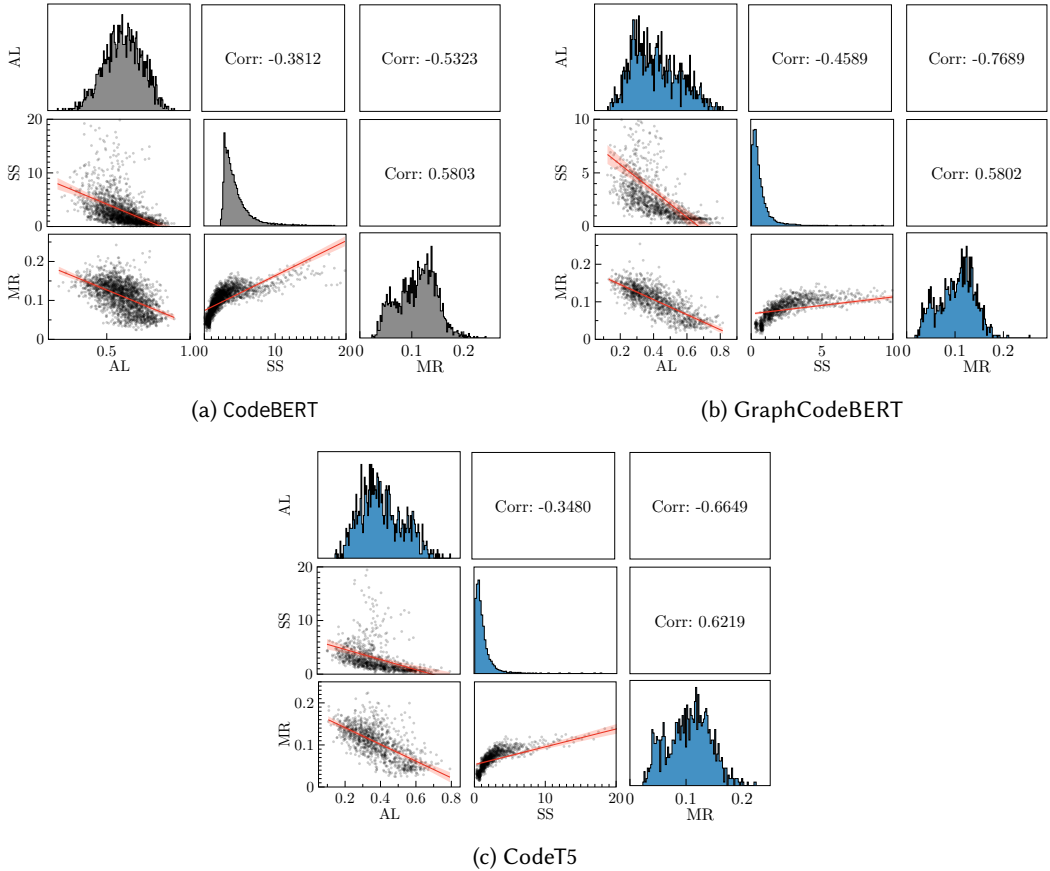


Fig. 15. Scatter plot of matrices (SPLOM) visualizes the correlation of AL, SS and MR values for CodeBERT, GraphCodeBERT and CodeT5 models on defect detection task.

A.4 Introduction of peer algorithms

In this section, we introduce the mechanism of each algorithm we adopted as peer algorithms.

- **MHM:** Zhang *et al.* [83] formalized the process of AE generation as a sampling problem. The problem can be decomposed into an iterative process consisting of three stages: 1) selection of the identifier to be renamed, 2) selection of substitutions and 3) acceptance or rejection decision. To generate AEs for models of code, They proposed Metropolis-Hastings Modifier (MHM), a Metropolis-Hastings sampling-based [25, 46] identifier renaming technique. This method is a black-box attack that randomly selects replacements for local variables and then strategically deciding whether to accept or reject these replacements. This decision is informed by both predicted labels and the corresponding confidence of the victim models, enabling more effective AE generation. MHM employs a pre-defined extensive collection of identifier names, from which the replacements are selected.
- **Greedy Attack:** Yang *et al.* [78] used identifier renaming as the AE generation technique and explored how to produce AEs that are natural. They defined a metric to measure the importance of identifier names in a code snippet and started to substitute identifiers with the highest importance. Greedy Attack greedily selects the replacements (out of all natural substitutes), from which the generated AE makes the victim model produce lower confidence on the ground truth label. If it fails to change the prediction results, Greedy Attack continues to replace the next identifier until all the identifiers are considered or an AE is obtained.
- **ALERT:** Yang *et al.* [78] think that finding appropriate substitutes to generate AEs is essentially a combinatorial optimization problem, whose objective is to find the optimal combination of identifiers and corresponding substitutes that minimizes the victim model's confidence on the ground truth label. Thus, they design an attack based on genetic algorithms, called ALERT to solve the problem that Greedy Attack may be stuck in a sub-optimal solution. If the Greedy Attack fails to find a successful adversarial example, they apply ALERT to search more comprehensively. ALERT first represents chromosomes as a list of identifiers pairs which means replacing the identifier by the replacement, and then initialized the population. Subsequently, it performs genetic operators to generate new solution, and keep solutions with larger fitness values in the population. In the end, the algorithm returns the solution with the highest fitness value.

B EXTENDED RESULTS

B.1 Relationship between objective functions

Fig. 15, Fig. 16, and Fig. 17 shows the scatter plot of matrices (SPLOM) visualizes the correlation of AAL, ASS and AMR values for CodeBERT, GraphCodeBERT and CodeT5 models on defect detection, clone detection and authorship attribution tasks, respectively. The correlation between AL and SS is positive, while the correlation between AL and MR is negative. The correlation between SS and MR is negative. The results indicate that the three objective functions are not independent, and the relationship between them is complex. This validates the rationale of our multi-objective optimization formulation.

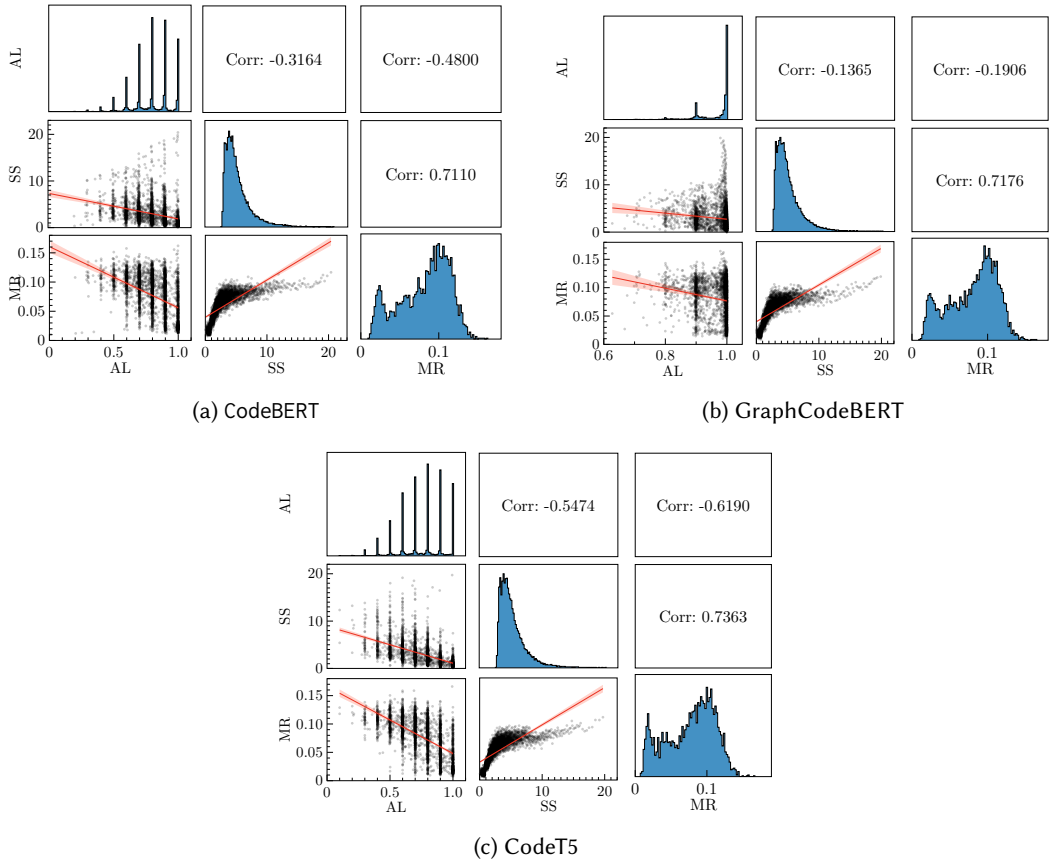


Fig. 16. Scatter plot of matrices (SPLOM) visualizes the correlation of AL, SS and MR values for CodeBERT, GraphCodeBERT and CodeT5 models on clone detection task.

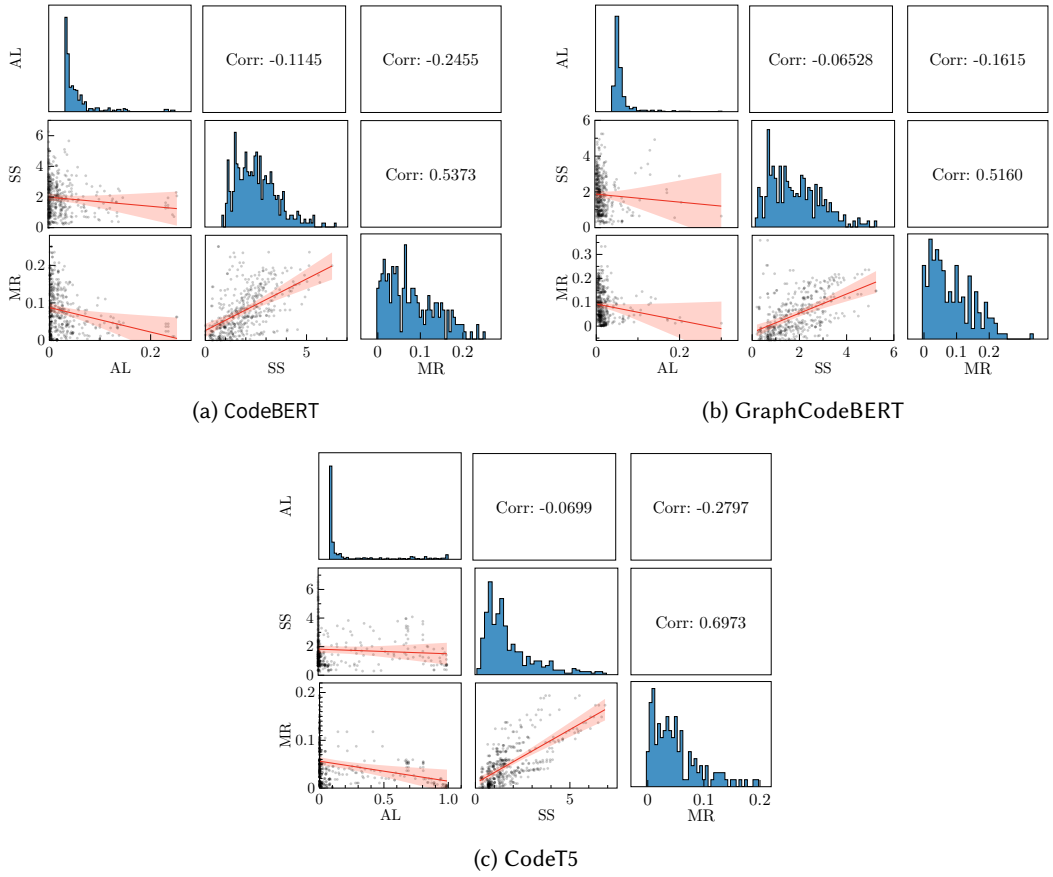


Fig. 17. Scatter plot of matrices (SPLOM) visualizes the correlation of AL, SS and MR values for CodeBERT, GraphCodeBERT and CodeT5 models on authorship attribution task.

B.2 Comparison Results

The Table 4 shows the comparison results of ASR, AAL, ASS AMR and AQC values obtained by MOAA and the other three selected peer algorithms on defect detection, clone detection, authorship attribution, code translation and code summarization tasks.

In addition to attack efficiency, MOAA is also competitive in maintaining the semantic similarity and modification rate of the generated AEs. Note that in Fig. 4, MOAA does not exhibit statistically significant differences in terms of ASS and AMR scores. This is because MOAA generates a set of AEs with diverse trade-offs between the three objectives. When calculating these metrics, we aggregate the values across the entire population, and thus the reported ASS and AMR represent the centric values of the whole population, which can be obscured by extreme AEs that solely focusing on optimizing the adversarial loss (see Fig. 18 for an example).

Fig. 18 visualizes the distribution of generated AEs by MOAA for two examples in defect detection dataset in the objective space. The ‘▲’ represents the mean values of the entire population. From the right half of the figure, we can observe the mean value of f_3 , i.e. MR, is influenced by the distribution of the population, resulting in an increase. Furthermore, Fig. 19 demonstrates the results for AL, SS and MR when selecting the minimum, median and mean values in the AE population generated by MOAA. The bar charts shows that the optimal values within the population generate by MOAA significantly surpass the baseline algorithms, indicating the superiority of our algorithm.

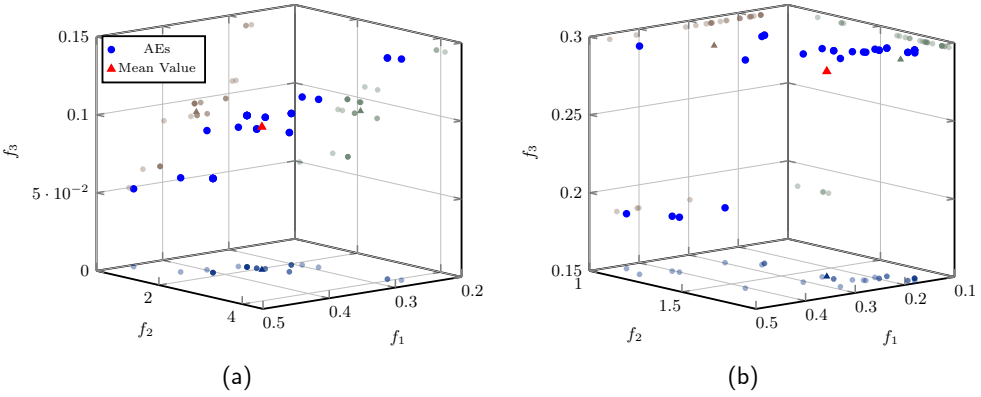


Fig. 18. Visualization of the distribution of generated AEs by MOAA for two examples in defect detection dataset in the objective space.

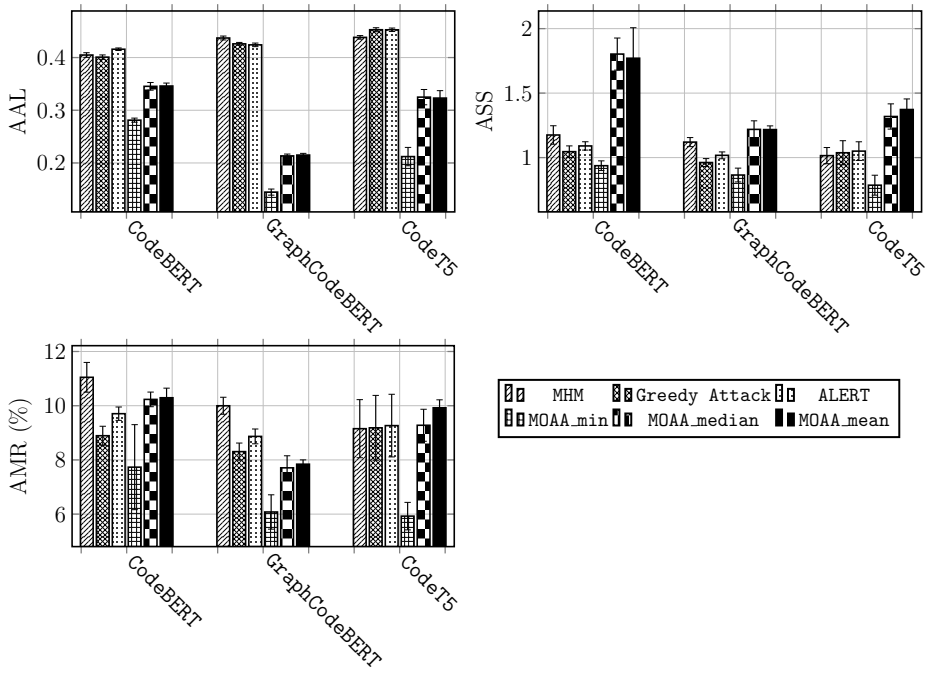


Fig. 19. The AAL, ASS, and AMR metrics when selecting the minimum, median, and average values from the population generated by MOAA for the defect detection task.

B.3 Ablation Study of Importance Score and Identifier Name Prediction

B.3.1 Methods. To investigate the usefulness of importance score and identifier name prediction, we develop three variations:

- **MOAA-v1:** This variant considers a vanilla NSGA-II without importance score and identifier name prediction. Instead, MOAA-v1 use randomly select an identifier name for mutation and replaced by a token provided by ALERT.
- **MOAA-v2:** This variant is similar to MOAA-v1 except using the importance score the same as MOAA to guide the mutation process.
- **MOAA-v3:** This variant is similar to MOAA except using random selection.

Note that the other parameter settings are kept the same as Section 5.2.

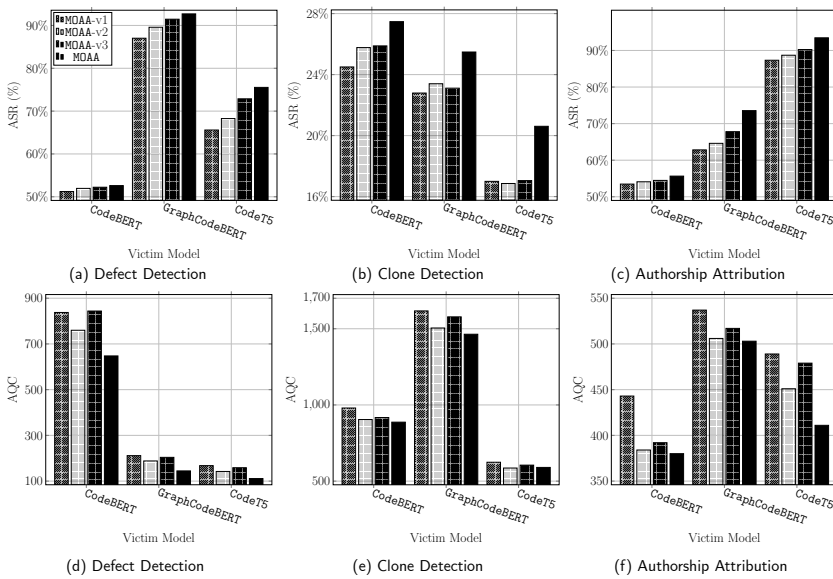


Fig. 20. Bar charts of the attack success rate and average query count of MOAA-v1, MOAA-v2, MOAA-v3 and MOAA on three victim models for the defect detection, clone detection and authorship attribution tasks.

B.3.2 Results and Analysis. Fig. 20 exhibits the result of the ablation study. We observe that MOAA consistently outperforms the other three variations in terms of ASR and AQC. This means that both the importance score and identifier name prediction are crucial for the success of MOAA.

B.4 Examples of Generated AEs

Table 5 presents examples of AEs generated by MHM, Greedy Attack, ALERT and MOAA for CodeBERT on the Defect Detection dataset. In the first example, only MOAA manages to generate a successful AE, despite MHM, Greedy Attack attempt to substitute all identifiers in the original code. This shows the superior efficiency of the search space utilized by MOAA. Regarding the second example, while all attackers achieve success, MOAA distinguishes itself by opting to replace ‘mr’ with ‘memory_map’, a choice informed by the function name ‘set_system_memory_map’. This strategic replacement demonstrates MOAA’s ability to grasp the broader context implied by the code, as opposed to the

other methods which rely solely on the superficial information provided by the identifier name ‘mr’, thus struggling to comprehend the code’s context.

Table 5. Examples of AEs generated by MHM, Greedy Attack, ALERT and MOAA for CodeBERT on Defect Detection dataset.

Attacker	Input	Output	Success
Original	<pre>int ff_schro_queue_push_back(FFSchroQueue *queue, void *p_data) { FFSchroQueueElement *p_new = av_malloc(sizeof(FFSchroQueueElement)); if (!p_new) return -1; p_new->data = p_data; if (!queue->p_head) queue->p_head = p_new; else queue->p_tail->next = p_new; queue->p_tail = p_new; ++queue->size; return 0; }</pre>	1	
MHM	<pre>queue->port; p_data->pfdat; p_new->p_fresh; size->name</pre>	1	×
Greedy Attack	<pre>queue->port; size->address; p_data->pfdat; p_new->pockNEW</pre>	1	×
ALERT	N/A	1	×
MOAA	queue->p_list	0	✓
Original	<pre>void set_system_memory_map(MemoryRegion *mr) { memory_region_transaction_begin(); address_space_memory.root = mr; memory_region_transaction_commit(); }</pre>	1	
MHM	mr->dr	0	✓
Greedy Attack	mr->mm	0	✓
ALERT	mr->fr	0	✓
MOAA	mr->memory_map	0	✓

B.5 Subjective Study

Fig. 21 shows the questionnaire used in the subjective study. We also evaluate the naturalness of the AEs generated by each algorithm by conducting a user study. Following previous works [67, 78], we invite 10 non-author participants who possess a Bachelor/Master degree in Computer Science. For this study, to accommodate the preferences of the participants, we select tasks of different programming languages for each of them. To alleviate recognition burden, we then filter the dataset to exclude code snippets longer than 200 tokens. To make comparison fair, we only select examples that can be successfully attacked by all four algorithms. We randomly pick up 100 examples from the remaining code snippets. For each of these examples, we take the AEs generated by each

algorithm¹¹. Therefore, we ask the participants to grade 400 $\langle \mathbf{x}, \mathbf{x}' \rangle$ pairs blindly in a 5-point rating scale. They are instructed to focus on whether the replaced identifiers are natural in the context of the code snippet.

In this subjective study, we follow the grading criteria as follows:

- 5 for highly natural and contextually appropriate identifier replacement;
- 4 for natural identifier replacements with minor contextual discrepancies;
- 3 for correct but less natural identifier replacements;
- 2 for awkward or contextually inappropriate identifier replacements;
- 1 for incorrect, nonsensical, or misleading identifier replacements.

Table 2 presents the population of AEs generated by MOAA when attacking CodeBERT model for the Defect Detection dataset. The diversity observed in these examples underscores a significant challenge for decision-makers, i.e. determining superiority among the AEs is not straightforward. This highlights the critical role of diversity when generating AEs.

Received 2023-09-28; accepted 2024-04-16

¹¹For MOAA we randomly select one AE from the generated population.

The Naturalness of Identifier

Welcome to our survey! Your feedback is invaluable in helping us understand how natural the replacement of certain identifiers in code snippets appears to experienced developers like you. Below, you will find some code excerpt where a specific identifier has been replaced. Your task is to assess the naturalness of this substitution based on your expertise and intuition about coding practices.

Thank you for participating in our study. Let's get started!

1. Read the following code:

```

1 AVFrame *avcodec_alloc_frame(void)
2 {
3     AVFrame *frame = av_mallocz(sizeof(AVFrame));
4     if (frame == NULL)
5         return NULL;
6     FF_DISABLE_DEPRECATION_WARNINGS
7     avcodec_get_frame_defaults(frame);
8     FF_ENABLE_DEPRECATION_WARNINGS
9     return frame;
10 }

```

Is **'buffer'** replace **'frame'** natural?

Very Unnatural	Unnatural	Neutral	Natural	Very Natural
★	★	★	★	★

2. Read the following code:

```

1 static void cpu_set_irq(void *opaque, int irq, int level)
2 {
3     CPUState *env = opaque;
4     if (level) {
5         CPUIRQ_DPRINTF("Raise CPU IRQ %d\n", irq);
6         env->halted = 0;
7         env->pil_in |= 1 << irq;
8         cpu_check_irqs(env);
9     } else {
10        CPUIRQ_DPRINTF("Lower CPU IRQ %d\n", irq);
11        env->pil_in &= ~(1 << irq);
12        cpu_check_irqs(env);
13    }
14 }

```

Is **'que'** replace **'env'** natural?

Very Unnatural	Unnatural	Neutral	Natural	Very Natural
★	★	★	★	★

Fig. 21. The questionnaire used in the subjective study.