

1177 A DETAILED INFORMATION

1178 A.1 Introduction of BLEU and CodeBLEU Metrics

1179 A.1.1 *BLEU*. The BLEU (Bilingual Evaluation Understudy) [51] metric is a widely used method
 1180 for evaluating the quality of text which has been machine-translated from one language to another.
 1181 Developed to assess the accuracy of machine translation outputs, BLEU compares the machine-
 1182 produced translations to one or more human reference translations. It quantifies translation quality
 1183 by calculating the precision of n-grams (sequences of n words) in the machine-generated text
 1184 relative to the reference texts, while also incorporating a penalty for overly brief translations. This
 1185 metric provides a score ranging from 0 to 1, where a score closer to 1 indicates a greater similarity
 1186 between the machine translation and the human reference, suggesting higher translation quality.
 1187 BLEU is praised for its simplicity and objectivity, making it a standard benchmark in the field of
 1188 natural language processing for comparing the performance of different translation systems.
 1189

1190 The overall BLEU score is calculated as:

$$1191 \quad 1192 \quad \text{BLEU} = BP \cdot \exp \left(\sum_{n=1}^N \frac{1}{N} \log p_n \right), \quad (12)$$

1193 where BP is the brevity penalty, N is the maximum n-gram length, and p_n is the precision of
 1194 n-grams. The brevity penalty is used to penalize overly short translations, and the precision of
 1195 n-grams is calculated as the ratio of the number of n-grams in the machine translation that appear
 1196 in the reference translations to the total number of n-grams in the machine translation. The BLEU
 1197 score is the geometric mean of the precision of n-grams, with the brevity penalty applied to the
 1198 score. The brevity penalty is calculated as:

$$1199 \quad 1200 \quad BP = \begin{cases} 1 & \text{if } c > r, \\ \exp \left(1 - \frac{r}{c} \right) & \text{if } c \leq r, \end{cases} \quad (13)$$

1201 For each n-gram level (e.g., unigram, bigram, trigram, etc.), the precision is calculated as:

$$1202 \quad 1203 \quad p_n = \frac{\sum_{\text{clip}_n} c_{\text{clip}_n}}{\sum_{\text{count}_n} c_{\text{count}_n}}, \quad (14)$$

1204 where clip_n is the number of n-grams that appear in the machine translation and the reference
 1205 translations, count_n is the number of n-grams in the machine translation, and c_{clip_n} and c_{count_n}
 1206 are the corresponding counts. The BLEU score is the weighted geometric mean of the precision
 1207 of n-grams, with the weights being the inverse of the number of n-grams. The weights are used
 1208 to balance the contributions of different n-gram levels to the overall score, with higher weights
 1209 assigned to longer n-grams. This is done to reflect the fact that longer n-grams are more informative
 1210 and carry more meaning, and thus should be given more importance in the evaluation.

1211 A.1.2 *CodeBLEU*. CodeBLEU [57] is an evaluation metric specifically designed for assessing
 1212 the quality of code generated by machine learning models in programming tasks. It extends the
 1213 principles of the BLEU metric, traditionally used in natural language processing for evaluating
 1214 machine translations, to the domain of source code generation. CodeBLEU takes into account
 1215 not only the syntactic accuracy by comparing n-grams between the generated code and the
 1216 reference code, but also incorporates semantic and structural aspects unique to programming
 1217 languages. This includes considering code abstract syntax trees (ASTs), data flow, and logical
 1218 control structures to better capture the functional correctness of the generated code relative to the
 1219 reference implementations. By integrating these dimensions, CodeBLEU aims to provide a more
 1220 comprehensive and meaningful assessment of code generation models, reflecting both the stylistic
 1221 and functional fidelity of the produced source code. This metric has become increasingly important
 1222

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.

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/>

- 1275 • Code summarization: The objective of code summarization is to generate a natural language
 1276 comment for a given code snippet. In this paper, we use the CodeSearchNet dataset [31] that
 1277 consists of six programming languages, including Python, Java, JavaScript, PHP, Ruby, and
 1278 Go. The data are collected from publicly available open-source non-fork GitHub repositories,
 1279 with each documentation representing in the first paragraph of the code. To enhance its
 1280 overall quality, we employ a filtering process according to the guidelines outlined in [41].
 1281 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.

- 1296 • Attack success rate (ASR): the percentage of the number of successful adversarial attacks
 1297 examples (denoted as N_{succ}) w.r.t. the total number of examples generated by the corre-
 1298 sponding algorithm (denoted as N_{total}):

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

1301 More specifically, for classification tasks, adversarial examples that lead to inconsistencies
 1302 between the victim model's predictions and the original classification results are considered
 1303 successful. For generation tasks, following the approach in NLP of setting a threshold [17, 60],
 1304 we consider adversarial examples successful when the BLEU [51] or CodeBLEU [57] scores
 1305 are below 50% of the original values. The higher the ASR is, the better performance of the
 1306 algorithm achieves.

- 1308 • Average adversarial loss (AAL): the average value of adversarial loss.

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

- 1312 • Average semantic similarity (ASS): the average value of semantic similarity.

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

- 1317 • Average modification rate (AMR): the average ratio of the perturbation tokens w.r.t. the
 1318 total number of tokens:

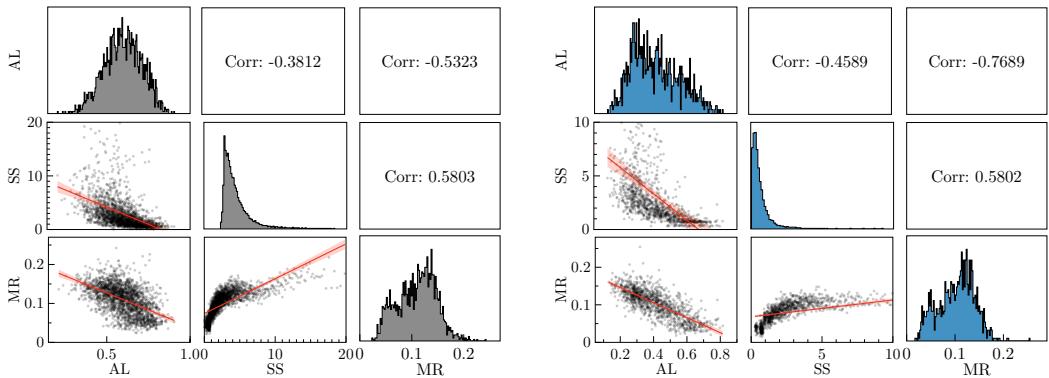
$$AMR = \frac{1}{N_{succ}} \sum_{i=1}^{N_{succ}} f_3(\mathbf{x}'). \quad (19)$$

- 1324 • Average query count (AQC): the average number of queries (denoted as N_{query}) w.r.t. the
 1325 victim model to find a successful adversarial example:

$$1327 \quad AQC = \frac{1}{N_{\text{succ}}} \sum_{i=1}^{N_{\text{succ}}} N_{\text{query}}. \quad (20)$$

1330 It is important to note that the primary goal of MOAA is to generate a diversity trade-off
 1331 AEs, in contrast to existing methods which only focus on finding an AE. To ensure a fair
 1332 comparison with AQC, we evaluate the efficiency of algorithms based on the number of
 1333 query required to generate the first AE.

1334 Recap from Section 2.2 that we have defined our objectives as to minimize AL, SS and MR.
 1335 Consequently, we would expect algorithms that yield AAL, ASS, AMR and AQC as low as possible.
 1336



1350 Fig. 15. Scatter plot of matrices (SPLOM) visualizes the correlation of AL, SS and MR values for CodeBERT,
 1351 GraphCodeBERT and CodeT5 models on defect detection task.

A.4 Introduction of peer algorithms

1369 In this section, we introduce the mechanism of each algorithm we adopted as peer algorithms.
 1370

- 1373 • **MHM:** Zhang *et al.* [83] formalized the process of AE generation as a sampling problem. The
 1374 problem can be decomposed into an iterative process consisting of three stages: 1) selection
 1375 of the identifier to be renamed, 2) selection of substitutions and 3) acceptance or rejection
 1376 decision. To generate AEs for models of code, They proposed Metropolis-Hastings Modifier
 1377 (MHM), a Metropolis-Hastings sampling-based [25, 46] identifier renaming technique. This
 1378 method is a black-box attack that randomly selects replacements for local variables and
 1379 then strategically deciding whether to accept or reject these replacements. This decision is
 1380 informed by both predicted labels and the corresponding confidence of the victim models,
 1381 enabling more effective AE generation. MHM employs a pre-defined extensive collection of
 1382 identifier names, from which the replacements are selected.
- 1383 • **Greedy Attack:** Yang *et al.* [78] used identifier renaming as the AE generation technique
 1384 and explored how to produce AEs that are natural. They defined a metric to measure
 1385 the importance of identifier names in a code snippet and started to substitute identifiers
 1386 with the highest importance. Greedy Attack greedily selects the replacements (out of all
 1387 natural substitutes), from which the generated AE makes the victim model produce lower
 1388 confidence on the ground truth label. If it fails to change the prediction results, Greedy
 1389 Attack continues to replace the next identifier until all the identifiers are considered or an
 1390 AE is obtained.
- 1391 • **ALERT:** Yang *et al.* [78] think that finding appropriate substitutes to generate AEs is es-
 1392 sentially a combinatorial optimization problem, whose objective is to find the optimal
 1393 combination of identifiers and corresponding substitutes that minimizes the victim model's
 1394 confidence on the ground truth label. Thus, they design an attack based on genetic algo-
 1395 rithms, called ALERT to solve the problem that Greedy Attack may be stuck in a sub-optimal
 1396 solution. If the Greedy Attack fails to find a successful adversarial example, they apply
 1397 ALERT to search more comprehensively. ALERT first represents chromosomes as a list of
 1398 identifiers pairs which means replacing the identifier by the replacement, and then initial-
 1399 ized the population. Subsequently, it performs genetic operators to generate new solution,
 1400 and keep solutions with larger fitness values in the population. In the end, the algorithm
 1401 returns the solution with the highest fitness value.

1402 B EXTENDED RESULTS

1404 B.1 Relationship between objective functions

1405 Fig. 15, Fig. 16, and Fig. 17 shows the scatter plot of matrices (SPLOM) visualizes the corre-
 1406 lation of AAL, ASS and AMR values for CodeBERT, GraphCodeBERT and CodeT5 models on defect
 1407 detection, clone detection and authorship attribution tasks, respectively. The correlation
 1408 between AL and SS is positive, while the correlation between AL and MR is negative. The correlation
 1409 between SS and MR is negative. The results indicate that the three objective functions are not
 1410 independent, and the relationship between them is complex. This validates the rationale of our
 1411 multi-objective optimization formulation.

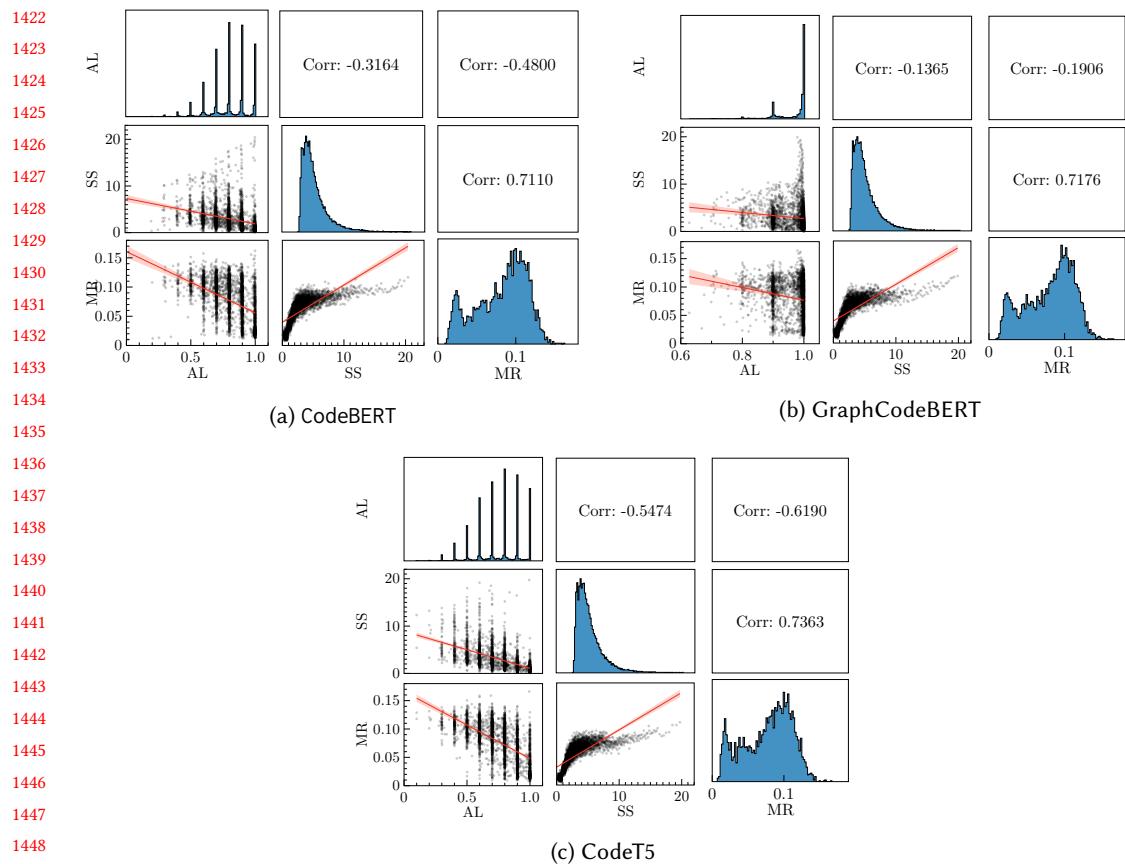
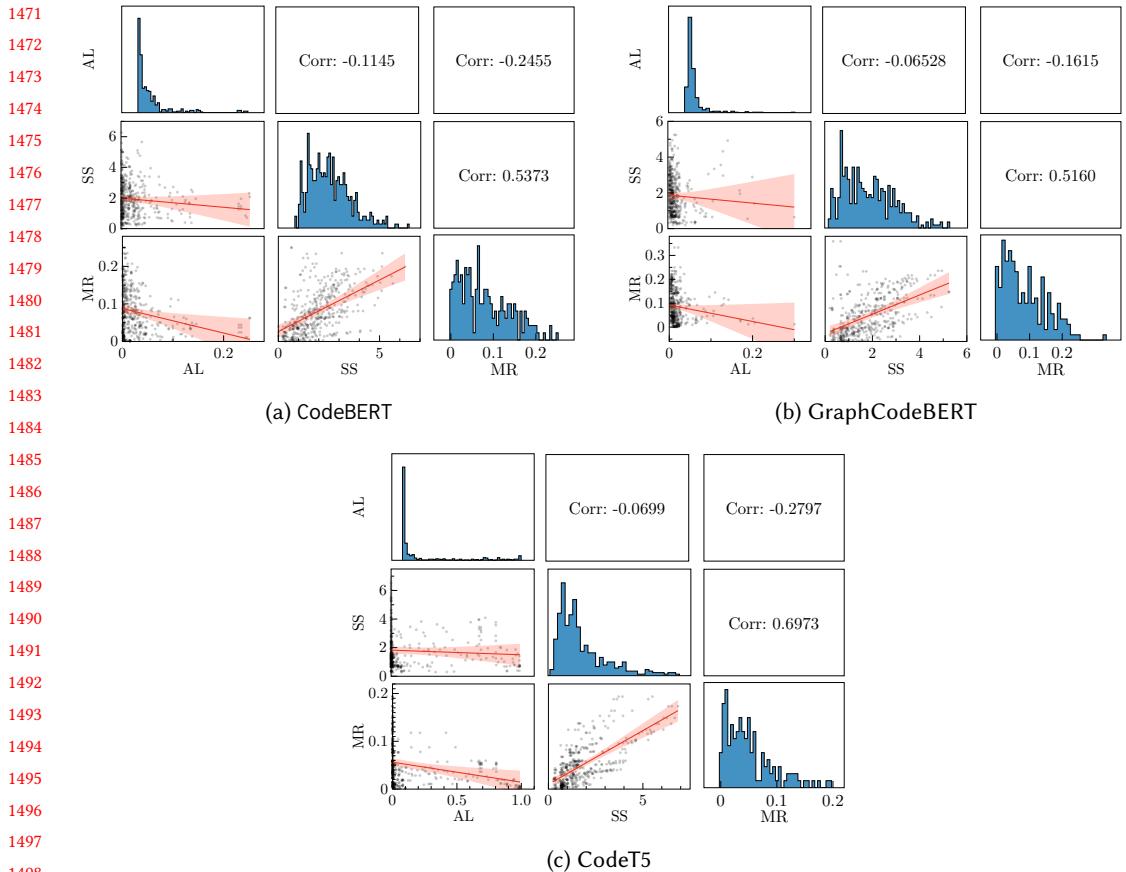


Fig. 16. Scatter plot of matrices (SPLOM) visualizes the correlation of AL, SS and MR values for CodeBERT, GraphCodeBERT and CodeT5 models on clone detection task.



1499 Fig. 17. Scatter plot of matrices (SPLOM) visualizes the correlation of AL, SS and MR values for CodeBERT,
1500 GraphCodeBERT and CodeT5 models on authorship attribution task.
1501

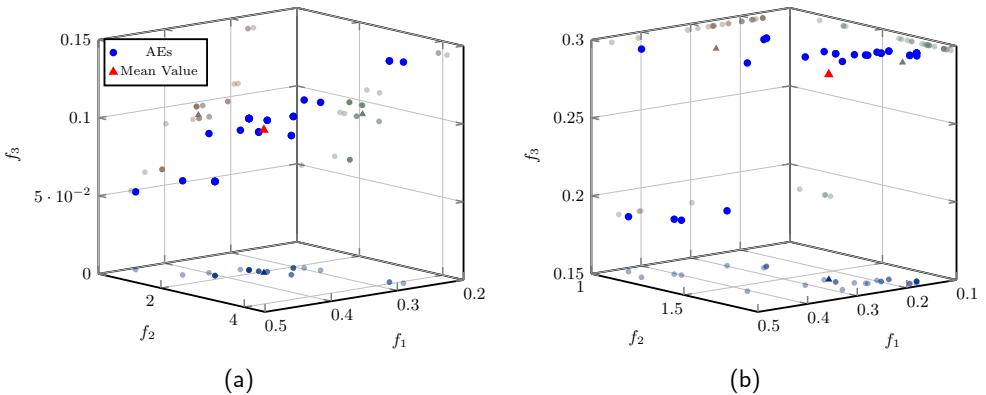
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519

1520 B.2 Comparison Results

1521 The Table 4 shows the comparison results of ASR, AAL, ASS AMR and AQC values obtained
 1522 by MOAA and the other three selected peer algorithms on defect detection, clone detection,
 1523 authorship attribution, code translation and code summarization tasks.

1524 In addition to attack efficiency, MOAA is also competitive in maintaining the semantic similarity
 1525 and modification rate of the generated AEs. Note that in Fig. 4, MOAA does not exhibit statistically
 1526 significant differences in terms of ASS and AMR scores. This is because MOAA generates a set of AEs
 1527 with diverse trade-offs between the three objectives. When calculating these metrics, we aggregate
 1528 the values across the entire population, and thus the reported ASS and AMR represent the centric
 1529 values of the whole population, which can be obscured by extreme AEs that solely focusing on
 1530 optimizing the adversarial loss (see Fig. 18 for an example).

1531 Fig. 18 visualizes the distribution of generated AEs by MOAA for two examples in defect detection
 1532 dataset in the objective space. The ‘▲’ represents the mean values of the entire population. From the
 1533 right half of the figure, we can observe the mean value of f_3 , i.e. MR, is influenced by the distribution
 1534 of the population, resulting in an increase. Furthermore, Fig. 19 demonstrates the results for AL, SS
 1535 and MR when selecting the minimum, median and mean values in the AE population generated
 1536 by MOAA. The bar charts shows that the optimal values within the population generate by MOAA
 1537 significantly surpass the baseline algorithms, indicating the superiority of our algorithm.
 1538



1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553 Fig. 18. Visualization of the distribution of generated AEs by MOAA for two examples in defect detection
 1554 dataset in the objective space.

Table 4. 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.

	Task	Model	Metric	Greedy	ATL(5d)	ALERT	MOAA	Task	Model	Metric	Greedy	ATL(5d)	ALERT	MOAA
Defect Detection	CodeBERT	ASR(%)	34.965(1.181) [†]	50.252(82.258) [†]	51.895(2.039) [†]	52.610(0.097)[†]	52.610(0.097) [†]	CodeBERT	CodeBERT	ASR(%)	78.326(2.967) [†]	80.823(2.013) [†]	86.424(2.114)[†]	94.388(2.097)[†]
		AAL	0.4048(0.0044) [†]	0.4158(0.0024) [†]	0.4158(0.0024) [†]	0.3460(0.0056)[†]	0.3460(0.0056) [†]			AAL	5.9282(0.6357) [†]	6.3346(0.4123) [†]	9.0114(0.0070) [†]	0.0114(0.0070) [†]
		ASS	1.17520(0.0717) [†]	1.0460(0.0447)[†]	1.0904(0.0324) [†]	1.7706(0.2364)	1.7706(0.2364)			ASS	1.9247(0.0164) [†]	2.0397(0.0748) [†]	2.1440(0.0317) [†]	1.8867(0.0264)[†]
		AMR(%)	11.0479(0.5489) [†]	8.8927(0.3491)[†]	9.765(0.2470) [†]	10.2937(0.3540)	10.2937(0.3540)			AMR(%)	13.1255(0.8292) [†]	13.3910(0.7708) [†]	13.3661(0.9694) [†]	13.336(1.3100)
GraphCodeBERT	GraphCodeBERT	AQC	1054.4211(44.2738) [†]	228.5147(5.7074)[†]	643.9245(16.9197) [†]	647.5264(15.1203)	647.5264(15.1203)			AQC	1208.216(11.3535)[†]	1329.3666(16.4602) [†]	1359.2694(17.0333) [†]	1708.442(21.7381)
		ASR(%)	52.2113(6.1609) [†]	66.3950(1.8409) [†]	70.4752(0.2003) [†]	70.4752(0.0838)[†]	70.4752(0.0838) [†]			ASR(%)	60.6947(1.9984) [†]	66.2294(1.5590) [†]	74.2297(1.4916) [†]	91.2442(1.4846)[†]
		AAL	0.4237(0.0027) [†]	0.4266(0.0027) [†]	0.4241(0.0034) [†]	0.2150(0.0032)[†]	0.2150(0.0032) [†]			AAL	6.9529(0.0774) [†]	6.9803(0.0217) [†]	6.646(0.0250) [†]	1.0689(0.0068)[†]
		ASS	1.1201(0.0542) [†]	8.3397(0.3142) [†]	8.8675(0.2732) [†]	9.2658(0.1597) [†]	9.2658(0.1597) [†]			ASS	1.3706(0.0031) [†]	1.3386(0.0073) [†]	16.1356(1.2277) [†]	15.8896(1.7852) [†]
CodeT5	CodeT5	AMR(%)	15.4910(0.9680) [†]	21.1648(0.4252) [†]	60.0875(0.2511) [†]	60.5757(0.1833) [†]	60.5757(0.1833) [†]			AMR(%)	120.406(0.5377) [†]	122.1743(0.3457) [†]	122.1743(0.3457) [†]	159.4993(2.2593)[†]
		ASR(%)	58.6711(1.6507) [†]	59.129(2.8469) [†]	62.7901(1.9411) [†]	75.5396(4.0098)[†]	75.5396(4.0098) [†]			ASR(%)	6.2970(2.5597) [†]	6.2981(2.2557) [†]	6.2981(2.2557) [†]	95.4993(2.2593)[†]
		AAL	0.4383(0.0034) [†]	0.4383(0.0034) [†]	0.4257(0.0033) [†]	0.3233(0.0141)[†]	0.3233(0.0141) [†]			AAL	7.1635(0.3053) [†]	6.9754(0.4176) [†]	6.9153(0.3109) [†]	0.0454(0.0062)[†]
		ASS	1.0152(0.0628)[†]	1.0374(0.0933) [†]	1.0520(0.0724) [†]	1.3728(0.0820) [†]	1.3728(0.0820) [†]			ASS	1.7395(0.0754) [†]	1.7293(0.3826)[†]	1.7989(0.3910) [†]	1.782(0.0227) [†]
GraphCodeBERT	GraphCodeBERT	AMR(%)	9.1522(1.0717) [†]	9.1843(1.1961) [†]	9.2658(1.1559) [†]	9.2280(0.5606) [†]	9.2280(0.5606) [†]			AMR(%)	13.6793(0.8142) [†]	13.4374(1.0781)[†]	13.9038(1.0768) [†]	13.7450(1.2203)
		AQC	120.1973(0.4255) [†]	269.4983(4.1553) [†]	468.4779(4.1553) [†]	111.0769(6.6249)[†]	111.0769(6.6249) [†]			AQC	102.2953(0.4255) [†]	68.4791(4.7953)[†]	69.874(0.4255) [†]	90.047(0.2748)
CodeT5	CodeT5	ASR(%)	9.47430(3.379) [†]	24.8216(0.7862) [†]	28.9126(0.9873)[†]	26.6464(1.2258) [†]	26.6464(1.2258) [†]			ASR(%)	71.9599(2.1537) [†]	78.8542(2.1126) [†]	83.2597(2.1489) [†]	91.1603(2.0348)[†]
		AAL	0.01370(0.0255) [†]	0.0118(0.0277)[†]	0.01400(0.0021) [†]	0.0608(0.0104) [†]	0.0608(0.0104) [†]			AAL	8.5247(1.1589) [†]	8.1258(0.9196) [†]	8.4839(0.7486) [†]	1.6432(0.0346)[†]
		ASS	1.0297(0.0614)[†]	1.406(0.2052) [†]	1.5253(0.0538) [†]	3.2655(0.1381) [†]	3.2655(0.1381) [†]			ASS	1.0395(0.0610) [†]	1.0172(0.0046)[†]	1.0172(0.0046) [†]	1.2173(0.0744)
		AMR(%)	4.43110(29.9439)[†]	7.2608(0.2271) [†]	8.3348(0.2251) [†]	9.9231(0.2942) [†]	9.9231(0.2942) [†]			AMR(%)	10.6472(0.0906) [†]	10.1908(0.0808) [†]	10.1842(0.0763) [†]	9.7968(0.1862)[†]
GraphCodeBERT	GraphCodeBERT	AQC	26.6464(1.2258) [†]	34.3950(1.2258) [†]	34.3950(1.2258) [†]	38.8514(0.1619)[†]	38.8514(0.1619) [†]			ASR(%)	75.406(2.5597) [†]	75.5397(2.5597) [†]	75.5397(2.5597) [†]	75.5397(2.5597)[†]
		ASR(%)	8.07420(3.379) [†]	35.4948(0.4241) [†]	37.7704(0.4241) [†]	38.8514(0.1619)[†]	38.8514(0.1619) [†]			ASR(%)	74.853(1.2567) [†]	75.777(1.2567) [†]	83.325(2.4148) [†]	91.6697(2.0277)[†]
		AAL	0.0622(0.0199)[†]	0.0706(0.0115) [†]	0.0657(0.0088) [†]	0.0837(0.0707) [†]	0.0837(0.0707) [†]			AAL	8.1251(0.6469) [†]	8.1251(0.6469) [†]	8.3119(0.7460) [†]	1.7161(0.4572)[†]
		ASS	2.0150(0.0929)[†]	2.1510(0.1083) [†]	2.1510(0.1083) [†]	3.5807(0.0723) [†]	3.5807(0.0723) [†]			ASS	1.0917(0.0462) [†]	1.0856(0.8349)[†]	1.0982(0.7230) [†]	1.1728(0.7549)
CodeT5	CodeT5	AMR(%)	4.6969(0.4176)[†]	10.2710(0.4624) [†]	11.408(0.4162) [†]	10.2514(0.1939) [†]	10.2514(0.1939) [†]			AMR(%)	10.1801(0.1246) [†]	10.0554(0.4416) [†]	10.0320(0.3695) [†]	9.862(0.3426)[†]
		AQC	22.4663(0.6287) [†]	60.9(0.4614) [†]	60.9(0.4614) [†]	87.293(0.41412)[†]	87.293(0.41412) [†]			AQC	105.2413(0.1599) [†]	78.039(0.41412)[†]	95.8815(0.74213) [†]	100.471(0.1516)
GraphCodeBERT	GraphCodeBERT	ASR(%)	13.30971(5.9351) [†]	15.0468(1.8145) [†]	21.4845(1.8099)[†]	20.6125(0.5927) [†]	20.6125(0.5927) [†]			ASR(%)	70.8053(2.4352) [†]	82.5873(2.0992) [†]	90.5189(1.6097)[†]	94.475(1.6097)[†]
		AAL	0.0446(0.0034) [†]	0.0172(0.0042)[†]	0.0177(0.0018) [†]	0.0375(0.0102) [†]	0.0375(0.0102) [†]			AAL	8.4889(0.2460) [†]	8.3756(0.6131) [†]	8.3665(0.7418) [†]	1.6792(0.3255)[†]
		ASS	1.45110(0.7171) [†]	0.8734(0.0633)[†]	1.3787(0.0660) [†]	3.0520(0.0729) [†]	3.0520(0.0729) [†]			ASS	1.1355(0.0042) [†]	1.0180(0.0264) [†]	1.0180(0.0264) [†]	1.1950(0.0072)[†]
		AMR(%)	6.8090(0.3439) [†]	6.8090(0.3439) [†]	6.8090(0.3439) [†]	6.8090(0.3439)[†]	6.8090(0.3439) [†]			AMR(%)	10.025(0.3275) [†]	10.0554(0.4416) [†]	10.0320(0.3695) [†]	9.862(0.3426)[†]
CodeT5	CodeT5	AQC	22.4663(0.6287) [†]	60.9(0.4614) [†]	60.9(0.4614) [†]	87.293(0.41412)[†]	87.293(0.41412) [†]			AQC	105.2413(0.1599) [†]	78.039(0.41412)[†]	95.8815(0.74213) [†]	100.471(0.1516)
		ASR(%)	29.6515(1.2053) [†]	30.7596(0.2537) [†]	34.0385(1.9939) [†]	37.5289(1.0267)[†]	37.5289(1.0267) [†]			ASR(%)	6.0817(2.1512) [†]	7.2310(2.0378) [†]	79.1855(0.0357) [†]	88.87(1.1068)[†]
		AAL	0.0152(0.0007)[†]	0.0579(0.0031) [†]	0.0605(0.0043) [†]	0.0464(0.0024) [†]	0.0464(0.0024) [†]			AAL	7.3475(0.0018) [†]	7.4288(0.0022) [†]	7.4140(0.0025) [†]	0.015(0.0017)[†]
		ASS	1.93830(21.6716) [†]	2.0739(0.0851) [†]	2.0739(0.0851) [†]	2.3739(0.0445)[†]	2.3739(0.0445) [†]			ASS	1.0235(0.1042) [†]	1.0230(0.0863) [†]	1.0230(0.0863) [†]	1.0424(0.0351)[†]
GraphCodeBERT	GraphCodeBERT	AMR(%)	10.1490(0.1041) [†]	10.1490(0.5464) [†]	10.6832(0.2184) [†]	10.6832(0.2184)[†]	10.6832(0.2184) [†]			AMR(%)	10.0326(0.2601) [†]	10.4525(0.2280) [†]	10.4525(0.2280) [†]	11.4227(0.0423)[†]
		AQC	101.1025(0.1599) [†]	60.9(0.4614) [†]	60.9(0.4614) [†]	75.673(0.41412)[†]	75.673(0.41412) [†]			AQC	105.2413(0.1599) [†]	78.039(0.41412)[†]	95.8815(0.74213) [†]	100.471(0.1516)
		ASR(%)	30.6430(0.3499) [†]	77.14240(0.4289) [†]	77.14240(0.4289) [†]	93.4615(1.0533)[†]	93.4615(1.0533) [†]			ASR(%)	72.4293(1.2211) [†]	79.3892(1.2211) [†]	81.4914(0.0459) [†]	93.7935(1.8473)[†]
		AAL	0.3342(0.0239)[†]	0.3252(0.0222) [†]	0.3250(0.0233) [†]	0.0456(0.0038) [†]	0.0456(0.0038) [†]			AAL	7.1366(0.1149) [†]	7.1250(0.0505) [†]	7.1176(0.0276) [†]	0.0914(0.0041)[†]
CodeT5	CodeT5	ASS	1.5294(0.0207) [†]	1.6612(0.2056) [†]	1.6643(0.0068) [†]	1.0160(0.0467) [†]	1.0160(0.0467) [†]			ASS	1.0767(0.0467) [†]	1.0211(0.0624)[†]	1.0236(0.0622) [†]	1.0709(0.0341)
		AMR(%)	7.74751(0.5695) [†]	8.8926(0.6478) [†]	8.8926(0.6478) [†]	9.5425(0.7816)[†]	9.5425(0.7816) [†]			AMR(%)	9.7834(0.3947) [†]	10.812(0.0758) [†]	10.8085(0.0593) [†]	10.8085(0.0593)
		AQC	120.882(22.7374) [†]	129.8179(0.5695) [†]	136.262(0.2487) [†]	141.5173(15.5111)[†]	141.5173(15.5111) [†]			AQC	113.1941(0.7813) [†]	104.186(2.4511)[†]	104.186(2.4511) [†]	113.8313(7.2029)[†]
		ASR(%)	29.1521(0.3497) [†]	34.8186(0.2498) [†]	35.9256(1.0662) [†]	43.9068(1.1027)[†]	43.9068(1.1027) [†]			ASR(%)	6.0317(1.1697) [†]	6.0317(1.1697) [†]	6.0317(1.1697) [†]	6.0317(1.1697)[†]
GraphCodeBERT	GraphCodeBERT	AAL	42.7681(0.3427) [†]	41.982(3.1617) [†]	41.9604(2.0505) [†]	37.6437(1.2953)[†]	37.6437(1.2953) [†]			ASS	9.4605(0.5683) [†]	9.4158(0.1463) [†]	9.466(0.1768) [†]	6.1562(0.0330)[†]
		ASS	1.4020(0.7874)[†]	1.7456(0.1932) [†]	1.7513(0.0882) [†]	1.8374(0.2153) [†]	1.8374(0.2153) [†]			AMR(%)	1.0230(0.0167) ^{†</}			

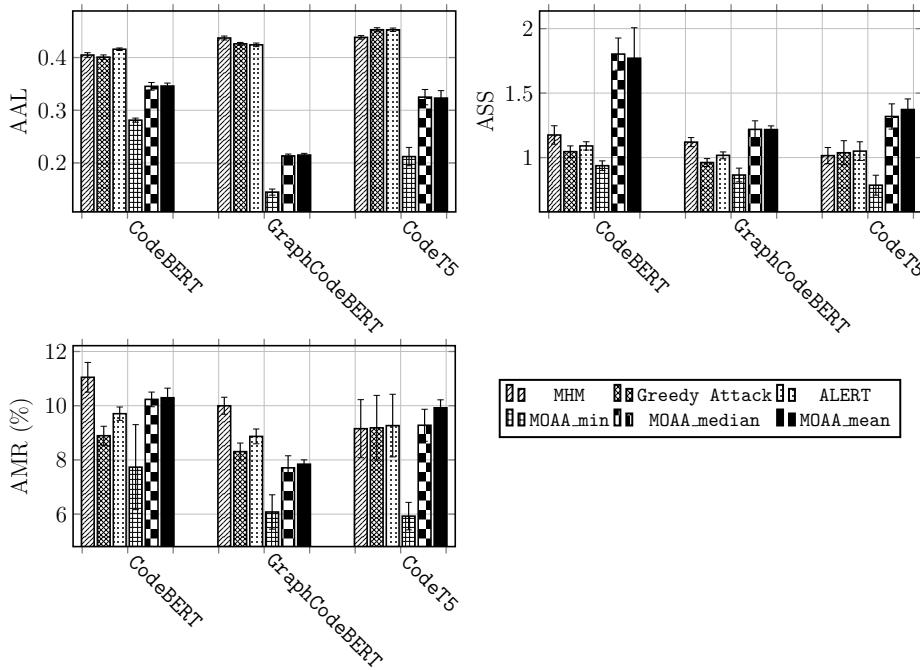


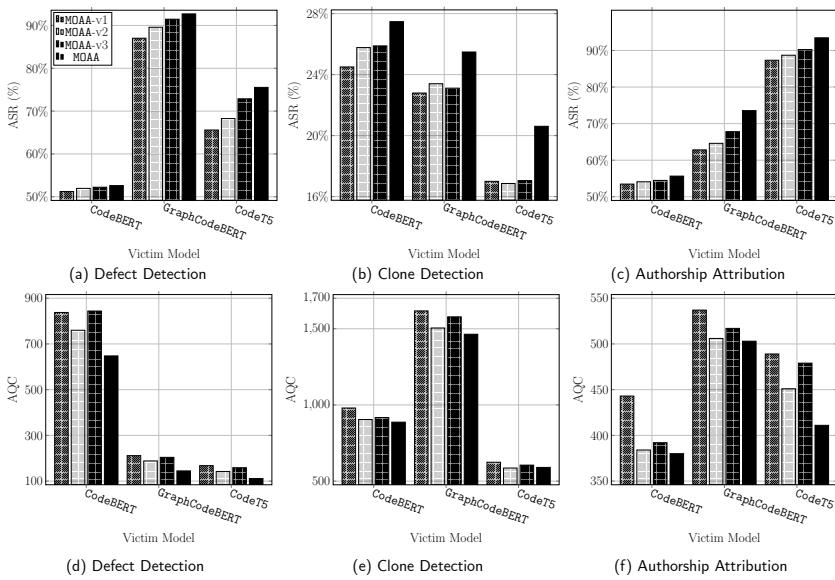
Fig. 19. The AAL, ASS, and AMR metrics when selecting the minimum, median, and average values from the population generated by MOAA for the defect detection task.

1667 B.3 Ablation Study of Importance Score and Identifier Name Prediction

1668 *B.3.1 Methods.* To investigate the usefulness of importance score and identifier name prediction,
 1669 we develop three variations:

- 1670 • **MOAA-v1:** This variant considers a vanilla NSGA-II without importance score and identifier
 1671 name prediction. Instead, MOAA-v1 use randomly select an identifier name for mutation and
 1672 replaced by a token provided by ALERT.
- 1673 • **MOAA-v2:** This variant is similar to MOAA-v1 except using the importance score the same as
 1674 MOAA to guide the mutation process.
- 1675 • **MOAA-v3:** This variant is similar to MOAA except using random selection.

1676 Note that the other parameter settings are kept the same as Section 5.2.



1688 Fig. 20. Bar charts of the attack success rate and average query count of MOAA-v1, MOAA-v2, MOAA-v3 and
 1689 MOAA on three victim models for the defect detection, clone detection and authorship attribution
 1690 tasks.

1702 *B.3.2 Results and Analysis.* Fig. 20 exhibits the result of the ablation study. We observe that MOAA
 1703 consistently outperforms the other three variations in terms of ASR and AQC. This means that
 1704 both the importance score and identifier name prediction are crucial for the success of MOAA.

1706 B.4 Examples of Generated AEs

1708 Table 5 presents examples of AEs generated by MHM, Greedy Attack, ALERT and MOAA for CodeBERT
 1709 on the Defect Detection dataset. In the first example, only MOAA manages to generate a successful
 1710 AE, despite MHM, Greedy Attack attempt to substitute all identifiers in the original code. This shows
 1711 the superior efficiency of the search space utilized by MOAA. Regarding the second example, while
 1712 all attackers achieve success, MOAA distinguishes itself by opting to replace ‘mr’ with ‘memory_map’.
 1713 This strategic replacement demonstrates MOAA’s ability to grasp the broader context implied by the code, as opposed to the
 1714

1716 other methods which rely solely on the superficial information provided by the identifier name
 1717 ‘mr’, thus struggling to comprehend the code’s context.

1718

1719 Table 5. Examples of AEs generated by MHM, Greedy Attack, ALERT and MOAA for CodeBERT on Defect
 1720 Detection dataset.

1721

1722 Attacker	1723 Input	1724 Output	1725 Success
1726 Original	1727 int ff_schro_queue_push_back(FFSchroQueue *queue, void *p_data) 1728 { 1729 FFSchroQueueElement *p_new = av_mallocz(sizeof(FFSchroQueueElement)); 1730 if (!p_new) 1731 return -1; 1732 p_new->data = p_data; 1733 if (!queue->p_head) 1734 queue->p_head = p_new; 1735 else 1736 queue->p_tail->next = p_new; 1737 queue->p_tail = p_new; 1738 ++queue->size; 1739 return 0; 1740 }	1741	1742
1743 MHM	1744 queue->port; 1745 p_data->pfdat; 1746 p_new->p_fresh; 1747 size->name	1748 1	1749 ×
1750 Greedy Attack	1751 queue->port; 1752 size->address; 1753 p_data->pfdat; 1754 p_new->pockNEW	1755 1	1756 ×
1757 ALERT	1758 N/A	1759 1	1760 ×
1761 MOAA	1762 queue->p_list	1763 0	1764 ✓
1765 Original	1766 void set_system_memory_map(MemoryRegion *mr) 1767 { 1768 memory_region_transaction_begin(); 1769 address_space_memory.root = mr; 1770 memory_region_transaction_commit(); 1771 }	1772 1	1773
	1774 MHM	1775 0	1776 ✓
	1777 Greedy Attack	1778 0	1779 ✓
	1780 ALERT	1781 0	1782 ✓
	1783 MOAA	1784 0	1785 ✓

1752

1753

1754

1755

B.5 Subjective Study

1756 Fig. 21 shows the questionnaire used in the subjective study. We also evaluate the naturalness of the
 1757 AEs generated by each algorithm by conducting a user study. Following previous works [67, 78],
 1758 we invite 10 non-author participants who possess a Bachelor/Master degree in Computer Science.
 1759 For this study, to accommodate the preferences of the participants, we select tasks of different
 1760 programming languages for each of them. To alleviate recognition burden, we then filter the dataset
 1761 to exclude code snippets longer than 200 tokens. To make comparison fair, we only select examples
 1762 that can be successfully attacked by all four algorithms. We randomly pick up 100 examples from
 1763 the remaining code snippets. For each of these examples, we take the AEs generated by each
 1764

1765 algorithm¹¹. Therefore, we ask the participants to grade 400 $\langle \mathbf{x}, \mathbf{x}' \rangle$ pairs blindly in a 5-point rating
1766 scale. They are instructed to focus on whether the replaced identifiers are natural in the context of
1767 the code snippet.

1768 In this subjective study, we follow the grading criteria as follows:

- 1769 • 5 for highly natural and contextually appropriate identifier replacement;
1770 • 4 for natural identifier replacements with minor contextual discrepancies;
1771 • 3 for correct but less natural identifier replacements;
1772 • 2 for awkward or contextually inappropriate identifier replacements;
1773 • 1 for incorrect, nonsensical, or misleading identifier replacements.

1774 Table 2 presents the population of AEs generated by MOAA when attacking CodeBERT model for
1775 the Defect Detection dataset. The diversity observed in these examples underscores a significant
1776 challenge for decision-makers, i.e. determining superiority among the AEs is not straightforward.
1777 This highlights the critical role of diversity when generating AEs.

1778

1779

1780

1781

1782

1783

1784

1785

1786

1787

1788

1789

1790

1791

1792

1793

1794

1795

1796

1797

1798

1799

1800

1801

1802

1803

1804

1805

1806

1807

1808

1809

1810

1811

1812 ¹¹For MOAA we randomly select one AE from the generated population.

1813

1814 The Naturalness of Identifier

1815
 1816
 1817 Welcome to our survey! Your feedback is invaluable in helping us understand how natural the replacement of certain identifiers in code snippets appears to
 1818 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
 1819 of this substitution based on your expertise and intuition about coding practices.

1820 Thank you for participating in our study. Let's get started!

1821

1822 1. Read the following code:

```
1823
1824 1 AVFrame *avcodec_alloc_frame(void)
1825 2 {
1826 3     AVFrame *frame = av_mallocz(sizeof(AVFrame));
1827 4     if (frame == NULL)
1828 5         return NULL;
1829 6     FF_DISABLE_DEPRECATED_WARNINGS
1830 7     avcodec_get_frame_defaults(frame);
1831 8     FF_ENABLE_DEPRECATED_WARNINGS
1832 9     return frame;
1833 10 }
```

1834 Is 'buffer' replace 'frame' natural?

Very Unnatural	Unnatural	Neutral	Natural	Very Natural
----------------	-----------	---------	---------	--------------



1835

1836 2. Read the following code:

```
1837
1838
1839 1 static void cpu_set_irq(void *opaque, int irq, int level)
1840 2 {
1841 3     CPUSState *env = opaque;
1842 4     if (level) {
1843 5         CPUIRQ_DPRINTF("Raise CPU IRQ %d\n", irq);
1844 6         env->halted = 0;
1845 7         env->pil_in |= 1 << irq;
1846 8         cpu_check_irqs(env);
1847 9     } else {
1848 10        CPUIRQ_DPRINTF("Lower CPU IRQ %d\n", irq);
1849 11        env->pil_in &= ~(1 << irq);
1850 12        cpu_check_irqs(env);
1851 13    }
1852 }
```

1853 Is 'que' replace 'env' natural?

Very Unnatural	Unnatural	Neutral	Natural	Very Natural
----------------	-----------	---------	---------	--------------



1854 Fig. 21. The questionnaire used in the subjective study.