

1 We thank the reviewers for their valuable feedback. We will address the comments and the concerns as follows.

2 **R1-Motivation.** Our motivation is to solve the binary source code matching problem, which is very important for
3 computer security. It can be leveraged for code vulnerability analysis and malware detection. Given the binary code,
4 the researchers want to explore its specific meaning. A general method is using a decompiler to translate the binary
5 code into pseudo code, which requires a lot of expert experience and has low accuracy. At the same time, the final
6 pseudo code has poor readability because of the missing variable names and the complex "goto" statements. The
7 situation is even serious at function-level because much more information is lost, so researchers expect to find the real
8 corresponding source code. Traditional methods use the code literals to search the source code, but with very low top1
9 score (40%). Our deep learning model could help capture the potential semantic features between the source/binary
10 code and achieve 90% top 1 score. The main contribution of our paper is using deep learning models to make this
11 problem from "cannot" to "can". Our research can benefit tens of thousands of reverse engineering researchers.

12 **R1&R2&R4-Noveltly.** In our paper, the technique is not deep in the modeling part. However, we strongly agree with
13 R1&R4 that this should not be taken as a limitation because getting amazing results with such a simple method proves
14 the capability of deep learning on this task, which is exactly what we want to show. Based on our CodeCMR framework,
15 we tried many useful novel methods, including Circle Loss, Cross-Batch Memory, Adversarial Loss, and Cross-lingual
16 Language Model. We could add a discussion part to talk about these methods, but we still think they are not the focus
17 of this paper because our goal is to use AI to solve valuable problems. Also, we propose the integer-LSTM for code
18 literal modeling, and a simple but useful norm weighted sampling method which helps improve 10% top1 score.

19 **R1&R3-Why Treating Source Code as Text.** Treating source code as text and getting such good results has surprised
20 many experts. In fact, it is a salient point of our paper. Before explanation, we would like to mention that "CFG" is short
21 for "control flow graph" (a graphical representation of binary code) in our paper rather than "context free grammars",
22 we are not sure whether this leads to the misunderstanding. The traditional approach is to parse the codes into graphical
23 inputs. There are two non-negligible indicators in this process: time and success rate, because function-level code loses
24 much information such as "#include" and "#define". For example, it takes about 20 seconds to parse a piece of source
25 code into code property graph with 84.1% success rate on 10,000 pieces. These scores are unacceptable since we want
26 to create a solid method which can be applied to our million-to-billion-level scenario rather than a toy model. Treating
27 source code as text has no time cost and achieves 100% success rate. Surprisingly, when we compare the results of our
28 method with parsing models using a successfully-parsed dataset, we find that our method's top1 score is only 0.5% -
29 1% less. Another advantage is that our method does not need much expert experience, so the AI/NLP researchers can
30 transfer their experience on this task easily, which could help promote the development of this field.

31 **R1&R3-Writing.** (1) We totally agree with R1&R3 that it is more canonical to explain the equations more thoroughly.
32 However, these sentences took up about a quarter of a page, and we had to abridge them because we want to keep other
33 sentences for the readability and consistency of the paper. This problem could be solved easily because we will have an
34 additional page on camera ready version. (2) For Table 1, some additional explaining sentences will be added on camera
35 ready version too. The second section of Table 1 ignores the code literals, and the third section only uses code literals.
36 Since we propose different models on code semantic feature and code literal feature, we should do ablation study to
37 explore the effect of each part. And in the fourth section the RANDOM line means the combination of the best models
38 on code semantic feature and code literal feature with random sampling method, which is DPCNN+HBMP+Integer-
39 LSTM(Integer)+Hier-LSTM(String). By comparing the RANDOM line with DPCNN+HBMP, Integer-LSTM(Integer)
40 and Hier-LSTM(String), we could observe that the code semantic feature is much more useful than the code literal
41 feature, which meets our assumption. To illustrate the effect of value s, we conduct the experiments on s=8 and s=10, it
42 turns out that s=5 is a decline point. We will add a figure to show this.

43 **R1&R3&R4-Dataset.** We decide to publish our dataset and introduce the details of how we build and how to use
44 it. We first extract files from many open-source libraries including openssl, glibc, etc. Then we compile them on 32
45 combinations of different compilers (gcc/clang), different platforms (x86/x64/arm/arm64) and different optimizations
46 (O0/O1/O2/O3). About 80,000 files are successfully compiled, then we split them into about 1,700,000 functions. We
47 randomly sample 50,000 functions and split them into trainset, validset, and testset (3:1:1). The functions are saved in a
48 .pkl file with 33 columns, in which the first column contains source code, and 2nd-33rd columns contain binary code.

49 **R3-Details.** (1) For "source to binary" task, it is feasible to compile and compute similarity scores. Our method is more
50 powerful. Firstly, we could save much time from compiling the code under many combinations, because our method
51 could achieve amazing results on different combinations (more details could be found in supplementary material Figure
52 1). Secondly, sometimes the function-level code could not be successfully parsed because of information loss, while our
53 method could always work. (2) The source code order is not useful on high-level. For example, "if A else B" equals to
54 "if B else A". We have tried positional encoding and LSTM to capture the order feature, but the score declines. (3) We
55 think the answer of the boarder impact question is yes, and we will add it on camera ready version. However, we believe
56 that our research has more positive outcomes because the majority of researchers are protectors rather than attackers.